

1 **Title:** Plant genome evolution in the genus *Eucalyptus* driven by structural
2 rearrangements that promote sequence divergence

3

4 **Running title:** Architectural evolution of plant genomes

5

6 Scott Ferguson^{1*}, Ashley Jones^{1*}, Kevin Murray^{1,2}, Rose Andrew³, Benjamin
7 Schwessinger¹, and Justin Borevitz¹

8

9 1. Research School of Biology, Australian National University, Canberra,
10 Australian Capital Territory, Australia

11 2. Weigel Department, Max Planck Institute for Biology Tübingen,
12 Tübingen, Germany

13 3. Botany & N.C.W. Beadle Herbarium, School of Environmental and Rural
14 Science, University of New England, Armidale, NSW 2351, Australia.

15 *. First authors

16

17 Corresponding authors

18 Scott Ferguson - scott.ferguson.papers@gmail.com

19 Ashley Jones - ashley.jones@anu.edu.au

20 Abstract

21 Genomes have a highly organised architecture (non-random organisation of
22 functional and non-functional genetic elements within chromosomes) that is
23 essential for many biological functions, particularly, gene expression and
24 reproduction. Despite the need to conserve genome architecture, a high level of
25 structural variation has been observed within species. As species separate and
26 diverge, genome architecture also diverges, becoming increasingly poorly
27 conserved as divergence time increases. However, within plant genomes, the
28 processes of genome architecture divergence are not well described. Here we
29 use long-read sequencing and *de novo* assembly of 33 phylogenetically diverse,
30 wild and naturally evolving *Eucalyptus* species, covering 1-50 million years of
31 diverging genome evolution to measure genome architectural conservation and
32 describe architectural divergence. The investigation of these genomes revealed
33 that following lineage divergence genome architecture is highly fragmented by
34 rearrangements. As genomes continue to diverge, the accumulation of mutations
35 and subsequent divergence beyond recognition of rearrangements becomes the
36 primary driver of genome divergence. The loss of syntenic regions also
37 contribute to genome divergence, but at a slower pace than rearrangements. We
38 hypothesise that duplications and translocations are potentially the greatest
39 contributors to *Eucalyptus* genome divergence.

40

41 Keywords

- 42 1. Plant genome evolution
- 43 2. Synteny
- 44 3. Genome rearrangements
- 45 4. Genome assembly

46 5. *Eucalyptus*

47

48 Introduction

49 Genomes from all kingdoms are highly organised, but vary greatly in their
50 structural architecture (Koonin 2009). Within eukaryotic genomes, genome
51 architecture refers to the non-random organisation of functional and non-
52 functional genetic elements within chromosomes (genes, regulatory regions,
53 small RNAs, transposons, pseudogenes, introns, centromeres, telomeres, etc.),
54 and is critical for many biological functions, in particular reproduction and gene
55 expression. However, the conservation and divergence of genome architecture
56 or structure among a group of radiating plant species that share a common
57 karyotype has not been well described.

58

59 For effective recombination during meiosis and the production of viable
60 reproducing offspring, the genome architecture of both parental haplotypes must
61 be highly similar. Changes to the genetic architecture can result in reproductive
62 isolation/incompatibility or non-viable gametes (Hardigan et al. 2020; Simakov et
63 al. 2020). Therefore, a common genome architecture within individuals of a
64 breeding population tends to be highly conserved, except at some loci with high
65 diversity (Jiao and Schneeberger 2020). Similarly, for expression of a gene to be
66 correctly regulated, it must be placed on a chromosome alongside the required
67 promoters, enhancers, and inhibitors. The 3D organisation of the surrounding
68 chromatin must permit physical access to allow transcription (Heng et al. 2004;
69 Dixon et al. 2016; Oudelaar and Higgs 2021).

70

71 Despite this functional need for structural conservation, some structural
72 differences are known to exist between genomes within species. The extent to

73 which reproductively compatible genomes are structurally different is an open
74 area of research; however, several studies have shown genomes with a
75 surprising amount of structural differences to be reproductively compatible (Lin
76 and Gokcumen 2019; Alonge et al. 2020; Jiao and Schneeberger 2020; Tang et
77 al. 2022). Between diverged species, genomes share less of their architecture
78 than genomes within species, but typically genome architecture is conserved in
79 proportion to phylogenetic distance (Luo et al. 2020; Weissensteiner et al. 2020;
80 Derežanin et al. 2022; Ruggieri et al. 2022), and becomes poorly conserved at
81 larger evolutionary distances (Koonin 2009).

82

83 However, genomes have often, but not always, been viewed as containers to
84 hold genes (Heng 2009; Marques et al. 2019). The legacy of a gene-centric
85 genome has persisted due to the Modern Synthesis (Crkvenjakov and Heng
86 2022) and the highly influential work of Dawkins (Dawkins 1976) and others.
87 Guided by an evolutionary view dominated by genes and gene variants, many
88 genomes from various species have been sequenced, and by identifying their
89 genes and gene variants have provided us with a better understanding of the
90 processes of evolution, divergence, and speciation (Rellstab et al. 2015; Schumer
91 et al. 2015; Meier et al. 2017). However, a heavily gene-centric view may also
92 limit our understanding (Heng 2009). This heavily gene variant-based view of
93 evolution was common until recent advances in long-read sequencing
94 technologies enabled genome-wide investigations into genome architecture
95 (Amarasinghe et al. 2020). Larger structural genome changes were thought to be
96 rare, and as such genomes have been treated as largely structurally static, with
97 individuals typically conceptualised as differing mostly by single nucleotide
98 polymorphisms (SNPs) (Feulner and De-Kayne 2017). Pan-genome studies by
99 using a collection of genes or sequences in a population or species (Bayer et al.

100 2020; Lei et al. 2021), have revealed a significant amount of structural variation
101 within genomes (Torkamaneh et al. 2021; Tang et al. 2022; Li et al. 2023).
102
103 Shared genome architecture is measured by synteny. Synteny is the
104 conservation of both the order and sequence of homologous chromosomes
105 between genomes (Passarge et al. 1999; Dawson et al. 2007; Heger and Ponting
106 2007). Synteny can refer both to individual genome regions, or in aggregate
107 when comparing whole genomes. A genome pair with a large proportion of
108 syntenic loci can be said to be more syntenic than a genome pair with a small
109 proportion of syntenic loci. Synteny can become disrupted by the loss, gain,
110 duplication, rearrangement, or divergence of existing sequences.
111 Rearrangements can occur as inversions, translocations, and duplications;
112 altering the order of sequences within chromosomes while maintaining gene
113 content and are often labelled as structural variants (SVs) (Rieseberg 2001).
114 Species-specific sequences resulting from the insertion, deletion, or localised
115 divergence of sequence, appear as unaligned regions when genomes are
116 analysed. The true origin of unaligned regions is more difficult to infer than
117 rearrangements or syntenic regions (Weisman et al. 2020).
118
119 Crucial to the study of plant genome evolution is study group choice. The ideal
120 study group would be naturally evolving, have low prezygotic reproductive
121 barriers, highly specious, and exist over a wide and variable evolutionary range.
122 *Eucalyptus* with over 800 wild and undomesticated species that exist across a
123 wide geographic and environmental range (Potts and Wiltshire 1997; Booth et al.
124 2015; Supple et al. 2018), retain a conserved karyotype (Grattapaglia et al.
125 2015; Butler et al. 2017), are pollinated by generalist pollinators (Pfeilsticker et
126 al. 2023), are capable of wide-ranging dispersal of genetic material (Bezemer et

127 al. 2016; Murray et al. 2019), and span 50 million years of divergent evolution
128 (Thornhill et al. 2019) make an ideal genus to study plant genome evolution.
129
130 Continuing our study into plant genome evolution (Ferguson et al. 2022b), we
131 generated long read sequences and assembled the genomes of 30
132 undomesticated *Eucalyptus* genomes and outgroups from two closely related
133 genera, *Angophora floribunda* and *Corymbia maculata*. Combined with three
134 previously and identically assembled *Eucalyptus* genomes (Ferguson et al.
135 2022b), we create a dataset covering approximately 1-50 million years of
136 diverging genome evolution, including all eight *Eucalyptus* subgenera (Thornhill
137 et al. 2019; Nicolle 2022, 6). Identifying all syntenic and rearranged regions
138 between all species pairs we demonstrate the rapid pace at which ancestral
139 genome architecture is lost. We further analysed our results to determine if
140 ancestral genome architecture was being lost to sequence rearrangement,
141 divergence beyond recognition, or insertions and deletions. Additionally, by
142 framing synteny, rearrangement, and unaligned loss or gain with phylogenetic
143 distance we sought to describe the overall pattern of genome evolution.

144

145 Results

146 Sequencing and Assembly

147 To investigate genome architecture, we performed nanopore long-read native
148 DNA sequencing and *de novo* genome assembly for 32 Eucalypt species (30
149 *Eucalyptus*, one *Angophora*, and one *Corymbia*; Table 1). All read libraries were
150 trimmed and filtered in preparation of assembly. Curated read libraries had an
151 average haploid coverage of 42.8x (range: 24.7x to 78.0x). For details of read
152 libraries and sequence length distributions, see Supplementary Table S1 and S2,

153 and Supplementary Figures S1 and S2. *Eucalyptus pauciflora* FAST5 files were
154 obtained from Wang et al (Wang et al. 2020), processed and assembled as per
155 our datasets, using a randomly selected 60x coverage of reads.

156

157 After assembling our trimmed and filtered read libraries, we curated our
158 genomes, removing contigs identified as contamination and assembly artefacts.
159 Additionally we filtered haplotigs from our primary assemblies to form pseudo-
160 haploid genomes (Supplementary Table S3). Our genomes, which have a known
161 and conserved haploid karyotype of 11 chromosomes (Ribeiro et al. 2016),
162 assembled into an average of 517 contigs (range: 120 to 1,755) (Table 1). At the
163 completion of our assembly pipeline our genomes had an average contig N50 of
164 3.65 Mbp (range: 614.30 kbp to 11.10 Mbp). Scaffolding contigs against *E.*
165 *grandis* (Myburg et al. 2014) greatly increased our genome contiguity, placing on
166 average 99.69% of our genomes into pseudo-chromosomes (range: 98.83% to
167 99.99%). We have found syntenic scaffolding within *Eucalyptus* to be suitable in
168 the absences of chromosome conformation data (Ferguson et al. 2022b) as
169 Eucalypts have a conserved karyotype (Healey et al. 2021; Low et al. 2022).
170 Additionally, within other genera closely related genomes have been found
171 suitable for scaffolding (Burns et al. 2021). Additionally, RaGOO provides
172 confidence scores for assigning contigs to a scaffold, ordering contigs within
173 scaffolds, and orienting contigs within scaffolds. Confidence scores achieved by
174 our genomes indicated scaffolding was satisfactory (Supplementary Figure S3).

175

176 The completeness of our genomes was evaluated with benchmarking universal
177 single-copy orthologs (BUSCO) (Manni et al. 2021) and LTR Assembly Index (LAI)
178 (Ou et al. 2018). A more complete genome will contain a high proportion of
179 single-copy BUSCO genes, and all our genomes were found to be highly BUSCO
180 complete (average: 97.01%; range: 95.44% to 98.11). LAI searches a genome for

181 long terminal repeat (LTR) sequences, and reports on the proportion that are
 182 intact. The LAI scores achieved by our genomes indicate that they are highly
 183 complete (average: 18.17; range: 14.50 to 23.85). Quality scores for all our
 184 genomes indicate that our genomes are of high quality, contiguity and
 185 completeness (Table 1 and Supplementary Table S4). For statistics and sequence
 186 distribution plots describing our genomes during and at the completion of
 187 assembly (Supplementary Tables S5 and S6, and Supplementary Figures S4 to
 188 S6).

189 **Table 1. Summary of de novo genome assembly, quality assessment,**
 190 **and annotation of 35 Eucalypt genomes.** Alphabetically ordered list of
 191 genomes assembled and associated statistics. Genomes for *E. albens*, *E.*
 192 *meliiodora*, and *E. sideroxylon*, have been previously reported, being assembled
 193 using the same pipeline (Ferguson et al. 2022b).

Species	Scaffolded genome size (Mbp)	% of genome in scaffolds	Scaffold N50 (Mbp)	Contig N50 (Mbp)	Contig count	BUSCO complete	LAI	TE %
<i>A. floribunda</i>	388.21	99.73%	36.02	4.02	224	96.82%	14.5	34.55%
<i>C. maculata</i>	403.82	99.90%	40.55	4.69	173	97.25%	15.92	36.26%
* <i>E. albens</i>	606.89	99.79%	56.93	2.55	674	96.47%	17.3	46.57%
<i>E. ANBG9806169</i>	507.93	99.65%	49.61	2.40	476	96.86%	22.16	44.00%
<i>E. brandiana</i>	507.08	99.82%	45.47	7.28	168	98.11%	23.85	44.21%
<i>E. caleyi</i>	589.32	99.53%	59.52	4.77	276	96.47%	18.24	46.00%
<i>E. camaldulensis</i>	558.45	99.87%	52.65	2.48	418	96.73%	16.99	45.31%
<i>E. cladocalyx</i>	544.08	99.68%	51.92	2.80	390	97.59%	18.53	45.85%
<i>E. cloeziana</i>	480.07	99.75%	44.75	1.74	625	97.12%	19.06	42.57%
<i>E. coolabah</i>	606.31	99.53%	53.56	1.29	935	95.44%	15.9	45.89%
<i>E. curtisii</i>	435.26	99.96%	40.29	2.96	288	97.29%	18.34	41.66%
<i>E. dawsonii</i>	706.90	99.35%	67.73	0.99	1,342	97.51%	17.01	45.88%
<i>E. decipiens</i>	590.95	99.50%	60.20	1.99	552	96.99%	18.87	46.95%
<i>E. erythrocorys</i>	539.20	99.99%	50.47	4.02	250	97.55%	20.18	47.07%
<i>E. fibrosa</i>	589.91	99.85%	55.66	6.45	192	96.73%	17.49	45.10%
<i>E. globulus</i>	545.02	99.28%	51.39	0.64	1,747	96.69%	17.46	44.29%
<i>E. grandis</i>	615.89	99.44%	58.49	0.61	1,747	96.09%	17.11	46.53%
<i>E. guilfoylei</i>	472.36	99.97%	44.61	4.25	209	98.02%	16.39	41.22%
<i>E. lansdowneana</i>	633.52	99.92%	59.67	2.35	489	97.12%	19.46	46.10%
<i>E. leucophloia</i>	568.48	99.38%	54.41	2.66	382	96.99%	17.91	44.37%
<i>E. marginata</i>	512.89	98.83%	50.56	1.01	989	96.17%	19.58	43.43%
* <i>E. meliiodora</i>	639.15	99.30%	60.83	1.87	564	98.67%	18.32	47.20%
<i>E. meliiodora x E. sideroxylon</i>	603.57	99.80%	57.05	6.22	281	97.72%	17.96	46.71%
<i>E. microcorys</i>	440.91	99.92%	41.20	4.00	233	97.21%	16.2	41.39%
<i>E. paniculata</i>	588.85	99.66%	55.38	3.70	330	97.12%	18.58	44.92%
<i>E. pauciflora</i>	494.03	99.88%	50.46	6.58	209	97.25%	20.29	43.10%

<i>E. polyanthemos</i>	603.28	99.56%	57.46	4.66	300	96.82%	17.52	45.55%
<i>E. pumila</i>	529.75	99.70%	48.19	2.49	473	97.38%	17.74	44.17%
<i>E. regnans</i>	494.97	99.84%	47.06	5.26	205	97.25%	20.18	43.06%
<i>E. shirleyi</i>	597.18	99.88%	56.34	6.91	181	97.29%	19.89	45.85%
*<i>E. sideroxylon</i>	592.133	99.87%	62.13	5.22	297	96.65%	18.68	46.57%
<i>E. tenuipes</i>	397.78	99.99%	35.74	3.43	207	96.39%	15.07	37.82%
<i>E. victrix</i>	557.16	99.85%	53.19	11.10	120	96.65%	18.71	44.34%
<i>E. viminalis</i>	558.71	99.11%	52.93	0.65	1,755	96.47%	16.5	44.57%
<i>E. virginea</i>	532.79	99.97%	56.15	2.39	376	97.08%	17.69	43.78%

194

195

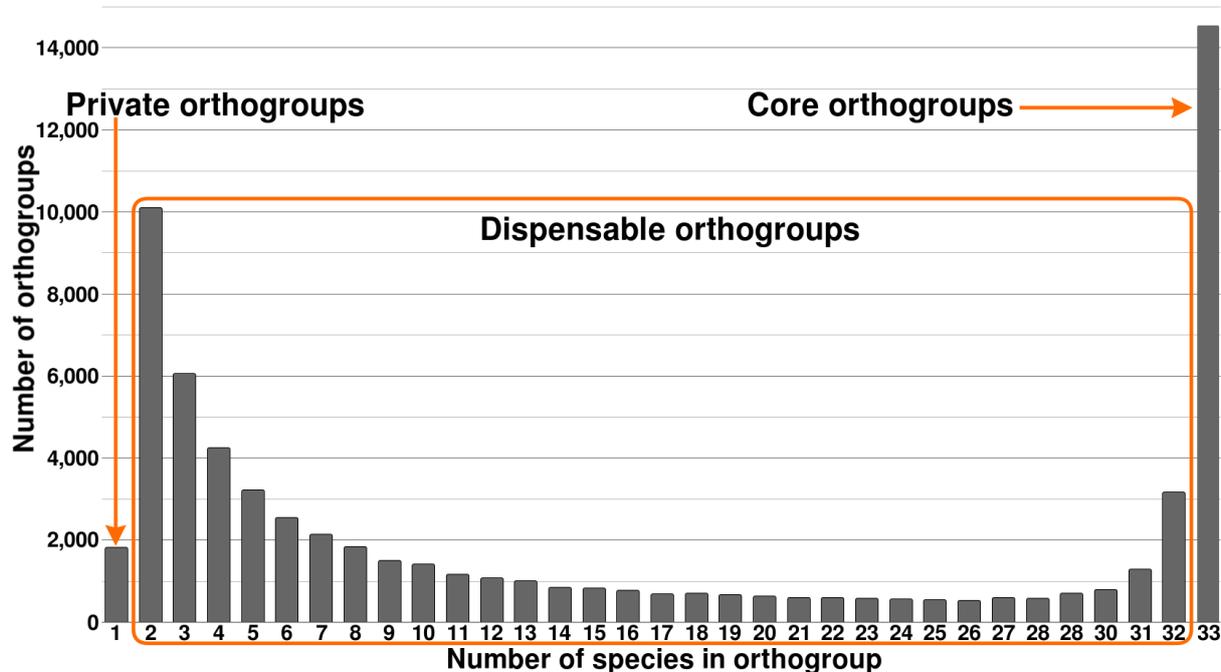
Genome annotation

196 As masking of repeats within genomes aids in gene annotation, we annotated
197 our genomes for both transposable elements (TEs) and simple repeats. Repeat
198 annotation was performed using *de novo* repeat libraries built for each genome.
199 Repeat annotation resulted in the classification of an average of 43.78% (range:
200 34.55% to 47.07%) of our genomes as TEs, and an average of 1.25% (range:
201 1.14% to 1.39%) as simple repeats (Table 1 and Supplementary Table S7). After
202 soft-masking all genomes, we trained species-specific gene HMM models and
203 subsequently annotated all genomes for genes. Gene models were trained on all
204 available gene transcripts for *A. thaliana* (Taxonomy ID: 3702) and Myrtaceae
205 (Taxonomy ID: 3931) found within the NCBI (Sayers et al. 2021). Annotation
206 predicted an average of 53,390 (range: 41,623 to 77,764) gene candidates
207 within our genomes (Supplementary Table S8). While the number of annotated
208 genes is consistent with plant gene number estimates (Sterck et al. 2007) there
209 is a wide variation between genomes. It is important to note the genes
210 annotated within these genomes will contain both false positives and false
211 negatives and are gene candidates, which in addition to real gene number
212 variations will contribute to the variation in the number of annotated genes.

213

214 *Eucalyptus* pan-genome

215 Due to the shared evolutionary history of our genomes, many gene candidates
216 will be homologues that have arisen pre-speciation (orthologs) or post-speciation
217 (paralogs) (Jensen 2001). To examine the evolutionary relationship between
218 *Eucalyptus* gene candidates we placed all highly similar primary (longest) gene
219 transcripts into orthogroups (OG). Of the 1,761,851 identified gene candidates
220 across our 33 *Eucalyptus* genomes, 1,726,511 (97.99%) were placed into one of
221 68,248 orthogroups (OG). The remaining 35,340 (2.01%) unique genes were not
222 placed within an OG as their sequences were too dissimilar (>40% transcript
223 identity and e-value < 0.001) to all other genes. On average each genome had
224 98.03% (range: 94.62% to 98.03%) of its gene candidates placed within an OG.
225 0.26% (4,551) of all gene candidates were found to occur within a genome-
226 specific OG. For detailed statistics on orthogrouping see Supplementary Tables
227 S8 and S9. Additionally, OGs were classified as core (present in all species),
228 dispensable (present in at least two species), and private (present in a single
229 species) (Figure 1). A total of 21.33% (14,552) of OG were core, likely
230 representing key *Eucalyptus* genes. Most OGs were dispensable, 76.00%
231 (51,858), which may be a source of phenotypic and adaptive variation within the
232 species. Only a very small number were private, 2.67% (1,821), potentially
233 representing highly species-specific genes and newly evolved genes.
234



235
236
237
238
239
240

Figure 1. Pangenome of 33 species of *Eucalyptus*. Shows the number of orthogroups shared by an increasing number of genomes. Private orthogroups are orthogroups that exist within a single genome, core exist in all, and dispensable orthogroups are those that exist in 2 to 32 ($n-1$) genomes.

241 *Eucalyptus* phylogeny

242 To describe the evolutionary patterns between our genomes we built a
243 phylogenetic tree from single-copy BUSCO genes. We additionally included the
244 genome of *C. calophylla* genome, which was identically assembled (Ahrens et al.
245 2021). From the initial BUSCO set of 2,326 genes, we selected only genes
246 present within 30 or more genomes, leaving 2,106 BUSCO genes across our 36
247 genomes or 72,516 total genes. For each gene, we generated a multi-sequence
248 alignment (MSA) with MAFFT, which we then trimmed and filtered, removing low
249 abundance regions and genes with overall poor alignments, leaving 1,674 gene
250 MSAs. Each MSA was used to construct a gene tree, subsequently all gene trees
251 were combined into a consensus species tree. The species tree was manually
252 rooted using the established relationship between *Angophora*, *Corymbia* and
253 *Eucalyptus* (Thornhill et al. 2019) (Figure 2 and Supplementary Figure S7). The
254 species tree in Newick format is available within Supplementary results.

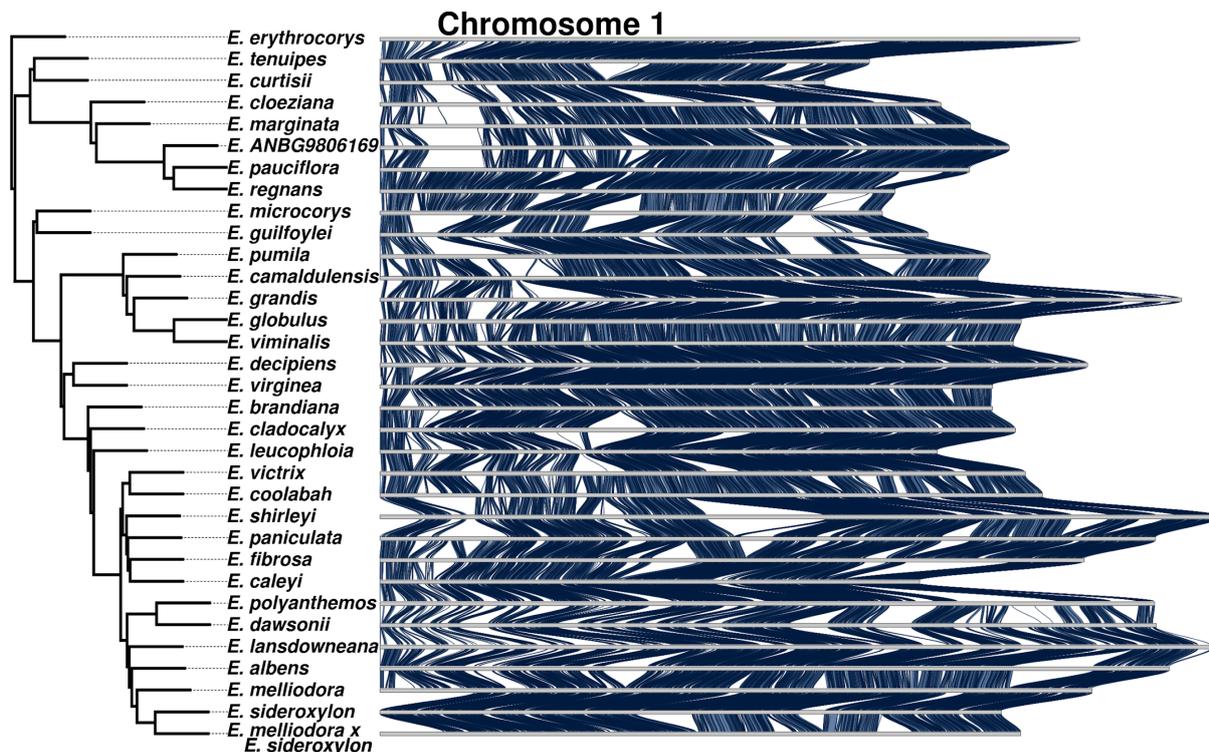
255

256 After constructing the species tree, *E. salubris* was unexpectedly found to be
257 grouped with *E. pauciflora* and *E. regnans*. If correctly placed, *E. salubris* would
258 be a sister lineage to the *Adnataria* group (*E. victrix* to *E. sideroxyton*) (Thornhill
259 et al. 2019). Morphological examination of the sample tree revealed that the tree
260 was incorrectly labelled. The correct species name is currently unknown, as such
261 we use its NCBI name *E. ANBG9806169*.

262

263 Genome conservation and loss

264 To resolve the syntenic and non-syntenic regions of our *Eucalyptus* genomes, we
265 performed one-to-one genome comparisons for all genome pairs. Whole genome
266 alignments for all comparisons were analysed with SyRI (Goel et al. 2019), and
267 subsequently, all genomic regions within both genomes of an alignment pair
268 were annotated as syntenic, rearranged (inversion, translocation, or duplication),
269 or unaligned (sequence that only exists in one genome, resulting from either an
270 insertion, deletion or sequence divergence). As repeat masking would inflate the
271 unaligned proportion of genome alignments and bias results, all genomes
272 remained unmasked. This analysis resulted in all genomes being annotated for
273 syntenic, rearranged, and unaligned regions 32 times, giving a total of 1,056
274 annotated genomes. A visual summary of shared synteny was plotted using our
275 phylogenetic ordering, see Figure 2 and Supplementary Figures S8 to S17.
276 Inspection of synteny plots indicated that syntenic regions exist across the
277 length of all chromosomes, however, synteny has become highly fragmented,
278 and is differently maintained. We acknowledge that large rearrangements
279 (exceeding contig size) may be under- or over-represented in our analysis due to
280 the current limitations of genome sequencing and scaffolding methods, see
281 discussion.



283
284
285
286
287
288
289

Figure 2. Synteny karyotype of Chromosome 1. Blue ribbons between karyotypes indicate the presence of syntenic sequences between species pairs. In all other regions synteny has become lost. Synteny is lost to either rearrangements (inverted, translocated, or duplicated), sequence divergence, loss or gain. Chromosomes are ordered by our phylogenetic tree.

290 Next, we calculated the proportion of sequences shared between genomes, how
291 sequences were shared (syntenic, inverted, translocated, and duplicated), and
292 the frequency at which rearrangements occurred between genomes. For this
293 analysis, we excluded all events less than 200 bp in length (the majority of which
294 were small unaligned annotations). The majority of sequence was shared
295 (syntenic and rearranged) between genomes, averaging 69.35% (range 46.67%
296 to 91.86%). Only four pairwise alignments had <50% shared sequence: *E.*
297 *coolabah*, *E. dawsonii*, *E. grandis*, and *E. melliadora*, all when compared to *E.*
298 *erythrocorys*. Synteny was the major contributor to shared sequence, averaging
299 39.32% (range: 21.34% to 60.44%). Rearrangements averaged: 30.24% (range:
300 16.97% to 49.49%). The remainder of sequence was annotated as unaligned,
301 averaging 30.43% (range: 8.08% to 53.32%) (Table 2). For a per-species

302 comparison breakdown of the percent of genome shared, syntenic, rearranged,
 303 and unaligned see Supplementary Tables S10 to S13.
 304
 305 Examination of the size and frequency of syntenic regions indicates that synteny
 306 between the 11 chromosomes of all genome pairs has, on average, fragmented
 307 into 12,153 (range: 6,657 to 18,810) regions with an average size of 17.97 kbp
 308 (range: 16.27 kbp to 24.26 kbp) (see Supplementary Figures S18 and S19 for per
 309 genome average event size and frequency plots). Rearrangements in total
 310 (inversions + duplications + translocations) contributed more to synteny loss
 311 than did unaligned regions, however, unaligned contributed more than any
 312 single rearrangement type. A more detailed examination of the relative size and
 313 frequency of syntenic, rearrangement, and unaligned events showed that
 314 syntenic regions were long and common, unaligned regions were short and
 315 common, inversions were long and very rare, duplications were shortest and
 316 very common, and translocations occurred at a moderate frequency and size.
 317 Syntenic regions are distributed over the entire length of all chromosomes
 318 between all genome pairs, however, synteny has become highly fragmented by
 319 rearrangements and unaligned regions.

320

321 **Table 2. Summary of synteny, rearranged, and unaligned statistics of**
 322 **all pairwise genome analyses.**

Alignment type	Average event size within each pairwise alignment (kbp)			Event counts			Percent of genome		
	Least	Most	Average	Least	Most	Average	Least	Most	Average
Syntenic	14.03	29.45	18.34	6,657	18,810	12,153	21.34%	60.44%	39.32%
Unaligned	3.28	15.62	7.58	11,808	35,983	22,365	8.08%	53.32%	30.43%
Inverted	29.48	477.33	125.18	81	209	148	0.87%	10.46%	3.28%
Translocated	6.70	22.98	11.65	4,322	15,975	9,521	9.26%	41.14%	19.92%
Duplicated	2.48	6.14	3.60	3,807	32,286	15,350	3.60%	38.97%	14.03%
Rearranged	-	-	-	8,274	46,066	25,020	16.97%	49.49%	30.24%
Total shared	-	-	-	-	-	-	46.67%	91.86%	69.35%

323

324 **Divergence time and genome conservation/loss**

325 To examine these trends of architecture change over increasing divergence time,
326 we examined the relationship between phylogenetic distance and genome
327 conservation and divergence (Figure 3). Not unexpectedly, we find that as
328 phylogenetic distances increase, the proportion of syntenic ($R^2 = 0.261$) and
329 rearranged ($R^2 = 0.356$) sequence decreased as lineages acquire unique
330 genomic variation. Similarly, as phylogenetic distances increase, the proportion
331 of genomes within duplications ($R^2 = 0.189$) and translocations ($R^2 = 0.240$)
332 decreased. Whereas, the portion of genomes unaligned quickly increases with
333 increasing phylogenetic distance ($R^2 = 0.536$). Inversions consistently occupied a
334 small proportion of genomes across all phylogenetic distances ($R^2 = 0.000$).

335

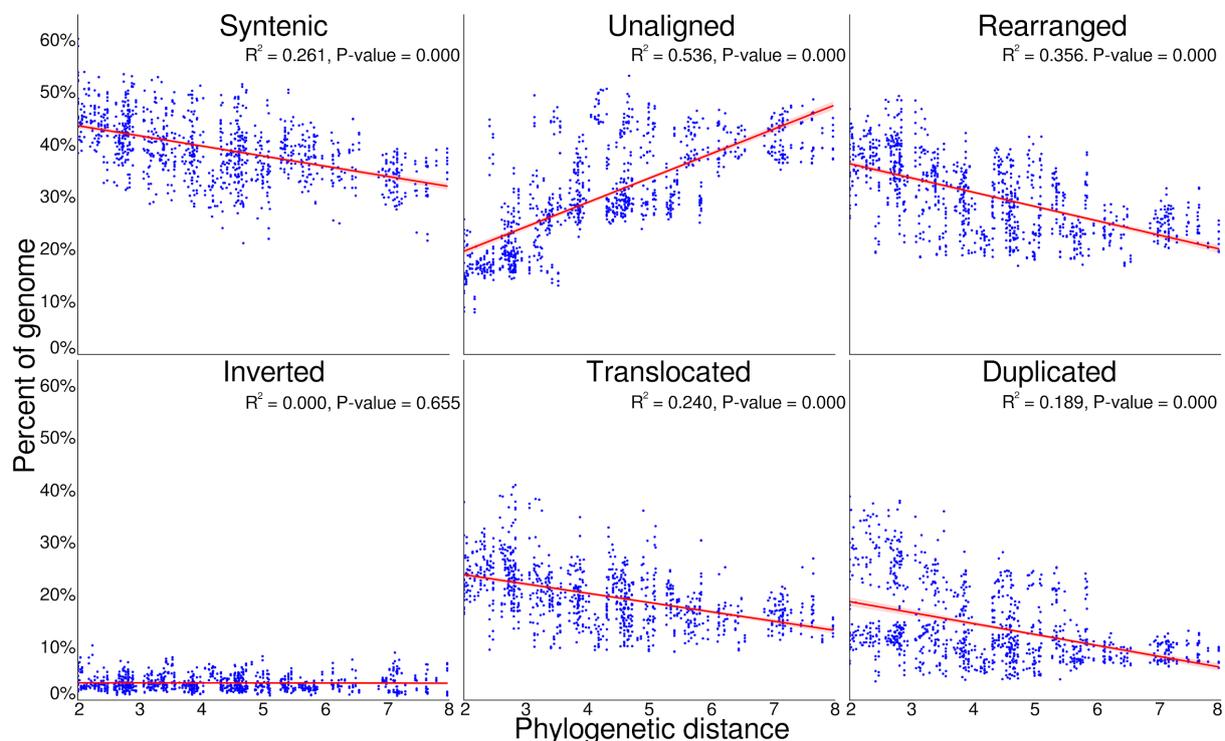
336 Unaligned sequences accumulate through the loss, gain, or divergence of
337 sequences. As genome sizes are similar (average: 552.75 Mbp; standard
338 deviation: 65.62 Mbp) sequence loss and gain are unlikely to fully explain the
339 rapid accumulation of unaligned sequences. Divergence beyond recognition is
340 likely the largest contributing factor. To test which regions were contributing to
341 the growth of unaligned sequences we gathered all alignment identity scores for
342 all syntenic, inverted, translocated, and duplicated regions, in each pairwise
343 alignment. Plotting identities against phylogenetic distance, we examined the
344 rate at which sequences diverge. Syntenic was observed to lose sequence
345 homology more rapidly (R^2 : 0.516), than duplicated (R^2 : 0.236), translocated (R^2 :
346 0.303), and inverted (R^2 : 0.260), Supplementary Figure S20. However, in all
347 cases the regression spanned a very small interval (syntenic: 91.58% - 93.14%,
348 duplicated: 91.48% - 92.63%, translocated: 91.56% - 92.81%, and inverted:

349 91.49% - 92.72%) and none approached our 80% sequence similarity threshold
 350 for alignments.

351

352 Overall, we find that the syntenic proportion of the genome decreases slowly
 353 with increasing divergence time, while the proportion rearranged as duplications
 354 and translocations decrease faster. The loss of homology between synteny,
 355 duplicated and translocated regions leads to a strong increase in the unaligned
 356 portion of the genome (insertions, deletions, and diverged sequences) as
 357 divergence time increases. The loss of duplications and translocations contribute
 358 more to the growth of unaligned than does synteny. We benchmarked our
 359 scaffolding with one species using Hi-C and found consistent results, with a
 360 limitation on the amount of inversions, which may be under reported and
 361 translocations which may contribute less to ongoing genome divergence than
 362 reported, Supplementary Results.

363



364

365

366

367

Figure 3. Pairwise genome conservation and loss, as phylogenetic distance increases. The proportion of both *Eucalyptus* genomes with an alignment pair that was identified as syntenic, rearranged, or unaligned, plotted

368 against the phylogenetic distance of the two genomes. The unaligned proportion
369 is the species-specific fraction of the genome between genome pairs, resulting
370 from either an insertion, deletion, differential inheritance, or sequence
371 divergence. When combined, the proportion of sequence that is syntenic,
372 unaligned, and rearranged equals 100% for each genome within an alignment
373 pair. The rearranged fraction is further broken down into inverted, translocated,
374 and duplicated regions. Phylogenetic distance was calculated as the sum of
375 branch lengths between each genome pair within phylogeny. P-value tests if the
376 slope of the regression line is nonzero.
377

378 Genome-specific and group-wide sequences

379 Unaligned sequences occupied on average 30.65% of each *Eucalyptus* genome
380 within each pairwise alignment. To determine if these sequences were unique to
381 a single genome or shared between multiple, all pairwise alignments for each
382 species were combined and the number of species sharing each base calculated.
383 Subsequently, genome regions that were unique to a genome, shared by
384 multiple genomes or shared by all genomes identified (Figure 4).

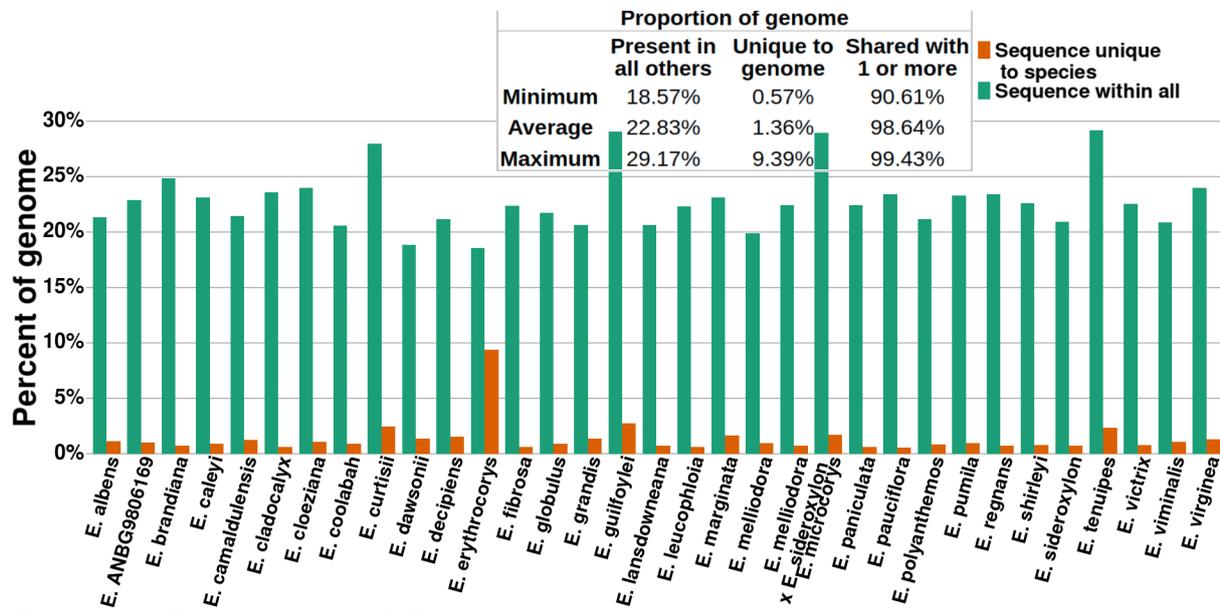
385

386 Genome-specific (unique) sequences occupied an average of 1.36% (241.55
387 Mbp) of the 33 *Eucalyptus* genomes, the remaining 98.64% of sequence was
388 shared by one or more genomes. The proportion of each genome shared by all
389 others averaged 22.83%. This finding mirrors our OG analysis where 2.67% of
390 groups were private, 76.00% dispensable, and 21.33% were core.

391

392 Of note is *E. erythrocorys*, whose genome had a significantly lower proportion of
393 genome-specific sequence and higher proportion of sequence shared by all other
394 genomes. *Eucalyptus erythrocorys* is the sister taxon of all our other genomes
395 within our *Eucalyptus* dataset. Given the age of the divergence between the *E.*
396 *erythrocorys* lineage and its sister lineage, this genome was expected to display
397 a unique pattern in this analysis, however, the extent to which *E. erythrocorys* is
398 different from all others was surprising.

399



400
401 **Figure 4. Proportion of *Eucalyptus* genomes unique and shared by all**
402 **others.** Sequence unique to species is the union of the genome that was
403 classified as unaligned within all pairwise alignments. Sequence within all is the
404 union of the genome that was classified syntenic and rearranged (i.e. common
405 between genomes) within all pairwise alignments.
406

407 Lineage-conserved rearrangements

408 The one-to-one analysis of our *Eucalyptus* genomes has described a genome
409 structure that has become highly fragmented by frequently occurring
410 rearrangements and unaligned regions. As genome structure is inherited by
411 offspring, some of the rearrangements discovered during our analysis are
412 assumed to exist within the genomes of monophyletic groups, i.e. a group of
413 species that have descended from a single ancestral species. Rearrangements
414 found within multiple genomes also help to confirm their validity. To search for
415 evidence of inherited rearrangements we analysed the *Adnataria* section, for
416 which we have the best coverage of genomes (13 genomes, as listed in Figure
417 5). Additionally, using only *Adnataria* genomes should maximise the occurrence
418 of retained rearrangements, as the phylogenetic distances within the *Adnataria*
419 group are relatively low, with many species still hybridising (Delaporte et al.
420 2001).

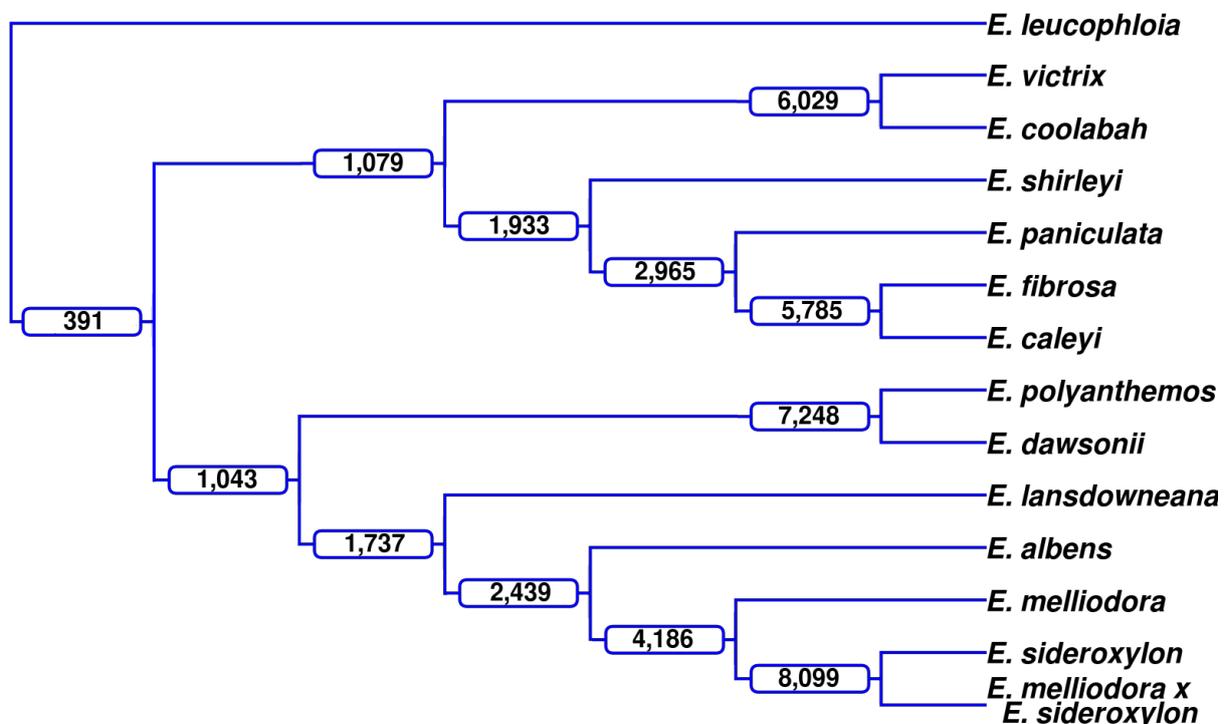
421

422 As all alignments and subsequent annotations are relative to the two species
423 involved, directly comparing the breakpoints of annotations to find common
424 rearrangements is not possible. Therefore, an outgroup genome, *E. leucophloia*,
425 a sister species of the *Adnataria* group, is used for comparison with each of the
426 13 selected genomes. The outgroup genome imposes a single set of genetic
427 coordinates and genome architecture, enabling comparisons of rearrangement
428 breakpoints and subsequent identification of shared rearrangements. Shared
429 rearrangements will contain the same sequence. While this method allows us to
430 find common inversions, translocations, and duplications, it doesn't allow us to
431 find unaligned (insertions, deletions, and highly diverged) regions between our
432 ingroup genomes, as genomes are not being directly compared.

433

434 Comparing the start and end breakpoints (± 50 bp) for events >1 kbp (250,693
435 total rearrangements across all *Adnataria* genomes) identified 58,388 (23.29%)
436 common rearrangements (rearrangements that exist within two or more
437 genomes). Of the 58,388 common rearrangements, 28,059 (48.06%) were
438 shared by two genomes, and 391 (0.67%) were shared by all. The number of
439 common rearrangements quickly decreased as the number of genomes
440 increased (Supplementary Figure S21). Lineage-conserved rearrangements were
441 identified by tracing common rearrangements through *Adnataria's* phylogeny
442 (Figure 5). As expected, more closely related genomes shared the largest
443 number of rearrangements, while more distant genomes shared less.
444 Additionally, as the number of descendant genomes of nodes increased, the
445 number of shared rearrangements also decreased. Inherited rearrangements
446 were identified within the *Adnataria* group. We repeated this analysis twice using
447 *E. brandiana* and *E. cladocalyx* as the outgroup genome achieving similar results
448 (Supplementary Figures S22 and S23).

449



450
 451 **Figure 5. Lineage-conserved rearrangements.** Using an outgroup genome,
 452 *E. leucophloia*, rearrangements were identified that were shared among
 453 members of a lineage, i.e. rearrangements with the same start and end points
 454 within the outgroup genome. At each branch in the dendrogram, the number of
 455 rearrangements shared by all taxa within that clade are labelled. e.g. *E. victrix*
 456 and *E. coolabah* share 6,029 rearrangements.
 457

458 Gene content of synteny, rearrangements, and

459 unaligned

460 To assess whether rearrangements that encompass genes are selected against,
 461 we calculated the proportion of genic (contains a gene/s) and non-genic
 462 (contains no gene/s) rearrangements, as well as syntenic and unaligned events
 463 per genome. Initially, all events too small to contain a gene and genes unplaced
 464 within an orthogroup were removed. A conservative event length of 1 kbp was
 465 used to filter out events, as events smaller than this are unlikely to contain a
 466 gene (Xu et al. 2006). Genes unplaced within an orthogroup are highly dissimilar
 467 to all other gene candidates and may be false positives resulting from incorrect
 468 annotation. The remaining rearrangement, synteny, and unaligned events were

469 examined for the presence of genes placed within an OG and subsequently
470 classed as genic or non-genic.

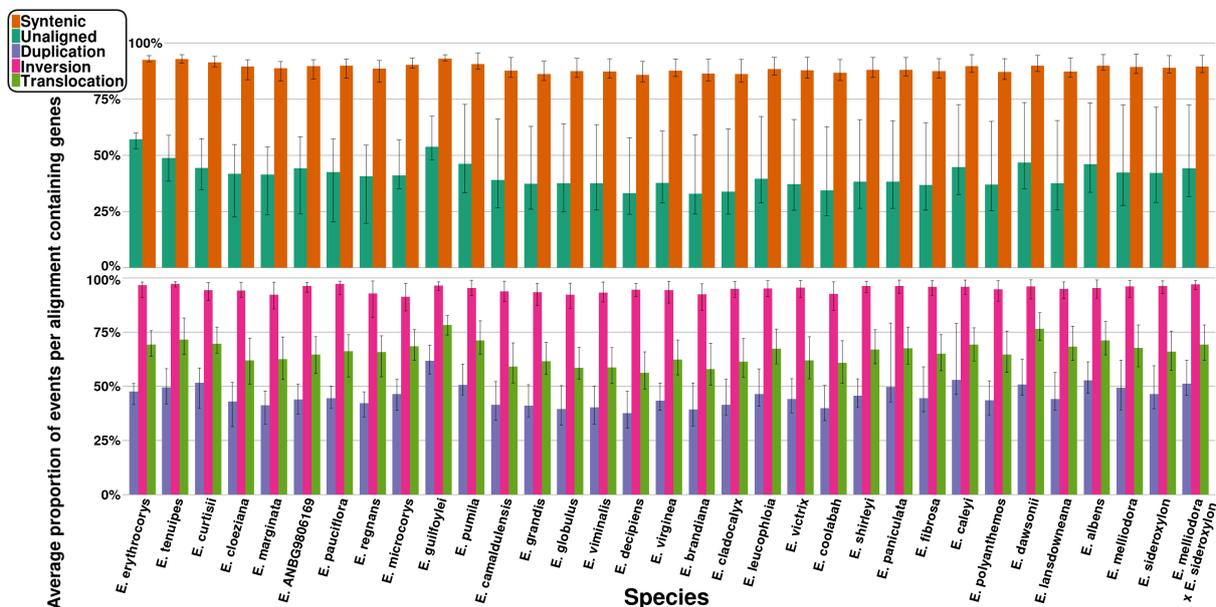
471

472 For each genome, when compared to all other genomes, we calculated the
473 average proportion of genic syntenic, inverted, translocated, duplicated, and
474 unaligned events, and plotted the results (Figure 6). An average of 88.80%
475 (range: 82.52% to 95.57%) genic syntenic events were observed across our
476 genomes. Genic unaligned averaged 41.13% (range: 19.76% to 73.48%), genic
477 inversions averaged 94.93% (range: 81.65% to 99.13%), genic translocations
478 averaged 65.70% (range: 48.77% to 83.98%), and genic duplications averaged
479 45.71% (range: 30.59% to 79.20%).

480

481 Additionally, we analysed the effects of divergence time on the proportion of
482 genic events for all rearrangement types, synteny, and unaligned
483 (Supplementary Figures S24). Results indicated that phylogenetic distance has
484 little to no impact on the proportion of genic syntenic events ($R^2 = 0.162$),
485 inverted events ($R^2 = 0.000$), translocation events ($R^2 = 0.004$), or duplication
486 events ($R^2 = 0.002$). Unaligned was the only event type whose genic proportion
487 was affected by phylogenetic distance. As phylogenetic distance increases,
488 unaligned events become more genic ($R^2 = 0.292$; p -value = 0.000).

489



490 **Figure 6.** Average proportion of genic events for each species genome. The
 491 proportion of genic events was calculated for each pairwise alignment, and
 492 averaged. Error bars indicate the minimum and maximum proportion of genic
 493 events found when aligned to all other genomes.
 494

495 Discussion

496 In this study we created a large collection of wild and naturally evolving high-
 497 quality *Eucalyptus* genomes covering 1-50 million years of divergent evolution.
 498 Using these genomes we find a pattern of genome evolution led by an initial
 499 rapid accumulation of rearrangements and subsequently a slow loss of both
 500 rearranged and syntenic sequences as lineage-specific mutations erode
 501 sequence homology. Rearrangements, likely due to their recombination effects
 502 and subsequent fixation/reduction of alleles (Faria and Navarro 2010), are lost
 503 more rapidly than syntenic regions. Translocations and duplications were the
 504 major disruptors of synteny and were rapidly lost as divergence times increased.
 505 Inversions did not contribute substantially to the loss of synteny or the loss of
 506 rearrangements, instead occurring at consistently low rates across all divergence
 507 times. As genome sizes remain constant and little species-specific sequence
 508 exists across our dataset, loss of existing sequence or gain of new sequence
 509 provides an unlikely explanation for the growth of unaligned sequences as

510 divergence times increase. Hi-C results provided confidence that our scaffolding
511 method was not highly influencing our conclusions. These results showed that
512 translocations initial contribution to genome divergence was significant,
513 however, they don't significantly contribute to ongoing genome divergence. But
514 as this assessment was an almost worst-case scenario for our scaffolding,
515 translocations are still likely a significant contributor to ongoing genome
516 divergence.

517

518 Duplications are a major contributor to functional and genome divergence
519 (Adams and Wendel 2005; Lynch and Conery 2000) especially within plant
520 lineages (Hanada et al. 2008; Van de Peer et al. 2009). We found that
521 duplications were highly abundant between all genomes, and that at smaller
522 divergence times, duplications contributed strongly to genome divergence. As
523 time since divergence increased, the contribution of duplications to genome
524 divergence lessened, becoming overshadowed by unaligned portions of the
525 genome. However, at all phylogenetic distances duplications were a major
526 contributor to genome divergence (occupying on average 14.03% of genomes
527 across an average of 15,350 events). The observed pattern of duplication loss as
528 time since divergence increased was unsurprising, as duplications, while highly
529 important to adaptation and evolution, are rarely conserved (Inoue et al. 2015;
530 Naseeb et al. 2017). Why some duplications are preserved while the majority are
531 lost is speculative; however, theory centres on neofunctionalisation,
532 subfunctionalisation, and novel function evolution (Braasch et al. 2016; Freeling
533 et al. 2015; Lien et al. 2016; Wu and Cox 2019). These hypotheses rely upon the
534 genic properties of duplications, i.e. if duplications don't gain novel function or
535 retain ancestral function, purifying selection will likely result in their removal (Wu
536 and Cox 2019). While duplications were the least genic of all rearrangement
537 types, a significant number (45.71%) were found to contain genes, likely

538 contributing to their preservation. Non-genic duplications, being less visible to
539 selection, are likely to experience increased evolutionary rates (mutations) and
540 genetic drift (Scannell and Wolfe 2008), eventually mutating beyond recognition,
541 and ultimately contributing to the unaligned proportion of alignments.

542 Duplications, both preserved and unpreserved, are likely one of the greatest
543 sources of genome divergence.

544

545 Chromosomal inversions, which are known to be associated with the
546 development of complex phenotypes, local adaptation, and speciation (Arostegui
547 et al. 2019; Lowry and Willis 2010; Twyford and Friedman 2015), were extremely
548 rare between all genomes (average 148 between genomes) and contributed less
549 than all other types of rearrangement to genome divergence. This observation
550 was consistent at all phylogenetic distances: as time since divergence increased,
551 the number of inversions remained constant. A similar finding was made by
552 Hirabayashi and Owens (Hirabayashi and Owens 2023). Inversions likely occur at
553 a high rate within plant genomes (Huang and Rieseberg 2020), however, a low
554 number of inversions was identified, suggesting that inversions are strongly
555 selected against and rarely maintained. To survive underdominant selection, a
556 novel inversion must provide enough selective advantages to outweigh its
557 disadvantages. Inversions may provide a selective advantage by rearranging
558 recombination loci, and linking alleles captured within their bounds. Inversion-
559 linked alleles can be strongly selected for, if adaptive, and rise to high
560 frequencies within populations (Rieseberg 2001; Harringmeyer and Hoekstra
561 2022). Additionally, adaptive alleles linked by inversions can be protected from
562 strong gene flow (Yeaman 2013). Alternatively, inversions may instead hinder
563 adaptation. If selective conditions were to alter, previously adaptive inversions
564 could prevent recombination from producing new allele combinations suitable for
565 the new conditions (Rieseberg 2001). Inversions, due to recombination

566 suppression, also reduce effective population size and increase genetic load, as
567 purifying selection can not purge linked deleterious mutations (Jay et al. 2021).
568 The inversions identified here, which are assumed to have survived selection,
569 were all very large, as expected (Wellenreuther and Bernatchez 2018), with the
570 majority (94.93%) containing genes. Inversions are rare and contribute little to
571 genome divergence, but are highly genic and likely play a significant role in
572 adaptation, evolution, and speciation processes.

573

574 Translocations can have similar genomic effects to inversions (Ortiz-Barrientos et
575 al. 2016), contributing to the development of complex phenotypes, local
576 adaptation, and speciation by disrupting recombination (Martin et al. 2020). As
577 for inversions, novel translocations that survive drift must provide enough
578 selective advantages to outweigh their disadvantages or be removed by
579 purifying or underdominant selection. Translocations were highly abundant
580 between recently diverged, phylogenetically close genomes. As time since
581 divergence increased, translocations reduced in frequency but remained
582 common. Translocations were the most common type of large rearrangement
583 (average size: 11.2 kbp), mirroring results obtained by Martin et al. (Martin et al.
584 2020). Translocations were much more abundant than inversions, especially
585 when genomes have recently diverged, suggesting that translocations are less
586 strongly selected against than inversions, despite having a similar effect on
587 recombination. Additionally, translocations, while highly genic (65.70%) were
588 less genic than inversions (94.93%). The different genomic pattern observed for
589 translocations and inversions is possibly due to the effects of local versus non-
590 local changes to recombination. Meiotic recombination may be more disrupted
591 when reordered recombination loci are close to their location of origin. If true,
592 purifying selection acts more strongly on inversions than on translocations.
593 Translocations are common and along with duplications are a major contributor

594 to genome divergence, possibly aiding in adaptation, evolution, and speciation
595 processes. However, as the effects and mechanisms of translocations have been
596 less studied than other rearrangements (Robberecht et al. 2013), it remains to
597 be seen if they are more likely to have functional/adaptive significance.

598

599 Although new long-read sequencing technologies have accelerated studies on
600 genome structural variations, identifying structural variations still presents
601 challenges. Here, we used RaGOO (Alonge et al. 2019) for reference-guided
602 scaffolding of our Mbp-sized contigs into chromosomes, as obtaining Hi-C data
603 from recalcitrant *Eucalyptus* tissue with high oil content is challenging. We assert
604 that this approach is most suitable given dataset limitations and the well-
605 conserved genome organisation observed in the *Eucalyptus* genus (Potts and
606 Wiltshire 1997; Booth et al. 2015; Grattapaglia et al. 2015; Butler et al. 2017;
607 Supple et al. 2018), as well as in closely related genera within the Myrtaceae
608 family including *Corymbia* (Healey et al. 2021), *Melaluca* (Voelker et al. 2021;
609 Chen et al. 2023), and *Syzygium* (Low et al. 2022; Ouadi et al. 2022). This
610 simplifies reference-guided scaffolding, unlike genera with variable karyotypes,
611 ploidy, and intentional introgressions, such as *Solanum* (Alonge et al. 2019;
612 Razifard et al. 2020). However, reference guided scaffolding may under-
613 represent macro-scale inversions (those larger than contig lengths), as observed
614 when comparing results obtained with reference-scaffolded and Hi-C scaffolded
615 genome assemblies of *E. melliodora*. Despite this limitation, our primary results
616 and conclusions remain unaffected, syntenic and rearrangements were
617 contained within contigs that are orders of magnitude longer (Table 1 and 2). In
618 *Eucalyptus*, we found inversions to be rare, approximately 1% of all structural
619 variations (Table 2), consistent with studies across diverse plant genera
620 (Hirabayashi and Owens 2023) and highly domesticated crop plants like maize
621 (Hufford et al. 2021). Obtaining Hi-C data in more species may help resolve

622 large-scale inversions, it can introduce errors (Alonge et al. 2019) and still
623 represents prediction and hypothesis. An alternative strategy is single-cell/single-
624 strand genome sequencing (Falconer et al. 2012), which was found to be one of
625 the most reliable methods to detect large-scale inversions in human genomes
626 (Chaisson et al. 2019). As long-read sequencing technologies advance, assembly
627 of telomere-to-telomere genomes, independent of Hi-C data and genome
628 scaffolding, will greatly enhance genome studies and overcome technical
629 challenges in structural variation discovery. These advances are exemplified by
630 long-read de novo assemblers such as hifiasm (UL) (Cheng et al. 2023)
631 and Verkko (Rautiainen et al. 2023).

632

633 To further investigate the potential importance of the syntenic, rearranged, and
634 unaligned genome regions identified in our study further research using genome-
635 wide association studies (GWAS) of phenotypes measured on seedlings in pots or
636 field trials, as well as landscape and genome-wide genotyping for Genome
637 Environment Association (GEA) scans for adaptive rearrangements are needed.
638 Within species derived rearrangements are predicted to be predominantly
639 neutral and exist at low frequencies, while others rising to higher frequencies
640 could be true lineage-specific adaptive rearrangements. With additional genomes
641 from populations, the frequency of rearrangements within each species could be
642 assessed. This would provide insight into the functional significance of the
643 widespread genomic rearrangements we have found and potentially identify
644 rearrangements conferring adaptive traits across the landscape.

645

646 *Eucalyptus* contains over 800 species that exist across a wide geographic and
647 environmental range, while retaining a largely conserved karyotype (Potts and
648 Wiltshire 1997; Booth et al. 2015; Grattapaglia et al. 2015; Butler et al. 2017;
649 Supple et al. 2018), which makes the genus ideal to study plant genome

650 evolution. Here we assembled representative genomes of 33 species, creating
651 one of the most comprehensive datasets to study plant genome evolution. These
652 genomes provide a genus-wide resource to study genome rearrangements, and
653 support future *Eucalyptus* research that require genomic references. Our findings
654 suggest that following divergence, genome architecture is highly fragmented
655 predominantly by rearrangements. As genomes continue to diverge, genome
656 architecture continues to be slowly lost. Additionally, as genomes diverge they
657 increasingly become unalignable due to the divergence of duplications and
658 translocations. Syntenic regions also contribute to the growing unalignable
659 proportion of genomes, but at a slower rate than rearrangements. Duplications
660 and translocations are potentially the greatest contributors to functional and
661 genome divergence, aiding in the development of complex phenotypes, and local
662 adaptation. Inversions occur at consistently low rates, contributing little to
663 genome architecture loss or accumulation of unalignable sequences. However,
664 inversions were highly genic, much more so than either duplications or
665 translocations, and likely also play a crucial role in the development of complex
666 phenotypes, and local adaptation. Genome architecture results from a complex
667 interaction of positive, neutral, and negative forces, all of which contribute to the
668 evolution, divergence, and adaptability of species (Koonin 2009; Huang and
669 Rieseberg 2020; Mérot et al. 2020). However, owing to technical limitations, the
670 evolution of genome architecture and its role within biology is not well
671 understood (Lynch et al. 2011; Cortés et al. 2018; Jiggins 2019). Here, by
672 describing the pattern of genome architecture as time since divergence
673 increases of 33 *Eucalyptus* genomes, we contribute to a better understanding of
674 the evolution of plant genomes. Rearrangements, along with polyploidy, TEs, and
675 other genome evolutionary mechanisms, play an important role in plant genome
676 evolution (Galindo-González et al. 2017; Marques et al. 2019; Meudt et al. 2021).

677 Further research in other plant lineages is required to assess the prominence of
678 rearrangements upon genome evolution.

679

680 Materials and Methods

681 Sampling

682 *Eucalyptus* species used in this study were collected throughout multiple
683 locations in Australia, which are detailed in supplementary results. The majority
684 of collected species are living collections with accession numbers at the
685 Australian National Botanic Gardens (Canberra, Australian Capital Territory
686 (ACT)) and Currency Creek Arboretum (Currency Creek, South Australia).
687 Additional samples were sourced from the Australian National University (Acton,
688 ACT) the National Arboretum Canberra (Molonglo Valley, ACT), the University of
689 Tasmania Herbarium (Sandy Bay, Tasmania) and within *Eucalyptus* woodlands of
690 southern Tasmania. Leaves were placed in plastic zip-lock bags, lightly sprayed
691 with water to keep them moist and transported to the lab as soon as possible,
692 where they were washed with water and stored at -80°C until DNA extraction.

693

694 DNA extraction, sequencing, and basecalling

695 To extract high-molecular weight DNA from recalcitrant *Eucalyptus* samples, we
696 developed two methods. Initially we combined a protocol to purify nuclei with
697 hexylene glycol (Bolger et al. 2014) with a magnetic bead-based DNA extraction
698 protocol (Mayjonade et al. 2016), which was further developed and is available
699 on Protocols.io in detail (Jones and Borevitz 2019). This was further optimised
700 and developed, which led to the second method of adopting a sorbitol pre-wash
701 of homogenate (Inglis et al. 2018) to wash crude nuclei instead of isolating pure
702 nuclei, followed by a magnetic bead-based DNA extraction, according to (Jones et

703 al. 2021). We found this method to be more time and resource efficient, hence
704 we switched to this method for all subsequent high-molecular weight DNA
705 extractions. For each *Eucalyptus* sample, the method which was used is listed
706 within the supplementary material, Table Supplementary Table S1, with the two
707 methods being referred to as nuclei and sorbitol respectively.

708

709 After isolating high-molecular weight DNA, we further purified and size selected
710 the DNA by using a PippinHT (Sage Science). The DNA was size selected for
711 fragments ≥ 20 kb or ≥ 40 kb depending on DNA yield and molecular weight,
712 which is listed in the supplementary material, Table S1, for each sample. Two
713 Oxford Nanopore Technologies long-read native DNA sequencing libraries were
714 prepared for each species according to the manufacturer's protocol 1D genomic
715 DNA by ligation (SQK-LSK109). *E. marginata* was an exception, which had one
716 ligation library as described but the second was a transposome library prep,
717 according to the manufacturer's protocol for rapid sequencing (SQK-RAD004).
718 Sequencing was performed on MinION Mk1B devices using two FLO-MIN106D
719 R9.4.1 flow cells per species. Sequencing output was improved when ONT Flow
720 Cell Wash Kits (EXP-WSH003 and EXP-WSH004) were made available, whereby
721 flow cells were washed when sequencing declined, primed again and more
722 library was loaded, according to the manufacturer's instructions. After
723 sequencing was complete, the FAST5 reads were basecalled with ONT Guppy
724 (versions: 3.3.0, 4.0.11, 4.0.14, and 4.0.15; See Supplementary Table S14 for per
725 species versions).

726

727 We complemented the long-read sequencing with highly accurate Illumina short-
728 read sequencing for later use in genome polishing of the long-read de novo
729 assemblies. Illumina short-read, whole-genome DNA sequencing libraries were
730 generated using a cost-optimised, transposome protocol based on Illumina

731 Nextera DNA prep methods (Jones et al. 2023). The pooled libraries were then
732 size selected for fragments with insert sizes between 350 and 600 bp with a
733 PippinHT (Sage Science). Multiplexed sequencing with other projects was
734 performed on a NovaSeq 6000 (Illumina), using a lane of an S4 flow cell with a
735 300 cycle kit (150 bp paired-end sequencing), at the Biomolecular Resource
736 Facility, Australian National University, Australian Capital Territory, Australia.
737

738 *De novo* assembly

739 *De novo* assembly and annotation was performed using the long-read *de novo*
740 plant assembly protocol developed by Ferguson et al. (Ferguson et al. 2022a).
741 Briefly, FASTQ reads are quality screened, removing DNA control strand,
742 sequencing adaptors and low quality read ends (the first and last 200 bp), short
743 reads (>1 kbp in length), and low quality reads (average quality <Q7), using the
744 NanoPack set of tools (De Coster et al. 2018). Curated reads are next assembled
745 using the long-read assembler Canu (versions: 1.9 and 2.0) (Koren et al. 2017),
746 which assembles high-quality *Eucalyptus* genomes (Ferguson et al. 2022b).
747 Assemblies were filtered of contamination (non-plant contigs), assembly artefact,
748 plasmid, and haplotig (contigs that span the same genomic region but originate
749 from different parental chromosomes) contigs using Blobtools (Laetsch and
750 Blaxter 2017) and Purge Haplotigs (version: 1.1.0) (Roach et al. 2018). Next, all
751 assemblies were long-read and then short-read polished, using assembly reads
752 and Illumina reads originating from the same individual as used for assembly.
753 Long-read polishing was performed with Racon (Vaser et al. 2017), and short-
754 read with Pilon (version: 1.3.1) (Walker et al. 2014). Long-read polishing made
755 use of the long-read aligner minimap2 (version: 2.17) (Li 2018), while short-read
756 polishing used (version: 0.7.17) (Li 2013). Next, assemblies were filtered to
757 remove all contigs less than 1 kbp in length. We chose this contig length

758 threshold so as to maximise genome contiguity while removing all contigs too
759 small to contain a gene. Finally, assemblies were scaffolded using homology to *E.*
760 *grandis* (Myburg et al. 2014). Scaffolding was performed with RaGOO (version
761 1.1) (Alonge et al. 2019) and minimap2.

762

763 After assembly all genomes were quality assessed using Benchmarking Universal
764 Single-Copy Orthologs (BUSCO; version 5; database: eudicots_odb10.2020-09-
765 10) (Manni et al. 2021), long terminal repeat assembly index (version: 2.9.0; LAI)
766 (Ou et al. 2018), and assembly statistics.

767

768 Transposon and gene annotation, and gene ortho

769 grouping

770 Genome repeat and gene annotation was also performed using the long-read *de*
771 *novo* plant assembly protocol developed by Ferguson et al. (Ferguson et al.
772 2022a). Firstly, *de novo* repeat libraries were created for each genome using
773 EDTA (version: 1.9.6) (Ou et al. 2019), subsequently all genomes were repeat
774 annotated with RepeatMasker (version: 4.0.9) (Smit et al. 2020). All genomes
775 were repeat soft-masked and subsequently annotated for genes. Gene
776 annotation was performed with BRAKER (version 2.1.5; Brůna *et al.*, 2021) using
777 GeneMark-EP (version: 4) (Brůna, Lomsadze and Borodovsky, 2020). Gene
778 transcript sequences for model training were obtained from the National Center
779 for Biotechnology Information (NCBI) (Sayers et al. 2021). Included in gene
780 training data were all Myrtaceae (Taxonomy ID: 3931) and *Arabidopsis thaliana*
781 (Taxonomy ID: 3702) transcripts. All gene candidates were grouped into
782 orthogroups using OrthoFinder (version: 2.5.4) (Emms and Kelly 2019). Using
783 DIAMOND (Buchfink et al. 2021) OrthoFinder aligned all gene transcripts,
784 grouping those with >40% identity and achieving an e-score < 0.001.

785

786 **Genome synteny, rearrangement and unaligned**787 **annotation**

788 Identification of all shared sequences began by aligning all pairwise
789 combinations of genomes with the MUMmer (version 3) (Kurtz et al. 2004) tool
790 NUCmer (parameters: --maxmatch -l 40 -b 500 -c 200). NUCmer first identifies all
791 shared 40-mers between genomes and their locations. Next, 40-mers within
792 500bp are clustered, creating a list of collinear blocks or alignments. Lastly,
793 using MUMmer's delta-filter tool alignments are filtered, removing all alignments
794 less than 200 bp in length and less than 80% similar. A low 80% alignment
795 similarity score was used as *Eucalyptus* are highly heterozygous (Murray et al.
796 2019), and a more stringent similarity score may incorrectly filter out real
797 alignments.

798

799 Having identified all shared sequences we next annotated all syntenic,
800 rearranged (inverted, translocated, and duplicated), and unaligned (sequence
801 that only exists in one genome, resulting from either an insertion, deletion or
802 sequence divergence) sequence between pairwise genomes using SyRI (version:
803 1.5) (Goel et al. 2019). SyRI's use of a directed acyclic graph results in genomes
804 being annotated for smaller regions, which when occurring in an unbroken series
805 of a single type, get combined. The resulting output includes both levels of
806 annotations, smaller more fragmented, and larger and more contiguous. We
807 make use of the larger and more continuous alignments. Additionally, we
808 combined inverted duplications with duplications, and inverted translocations
809 with translocations.

810

811 Phylogeny

812 Using highly conserved and single-copy genes, BUSCO genes, we built a eucalypt
813 phylogenetic tree describing the evolutionary relationships between all genomes
814 included in this study. The phylogenetic tree included four previously and
815 identically assembled genomes for *E. albens*, *E. melliodora*, *E. sideroxylon*, and
816 *C. calophylla*, creating a dataset of 36 genomes. To begin, FASTA sequences for
817 all single-copy BUSCO genes found within 30+ genomes were collected. Using
818 macse (version 2.03) (Ranwez *et al.*, 2018) multi-sequence alignment was
819 performed individually on all genes. As errors within gene MSAs will subsequently
820 lead to errors in phylogenetic inferences, we trimmed and filtered all gene MSAs.
821 Gene sequence errors were detected and removed using HmmCleaner version:
822 0.180750; (Di Franco *et al.* 2019). HmmCleaner uses a profile-hidden Markov
823 model to identify sequence segments that poorly fit the gene MSA and
824 subsequently removes them. Errors resulting from poor alignments were
825 removed using report2AA (parameters: -min_NT_to_keep_seq 30, -
826 min_seq_to_keep_site 4,-dist_isolate_AA 3, -min_homology_to_keep_seq 0.5, -
827 min_percent_NT_at_ends 0.7) from the macse program. report2AA removed sites
828 within MSAs that included less than 30 genomes, had less than 4 informative
829 characters, or had isolated sites (site was more than 3 characters away from the
830 next non-gap character). Additionally, report2AA removed genome from MSAs
831 that had less than 50% homology to another genome within the MSA, and
832 trimmed both MSA ends that had less than 70% of aligned sites as nucleotides
833 (i.e. 26+ genomes had to have a non-gap character). Additionally, as a result of
834 filtering and trimming MSAs of low quality are removed.

835

836 Individual gene trees were constructed for all filtered and trimmed MSAs using
837 IQ-TREE (version: 1.6.12) (Nguyen *et al.*, 2015) . Finally, all gene trees were

838 concatenated into a single file, from which a species tree was generated using
839 Astral III (version 5.7.3) (Zhang et al. 2018). The resulting species tree was
840 manually rooting at the *Angophora/Corymbia* and *Eucalyptus* branch, using
841 Figtree (version: 1.4.4) (Rambaut, no date).

842

843 Data access

844 Sequencing data and reference genomes generated in this study have been
845 submitted to the NCBI BioProject database
846 (<https://www.ncbi.nlm.nih.gov/bioproject/>) under accession number
847 PRJNA509734. Gene predictions, repeat annotations, and SyRI annotations
848 generated in this study are available on FigShare
849 ([https://figshare.com/projects/Plant_genome_evolution_in_the_genus_Eucalyptus_](https://figshare.com/projects/Plant_genome_evolution_in_the_genus_Eucalyptus_driven_by_structural_rearrangements_that_promote_sequence_divergence/97010)
850 [driven_by_structural_rearrangements_that_promote_sequence_divergence/](https://figshare.com/projects/Plant_genome_evolution_in_the_genus_Eucalyptus_driven_by_structural_rearrangements_that_promote_sequence_divergence/97010)
851 [97010](https://figshare.com/projects/Plant_genome_evolution_in_the_genus_Eucalyptus_driven_by_structural_rearrangements_that_promote_sequence_divergence/97010)). All of the analysis scripts used in this study are available at GitHub
852 (<https://github.com/fergsc/33-Eucs>) and as Supplemental Scripts.

853

854 Competing interest statement

855 The authors declare that they have no competing interests.

856

857 Acknowledgements

858 We would like to thank the Australian National Botanic Gardens in Canberra,
859 Australia for providing plant samples and associated metadata. This research
860 acknowledges the support provided by the Director of National Parks, the park
861 staff of the Australian National Botanic Gardens, and Parks Australia. The views

862 expressed in this document do not necessarily represent the views of the
863 Australian Government.

864

865 We thank Dean Nicolle, owner of the Currency Creek Arboretum, South Australia,
866 for providing samples and support for this project. We also thank Yoav Daniel
867 Bar-Ness, Giant Tree Expeditions, for collecting *E. regnans* (the Centurion). We
868 also thank Tamera Beath, David Stanley, Cynthia Torkel, Rob Lanfear, Wang
869 Weiwen and Brad Potts for their support and collecting samples.

870

871 We would like to thank the ACRF Biomolecular Resource Facility at the John
872 Curtin School of Medical Research, ANU in Canberra, Australia, where Oxford
873 Nanopore Technologies PromethION and Illumina NovaSeq 6000 sequencing was
874 conducted. This research acknowledges the support provided by NCRIS-enabled
875 Bioplatforms Australia infrastructure.

876

877 Computational resources were provided by the Australian Government through
878 the National Computational Infrastructure (NCI) under the ANU Merit Allocation
879 Scheme.

880

881 **Funding information**

882 This research was supported by the Australian Research Council (Project code:
883 CE140100008 and DP150103591). S.F. was supported by Australian Government
884 Research Training Program scholarships.

885

886 **Author contributions**

887 S.F. led the project and ran all the analysis. A.J. managed, developed, and
888 performed DNA sampling and sequencing. The project was conceived and

889 designed by all authors. S.F. wrote the first manuscript draft. All authors
 890 contributed to writing and review of the final manuscript.

891

892 References

- 893 Adams KL, Wendel JF. 2005. Polyploidy and genome evolution in plants. *Curr*
 894 *Opin Plant Biol* **8**: 135–141.
- 895 Ahrens CW, Murray K, Mazanec RA, Ferguson S, Bragg J, Jones A, Tissue DT,
 896 Byrne M, Borevitz JO, Rymer PD. 2021. Genomic constraints to drought
 897 adaptation. *bioRxiv*.
 898 <https://www.biorxiv.org/content/early/2021/08/08/2021.08.07.455511>.
- 899 Alonge M, Soyk S, Ramakrishnan S, Wang X, Goodwin S, Sedlazeck FJ, Lippman
 900 ZB, Schatz MC. 2019. RaGOO: fast and accurate reference-guided
 901 scaffolding of draft genomes. *Genome Biol* **20**: 224.
- 902 Alonge M, Wang X, Benoit M, Soyk S, Pereira L, Zhang L, Suresh H, Ramakrishnan
 903 S, Maumus F, Ciren D, et al. 2020. Major Impacts of Widespread Structural
 904 Variation on Gene Expression and Crop Improvement in Tomato. *Cell* **182**:
 905 145-161.e23.
- 906 Amarasinghe SL, Su S, Dong X, Zappia L, Ritchie ME, Gouil Q. 2020.
 907 Opportunities and challenges in long-read sequencing data analysis.
 908 *Genome Biol* **21**: 30.
- 909 Arostegui MC, Quinn TP, Seeb LW, Seeb JE, McKinney GJ. 2019. Retention of a
 910 chromosomal inversion from an anadromous ancestor provides the
 911 genetic basis for alternative freshwater ecotypes in rainbow trout. *Mol Ecol*
 912 **28**: 1412–1427.
- 913 Bayer PE, Golicz AA, Scheben A, Batley J, Edwards D. 2020. Plant pan-genomes
 914 are the new reference. *Nat Plants* **6**: 914–920.
- 915 Bezemer N, Krauss SL, Phillips RD, Roberts DG, Hopper SD. 2016. Paternity
 916 analysis reveals wide pollen dispersal and high multiple paternity in a
 917 small isolated population of the bird-pollinated *Eucalyptus caesia*
 918 (Myrtaceae). *Heredity* **117**: 460–471.
- 919 Bolger A, Scossa F, Bolger ME, Lanz C, Maumus F, Tohge T, Quesneville H,
 920 Alseekh S, Sørensen I, Lichtenstein G, et al. 2014. The genome of the
 921 stress-tolerant wild tomato species *Solanum pennellii*. *Nat Genet* **46**:
 922 1034–1038.
- 923 Booth TH, Broadhurst LM, Pinkard E, Prober SM, Dillon SK, Bush D, Pinyopusarerk
 924 K, Doran JC, Ivkovich M, Young AG. 2015. Native forests and climate
 925 change: Lessons from eucalypts. *For Ecol Manag* **347**: 18–29.
- 926 Braasch I, Gehrke AR, Smith JJ, Kawasaki K, Manousaki T, Pasquier J, Amores A,
 927 Desvignes T, Batzel P, Catchen J, et al. 2016. The spotted gar genome
 928 illuminates vertebrate evolution and facilitates human-teleost
 929 comparisons. *Nat Genet* **48**: 427–437.

- 930 Brůna T, Hoff KJ, Lomsadze A, Stanke M, Borodovsky M. 2021. BRAKER2:
931 automatic eukaryotic genome annotation with GeneMark-EP+ and
932 AUGUSTUS supported by a protein database. *NAR Genomics Bioinforma* **3**.
933 [https://academic.oup.com/nargab/article/doi/10.1093/nargab/lqaa108/606](https://academic.oup.com/nargab/article/doi/10.1093/nargab/lqaa108/6066535)
934 6535 (Accessed August 5, 2021).
- 935 Brůna T, Lomsadze A, Borodovsky M. 2020. GeneMark-EP+: eukaryotic gene
936 prediction with self-training in the space of genes and proteins. *NAR*
937 *Genomics Bioinforma* **2**: lqaa026.
- 938 Buchfink B, Reuter K, Drost H-G. 2021. Sensitive protein alignments at tree-of-life
939 scale using DIAMOND. *Nat Methods* **18**: 366–368.
- 940 Burns R, Mandáková T, Gunis J, Soto-Jiménez LM, Liu C, Lysak MA, Novikova PY,
941 Nordborg M. 2021. Gradual evolution of allopolyploidy in *Arabidopsis*
942 *suecica*. *Nat Ecol Evol* **5**: 1367–1381.
- 943 Butler JB, Vaillancourt RE, Potts BM, Lee DJ, King GJ, Baten A, Shepherd M,
944 Freeman JS. 2017. Comparative genomics of *Eucalyptus* and *Corymbia*
945 reveals low rates of genome structural rearrangement. *BMC Genomics* **18**:
946 397.
- 947 Chaisson MJP, Sanders AD, Zhao X, Malhotra A, Porubsky D, Rausch T, Gardner
948 EJ, Rodriguez OL, Guo L, Collins RL, et al. 2019. Multi-platform discovery of
949 haplotype-resolved structural variation in human genomes. *Nat Commun*
950 **10**: 1784.
- 951 Chen SH, Martino AM, Luo Z, Schwessinger B, Jones A, Tolessa T, Bragg JG,
952 Tobias PA, Edwards RJ. 2023. A high-quality pseudo-phased genome for
953 *Melaleuca quinquenervia* shows allelic diversity of NLR-type resistance
954 genes. *GigaScience* **12**: giad102.
- 955 Cheng H, Asri M, Lucas J, Koren S, Li H. 2023. Scalable telomere-to-telomere
956 assembly for diploid and polyploid genomes with double graph.
- 957 Cortés AJ, Skeen P, Blair MW, Chacón-Sánchez MI. 2018. Does the Genomic
958 Landscape of Species Divergence in *Phaseolus* Beans Coerce Parallel
959 Signatures of Adaptation and Domestication? *Front Plant Sci* **9**.
960 <https://www.frontiersin.org/articles/10.3389/fpls.2018.01816> (Accessed
961 March 27, 2023).
- 962 Crkvenjakov R, Heng HH. 2022. Further illusions: On key evolutionary
963 mechanisms that could never fit with Modern Synthesis. *Prog Biophys Mol*
964 *Biol* **169-170**: 3–11.
- 965 Dawkins R. 1976. *The selfish gene*. Oxford University Press.
- 966 Dawson DA, Åkesson M, Burke T, Pemberton JM, Slate J, Hansson B. 2007. Gene
967 Order and Recombination Rate in Homologous Chromosome Regions of
968 the Chicken and a Passerine Bird. *Mol Biol Evol* **24**: 1537–1552.
- 969 De Coster W, D’Hert S, Schultz DT, Cruts M, Van Broeckhoven C. 2018.
970 NanoPack: visualizing and processing long-read sequencing data ed. B.
971 Berger. *Bioinformatics* **34**: 2666–2669.
- 972 Delaporte KL, Conran JG, Sedgley M. 2001. Interspecific Hybridization within
973 *Eucalyptus* (Myrtaceae): Subgenus *Symphomyrtus* , Sections *Bisectae*

- 974 and *Adnataria*. *Int J Plant Sci* **162**: 1317–1326.
- 975 Derežanin L, Blažytė A, Dobrynin P, Duchêne DA, Grau JH, Jeon S, Kliver S, Koepfli
976 K-P, Meneghini D, Preick M, et al. 2022. Multiple types of genomic variation
977 contribute to adaptive traits in the mustelid subfamily Guloninae. *Mol Ecol*
978 **31**: 2898–2919.
- 979 Di Franco A, Poujol R, Baurain D, Philippe H. 2019. Evaluating the usefulness of
980 alignment filtering methods to reduce the impact of errors on evolutionary
981 inferences. *BMC Evol Biol* **19**: 21.
- 982 Dixon JR, Gorkin DU, Ren B. 2016. Chromatin Domains: The Unit of Chromosome
983 Organization. *Mol Cell* **62**: 668–680.
- 984 Emms DM, Kelly S. 2019. OrthoFinder: phylogenetic orthology inference for
985 comparative genomics. *Genome Biol* **20**: 238.
- 986 Falconer E, Hills M, Naumann U, Poon SSS, Chavez EA, Sanders AD, Zhao Y, Hirst
987 M, Lansdorp PM. 2012. DNA template strand sequencing of single-cells
988 maps genomic rearrangements at high resolution. *Nat Methods* **9**: 1107–
989 1112.
- 990 Faria R, Navarro A. 2010. Chromosomal speciation revisited: rearranging theory
991 with pieces of evidence. *Trends Ecol Evol* **25**: 660–669.
- 992 Ferguson S, Jones A, Borevitz J. 2022a. Plant assemble - Plant de novo genome
993 assembly, scaffolding and annotation for genomic studies. *protocols.io*.
994 <https://dx.doi.org/10.17504/protocols.io.81wgb6zk3lpk/v1> (Accessed
995 August 4, 2022).
- 996 Ferguson S, Jones A, Murray K, Schwessinger B, Borevitz JO. 2022b. Interspecies
997 genome divergence is predominantly due to frequent small scale
998 rearrangements in *Eucalyptus*. *Mol Ecol* **n/a**.
999 <https://onlinelibrary.wiley.com/doi/abs/10.1111/mec.16608> (Accessed
1000 August 8, 2022).
- 1001 Feulner PGD, De-Kayne R. 2017. Genome evolution, structural rearrangements
1002 and speciation. *J Evol Biol* **30**: 1488–1490.
- 1003 Freeling M, Scanlon MJ, Fowler JE. 2015. Fractionation and subfunctionalization
1004 following genome duplications: mechanisms that drive gene content and
1005 their consequences. *Curr Opin Genet Dev* **35**: 110–118.
- 1006 Galindo-González L, Mhiri C, Deyholos MK, Grandbastien M-A. 2017. LTR-
1007 retrotransposons in plants: Engines of evolution. *Gene* **626**: 14–25.
- 1008 Goel M, Sun H, Jiao W-B, Schneeberger K. 2019. SyRI: finding genomic
1009 rearrangements and local sequence differences from whole-genome
1010 assemblies. *Genome Biol* **20**: 277.
- 1011 Grattapaglia D, Mamani EMC, Silva-Junior OB, Faria DA. 2015. A novel genome-
1012 wide microsatellite resource for species of *Eucalyptus* with linkage-to-
1013 physical correspondence on the reference genome sequence. *Mol Ecol*
1014 *Resour* **15**: 437–448.
- 1015 Hanada K, Zou C, Lehti-Shiu MD, Shinozaki K, Shiu S-H. 2008. Importance of
1016 lineage-specific expansion of plant tandem duplicates in the adaptive
1017 response to environmental stimuli. *Plant Physiol* **148**: 993–1003.

- 1018 Hardigan MA, Feldmann MJ, Lorant A, Bird KA, Famula R, Acharya C, Cole G,
1019 Edger PP, Knapp SJ. 2020. Genome Synteny Has Been Conserved Among
1020 the Octoploid Progenitors of Cultivated Strawberry Over Millions of Years
1021 of Evolution. *Front Plant Sci* **10**: 1789.
- 1022 Harringmeyer OS, Hoekstra HE. 2022. Chromosomal inversion polymorphisms
1023 shape the genomic landscape of deer mice. *Nat Ecol Evol* **6**: 1965–1979.
- 1024 Healey AL, Shepherd M, King GJ, Butler JB, Freeman JS, Lee DJ, Potts BM, Silva-
1025 Junior OB, Baten A, Jenkins J, et al. 2021. Pests, diseases, and aridity have
1026 shaped the genome of *Corymbia citriodora*. *Commun Biol* **4**: 537.
- 1027 Heger A, Ponting CP. 2007. Evolutionary rate analyses of orthologs and paralogs
1028 from 12 *Drosophila* genomes. *Genome Res* **17**: 1837–1849.
- 1029 Heng HHQ. 2009. The genome-centric concept: resynthesis of evolutionary
1030 theory. *BioEssays* **31**: 512–525.
- 1031 Heng HHQ, Goetze S, Ye CJ, Liu G, Stevens JB, Bremer SW, Wykes SM, Bode J,
1032 Krawetz SA. 2004. Chromatin loops are selectively anchored using
1033 scaffold/matrix-attachment regions. *J Cell Sci* **117**: 999–1008.
- 1034 Hirabayashi K, Owens GL. 2023. The rate of chromosomal inversion fixation in
1035 plant genomes is highly variable. *Evolution* **77**: 1117–1130.
- 1036 Huang K, Rieseberg LH. 2020. Frequency, Origins, and Evolutionary Role of
1037 Chromosomal Inversions in Plants. *Front Plant Sci* **11**.
1038 <https://www.frontiersin.org/articles/10.3389/fpls.2020.00296> (Accessed
1039 December 21, 2022).
- 1040 Hufford MB, Seetharam AS, Woodhouse MR, Chougule KM, Ou S, Liu J, Ricci WA,
1041 Guo T, Olson A, Qiu Y, et al. 2021. De novo assembly, annotation, and
1042 comparative analysis of 26 diverse maize genomes. 9.
- 1043 Inglis PW, Pappas M de CR, Resende LV, Grattapaglia D. 2018. Fast and
1044 inexpensive protocols for consistent extraction of high quality DNA and
1045 RNA from challenging plant and fungal samples for high-throughput SNP
1046 genotyping and sequencing applications. *PLOS ONE* **13**: e0206085.
- 1047 Inoue J, Sato Y, Sinclair R, Tsukamoto K, Nishida M. 2015. Rapid genome
1048 reshaping by multiple-gene loss after whole-genome duplication in teleost
1049 fish suggested by mathematical modeling. *Proc Natl Acad Sci* **112**: 14918–
1050 14923.
- 1051 Jay P, Chouteau M, Whibley A, Bastide H, Parrinello H, Llaurens V, Joron M. 2021.
1052 Mutation load at a mimicry supergene sheds new light on the evolution of
1053 inversion polymorphisms. *Nat Genet* **53**: 288–293.
- 1054 Jensen RA. 2001. Orthologs and paralogs - we need to get it right. *Genome Biol*
1055 **2**: interactions1002.1-interactions1002.3.
- 1056 Jiao W-B, Schneeberger K. 2020. Chromosome-level assemblies of multiple
1057 *Arabidopsis* genomes reveal hotspots of rearrangements with altered
1058 evolutionary dynamics. *Nat Commun* **11**: 989.
- 1059 Jiggins CD. 2019. Can genomics shed light on the origin of species? *PLOS Biol* **17**:
1060 e3000394.

- 1061 Jones A, Borevitz J. 2019. Nuclear DNA purification from recalcitrant plant species
1062 for long-read sequencing. *protocols.io*.
1063 [https://www.protocols.io/view/nuclear-dna-purification-from-recalcitrant-](https://www.protocols.io/view/nuclear-dna-purification-from-recalcitrant-plant-s-28bghsn)
1064 [plant-s-28bghsn](https://www.protocols.io/view/nuclear-dna-purification-from-recalcitrant-plant-s-28bghsn) (Accessed January 11, 2022).
- 1065 Jones A, Stanley D, Ferguson S, Schwessinger B, Borevitz J, Warthmann N. 2023.
1066 Cost-conscious generation of multiplexed short-read DNA libraries for
1067 whole-genome sequencing. *PLOS ONE* **18**: 1-7.
- 1068 Jones A, Torkel C, Stanley D, Nasim J, Borevitz J, Schwessinger B. 2021. High-
1069 molecular weight DNA extraction, clean-up and size selection for long-read
1070 sequencing ed. M. Eppinger. *PLOS ONE* **16**: e0253830.
- 1071 Koonin EV. 2009. Evolution of genome architecture. *Int J Biochem Cell Biol* **41**:
1072 298-306.
- 1073 Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. 2017. Canu:
1074 scalable and accurate long-read assembly via adaptive *k*-mer weighting
1075 and repeat separation. *Genome Res* **27**: 722-736.
- 1076 Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, Salzberg SL.
1077 2004. Versatile and open software for comparing large genomes. *Genome*
1078 *Biol* **9**.
- 1079 Laetsch DR, Blaxter ML. 2017. BlobTools: Interrogation of genome assemblies.
1080 *F1000Research* **6**: 1287.
- 1081 Lei L, Goltsman E, Goodstein D, Wu GA, Rokhsar DS, Vogel JP. 2021. Plant Pan-
1082 Genomics Comes of Age. *Annu Rev Plant Biol* **72**: 411-435.
- 1083 Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with
1084 BWA-MEM. *ArXiv13033997 Q-Bio*. <http://arxiv.org/abs/1303.3997>
1085 (Accessed August 5, 2021).
- 1086 Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences ed. I. Birol.
1087 *Bioinformatics* **34**: 3094-3100.
- 1088 Li N, He Q, Wang J, Wang B, Zhao J, Huang S, Yang T, Tang Y, Yang S, Aisimutuola
1089 P, et al. 2023. Super-pangenome analyses highlight genomic diversity and
1090 structural variation across wild and cultivated tomato species. *Nat Genet*
1091 **55**: 852-860.
- 1092 Lien S, Koop BF, Sandve SR, Miller JR, Kent MP, Nome T, Hvidsten TR, Leong JS,
1093 Minkley DR, Zimin A, et al. 2016. The Atlantic salmon genome provides
1094 insights into rediploidization. *Nature* **533**: 200-205.
- 1095 Lin Y-L, Gokcumen O. 2019. Fine-Scale Characterization of Genomic Structural
1096 Variation in the Human Genome Reveals Adaptive and Biomedically
1097 Relevant Hotspots. *Genome Biol Evol* **11**: 1136-1151.
- 1098 Low YW, Rajaraman S, Tomlin CM, Ahmad JA, Ardi WH, Armstrong K, Athen P,
1099 Berhaman A, Bone RE, Cheek M, et al. 2022. Genomic insights into rapid
1100 speciation within the world's largest tree genus *Syzygium*. *Nat Commun*
1101 **13**: 5031.
- 1102 Lowry DB, Willis JH. 2010. A Widespread Chromosomal Inversion Polymorphism
1103 Contributes to a Major Life-History Transition, Local Adaptation, and
1104 Reproductive Isolation. *PLOS Biol* **8**: e1000500.

- 1105 Luo X, Xu L, Wang Y, Dong J, Chen Y, Tang M, Fan L, Zhu Y, Liu L. 2020. An ultra-
1106 high-density genetic map provides insights into genome synteny,
1107 recombination landscape and taproot skin colour in radish (*Raphanus*
1108 *sativus* L.). *Plant Biotechnol J* **18**: 274–286.
- 1109 Lynch M, Bobay L-M, Catania F, Gout J-F, Rho M. 2011. The Repatterning of
1110 Eukaryotic Genomes by Random Genetic Drift. *Annu Rev Genomics Hum*
1111 *Genet* **12**: 347–366.
- 1112 Lynch M, Conery JS. 2000. The Evolutionary Fate and Consequences of Duplicate
1113 Genes. *Science* **290**: 1151–1155.
- 1114 Manni M, Berkeley MR, Seppey M, Simão FA, Zdobnov EM. 2021. BUSCO Update:
1115 Novel and Streamlined Workflows along with Broader and Deeper
1116 Phylogenetic Coverage for Scoring of Eukaryotic, Prokaryotic, and Viral
1117 Genomes ed. J. Kelley. *Mol Biol Evol* **38**: 4647–4654.
- 1118 Marques DA, Meier JI, Seehausen O. 2019. A Combinatorial View on Speciation
1119 and Adaptive Radiation. *Trends Ecol Evol* **34**: 531–544.
- 1120 Martin G, Baurens F-C, Hervouet C, Salmon F, Delos J-M, Labadie K, Perdereau A,
1121 Mournet P, Blois L, Dupouy M, et al. 2020. Chromosome reciprocal
1122 translocations have accompanied subspecies evolution in bananas. *Plant J*
1123 **104**: 1698–1711.
- 1124 Mayjonade B, Gouzy J, Donnadieu C, Pouilly N, Marande W, Callot C, Langlade N,
1125 Muñoz S. 2016. Extraction of high-molecular-weight genomic DNA for long-
1126 read sequencing of single molecules. *BioTechniques* **61**: 203–205.
- 1127 Meier JI, Marques DA, Mwaiko S, Wagner CE, Excoffier L, Seehausen O. 2017.
1128 Ancient hybridization fuels rapid cichlid fish adaptive radiations. *Nat*
1129 *Commun* **8**: 14363.
- 1130 Mérot C, Oomen RA, Tigano A, Wellenreuther M. 2020. A Roadmap for
1131 Understanding the Evolutionary Significance of Structural Genomic
1132 Variation. *Trends Ecol Evol* **35**: 561–572.
- 1133 Meudt HM, Albach DC, Tanentzap AJ, Igea J, Newmarch SC, Brandt AJ, Lee WG,
1134 Tate JA. 2021. Polyploidy on Islands: Its Emergence and Importance for
1135 Diversification. *Front Plant Sci* **12**.
1136 <https://www.frontiersin.org/articles/10.3389/fpls.2021.637214> (Accessed
1137 October 12, 2022).
- 1138 Murray KD, Janes JK, Jones A, Bothwell HM, Andrew RL, Borevitz JO. 2019.
1139 Landscape drivers of genomic diversity and divergence in woodland
1140 *Eucalyptus*. *Mol Ecol* **28**: 5232–5247.
- 1141 Myburg AA, Grattapaglia D, Tuskan GA, Hellsten U, Hayes RD, Grimwood J,
1142 Jenkins J, Lindquist E, Tice H, Bauer D, et al. 2014. The genome of
1143 *Eucalyptus grandis*. *Nature* **510**: 356–362.
- 1144 Naseeb S, Ames RM, Delneri D, Lovell SC. 2017. Rapid functional and
1145 evolutionary changes follow gene duplication in yeast. *Proc R Soc B Biol*
1146 *Sci* **284**: 20171393.
- 1147 Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: A Fast and
1148 Effective Stochastic Algorithm for Estimating Maximum-Likelihood

- 1149 Phylogenies. *Mol Biol Evol* **32**: 268–274.
- 1150 Nicolle D. 2022. Classification of the eucalypts (Angophora, Corymbia and
1151 Eucalyptus) Version 6. [http://www.dn.com.au/Classification-Of-The-](http://www.dn.com.au/Classification-Of-The-Eucalypts.pdf)
1152 [Eucalypts.pdf](http://www.dn.com.au/Classification-Of-The-Eucalypts.pdf).
- 1153 Ortiz-Barrientos D, Engelstädter J, Rieseberg LH. 2016. Recombination Rate
1154 Evolution and the Origin of Species. *Trends Ecol Evol* **31**: 226–236.
- 1155 Ou S, Chen J, Jiang N. 2018. Assessing genome assembly quality using the LTR
1156 Assembly Index (LAI). *Nucleic Acids Res*.
1157 [https://academic.oup.com/nar/advance-article/doi/10.1093/nar/gky730/50](https://academic.oup.com/nar/advance-article/doi/10.1093/nar/gky730/5068908)
1158 [68908](https://academic.oup.com/nar/advance-article/doi/10.1093/nar/gky730/5068908) (Accessed August 5, 2021).
- 1159 Ou S, Su W, Liao Y, Chougule K, Agda JRA, Hellinga AJ, Lugo CSB, Elliott TA, Ware
1160 D, Peterson T, et al. 2019. Benchmarking transposable element annotation
1161 methods for creation of a streamlined, comprehensive pipeline. *Genome*
1162 *Biol* **20**: 275.
- 1163 Ouadi S, Sierro N, Goepfert S, Bovet L, Glauser G, Vallat A, Peitsch MC, Kessler F,
1164 Ivanov NV. 2022. The clove (*Syzygium aromaticum*) genome provides
1165 insights into the eugenol biosynthesis pathway. *Commun Biol* **5**: 1–13.
- 1166 Oudelaar AM, Higgs DR. 2021. The relationship between genome structure and
1167 function. *Nat Rev Genet* **22**: 154–168.
- 1168 Passarge E, Horsthemke B, Farber RA. 1999. Incorrect use of the term synteny.
1169 *Nat Genet* **23**: 387–387.
- 1170 Pfeilsticker TR, Jones RC, Steane DA, Vaillancourt RE, Potts BM. 2023. Molecular
1171 insights into the dynamics of species invasion by hybridisation in
1172 Tasmanian eucalypts. *Mol Ecol* **32**: 2913–2929.
- 1173 Potts BM, Wiltshire RJE. 1997. Eucalypt genetics and genecology. In (eds. J.
1174 Williams and J. Woinarski), pp. 56–91, Cambridge University Press
1175 <https://eprints.utas.edu.au/7460/> (Accessed November 30, 2021).
- 1176 Rambaut A. FigTree. <http://tree.bio.ed.ac.uk/software/figtree/> (Accessed
1177 November 28, 2022).
- 1178 Ranwez V, Douzery EJP, Cambon C, Chantret N, Delsuc F. 2018. MACSE v2:
1179 Toolkit for the Alignment of Coding Sequences Accounting for Frameshifts
1180 and Stop Codons. *Mol Biol Evol* **35**: 2582–2584.
- 1181 Rautiainen M, Nurk S, Walenz BP, Logsdon GA, Porubsky D, Rhie A, Eichler EE,
1182 Phillippy AM, Koren S. 2023. Telomere-to-telomere assembly of diploid
1183 chromosomes with Verkko. *Nat Biotechnol* **41**: 1474–1482.
- 1184 Razifard H, Ramos A, Della Valle AL, Bodary C, Goetz E, Manser EJ, Li X, Zhang L,
1185 Visa S, Tieman D, et al. 2020. Genomic Evidence for Complex
1186 Domestication History of the Cultivated Tomato in Latin America. *Mol Biol*
1187 *Evol* **37**: 1118–1132.
- 1188 Rellstab C, Gugerli F, Eckert AJ, Hancock AM, Holderegger R. 2015. A practical
1189 guide to environmental association analysis in landscape genomics. *Mol*
1190 *Ecol* **24**: 4348–4370.
- 1191 Ribeiro T, Barrela RM, Bergès H, Marques C, Loureiro J, Morais-Cecílio L, Paiva

- 1192 JAP. 2016. Advancing Eucalyptus Genomics: Cytogenomics Reveals
1193 Conservation of Eucalyptus Genomes. *Front Plant Sci* **7**.
1194 <http://journal.frontiersin.org/Article/10.3389/fpls.2016.00510/abstract>
1195 (Accessed November 30, 2021).
- 1196 Rieseberg LH. 2001. Chromosomal rearrangements and speciation. **8**.
- 1197 Roach MJ, Schmidt SA, Borneman AR. 2018. Purge Haplotigs: allelic contig
1198 reassignment for third-gen diploid genome assemblies. *BMC*
1199 *Bioinformatics* **19**: 460.
- 1200 Robberecht C, Voet T, Esteki MZ, Nowakowska BA, Vermeesch JR. 2013.
1201 Nonallelic homologous recombination between retrotransposable elements
1202 is a driver of de novo unbalanced translocations. *Genome Res* **23**: 411-
1203 418.
- 1204 Ruggieri AA, Livraghi L, Lewis JJ, Evans E, Cicconardi F, Hebberecht L, Ortiz-Ruiz
1205 Y, Montgomery SH, Ghezzi A, Rodriguez-Martinez JA, et al. 2022. A
1206 butterfly pan-genome reveals that a large amount of structural variation
1207 underlies the evolution of chromatin accessibility. *Genome Res* **32**: 1862-
1208 1875.
- 1209 Sayers EW, Beck J, Bolton EE, Bourexis D, Brister JR, Canese K, Comeau DC, Funk
1210 K, Kim S, Klimke W, et al. 2021. Database resources of the National Center
1211 for Biotechnology Information. *Nucleic Acids Res* **49**: D10-D17.
- 1212 Scannell DR, Wolfe KH. 2008. A burst of protein sequence evolution and a
1213 prolonged period of asymmetric evolution follow gene duplication in yeast.
1214 *Genome Res* **18**: 137-147.
- 1215 Schumer M, Cui R, Rosenthal GG, Andolfatto P. 2015. Reproductive Isolation of
1216 Hybrid Populations Driven by Genetic Incompatibilities. *PLoS Genet* **11**:
1217 e1005041.
- 1218 Simakov O, Marlétaz F, Yue J-X, O'Connell B, Jenkins J, Brandt A, Calef R, Tung C-
1219 H, Huang T-K, Schmutz J, et al. 2020. Deeply conserved synteny resolves
1220 early events in vertebrate evolution. *Nat Ecol Evol* **4**: 820-830.
- 1221 Smit A, Hubley R, Green P. 2020. RepeatMasker Open-4.0.
1222 <<http://www.repeatmasker.org>> (Accessed February 11, 2020).
- 1223 Sterck L, Rombauts S, Vandepoele K, Rouze P, Vandeppeer Y. 2007. How many
1224 genes are there in plants (... and why are they there)? *Curr Opin Plant Biol*
1225 **10**: 199-203.
- 1226 Supple MA, Bragg JG, Broadhurst LM, Nicotra AB, Byrne M, Andrew RL, Widdup A,
1227 Aitken NC, Borevitz JO. 2018. Landscape genomic prediction for
1228 restoration of a *Eucalyptus* foundation species under climate change ed.
1229 D.J. Kliebenstein. *eLife* **7**: e31835.
- 1230 Tang D, Jia Y, Zhang J, Li H, Cheng L, Wang P, Bao Z, Liu Z, Feng S, Zhu X, et al.
1231 2022. Genome evolution and diversity of wild and cultivated potatoes.
1232 *Nature* **606**: 535-541.
- 1233 Thornhill AH, Crisp MD, Külheim C, Lam KE, Nelson LA, Yeates DK, Miller JT. 2019.
1234 A dated molecular perspective of eucalypt taxonomy, evolution and
1235 diversification. *Aust Syst Bot* **32**: 29-48.

- 1236 Torkamaneh D, Lemay M-A, Belzile F. 2021. The pan-genome of the cultivated
1237 soybean (PanSoy) reveals an extraordinarily conserved gene content.
1238 *Plant Biotechnol J* **19**: 1852–1862.
- 1239 Twyford AD, Friedman J. 2015. Adaptive divergence in the monkey flower
1240 *Mimulus guttatus* is maintained by a chromosomal inversion. *Evolution* **69**:
1241 1476–1486.
- 1242 Van de Peer Y, Maere S, Meyer A. 2009. The evolutionary significance of ancient
1243 genome duplications. *Nat Rev Genet* **10**: 725–732.
- 1244 Vaser R, Sović I, Nagarajan N, Šikić M. 2017. Fast and accurate de novo genome
1245 assembly from long uncorrected reads. *Genome Res* **27**: 737–746.
- 1246 Voelker J, Shepherd M, Mauleon R, Shepherd M, Mauleon R. 2021. A high-quality
1247 draft genome for *Melaleuca alternifolia* (tea tree): a new platform for
1248 evolutionary genomics of myrtaceous terpene-rich species. *Gigabyte*
1249 **2021**: 1–15.
- 1250 Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA,
1251 Zeng Q, Wortman J, Young SK, et al. 2014. Pilon: An Integrated Tool for
1252 Comprehensive Microbial Variant Detection and Genome Assembly
1253 Improvement ed. J. Wang. *PLoS ONE* **9**: e112963.
- 1254 Wang W, Das A, Kainer D, Schalamun M, Morales-Suarez A, Schwessinger B,
1255 Lanfear R. 2020. The draft nuclear genome assembly of *Eucalyptus*
1256 *pauciflora*: a pipeline for comparing de novo assemblies. *GigaScience* **9**:
1257 giz160.
- 1258 Weisman CM, Murray AW, Eddy SR. 2020. Many, but not all, lineage-specific
1259 genes can be explained by homology detection failure. *PLOS Biol* **18**:
1260 e3000862.
- 1261 Weissensteiner MH, Bunikis I, Catalán A, Francoijs K-J, Knief U, Heim W, Peona V,
1262 Pophaly SD, Sedlazeck FJ, Suh A, et al. 2020. Discovery and population
1263 genomics of structural variation in a songbird genus. *Nat Commun* **11**:
1264 3403.
- 1265 Wellenreuther M, Bernatchez L. 2018. Eco-Evolutionary Genomics of
1266 Chromosomal Inversions. *Trends Ecol Evol* **33**: 427–440.
- 1267 Wu B, Cox MP. 2019. Greater genetic and regulatory plasticity of retained
1268 duplicates in *Epichloë* endophytic fungi. *Mol Ecol* **28**: 5103–5114.
- 1269 Xu L, Chen H, Hu X, Zhang R, Zhang Z, Luo ZW. 2006. Average Gene Length Is
1270 Highly Conserved in Prokaryotes and Eukaryotes and Diverges Only
1271 Between the Two Kingdoms. *Mol Biol Evol* **23**: 1107–1108.
- 1272 Yeaman S. 2013. Genomic rearrangements and the evolution of clusters of
1273 locally adaptive loci. *Proc Natl Acad Sci* **110**: E1743–E1751.
- 1274 Zhang C, Rabiee M, Sayyari E, Mirarab S. 2018. ASTRAL-III: polynomial time
1275 species tree reconstruction from partially resolved gene trees. *BMC*
1276 *Bioinformatics* **19**: 153.



Plant genome evolution in the genus *Eucalyptus* driven by structural rearrangements that promote sequence divergence

Scott Ferguson, Ashley Jones, Kevin Murray, et al.

Genome Res. published online April 8, 2024

Access the most recent version at doi:[10.1101/gr.277999.123](https://doi.org/10.1101/gr.277999.123)

Supplemental Material <http://genome.cshlp.org/content/suppl/2024/05/03/gr.277999.123.DC1>

P<P Published online April 8, 2024 in advance of the print journal.

Accepted Manuscript Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.

Creative Commons License This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see <https://genome.cshlp.org/site/misc/terms.xhtml>). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at <http://creativecommons.org/licenses/by-nc/4.0/>.

Email Alerting Service Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).



The NEW Vortex Mixer

USC
SCIENTIFIC
SINCE 1973

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>
