

1 **The heterozygous pineapple genome demonstrates the importance of haplotype-resolved**  
2 **plant genomes**

3 Jingping Fang<sup>1,2†</sup>, Patrick J. Mason<sup>2,3†</sup>, Garth Sanewski<sup>4</sup>, Matthew Webb<sup>4</sup>, Linwei Zhou<sup>1</sup>,  
4 Jinbin Wang<sup>1</sup>, Rhys G. R. Copeland<sup>5</sup>, Ying Zhuang<sup>1</sup>, Annapurna Chitikineni<sup>5</sup>, Vanika Garg<sup>5</sup>,  
5 Natalie Dillon<sup>4</sup>, Parwinder Kaur<sup>6</sup>, Rajeev K. Varshney<sup>5</sup>, Robert J. Henry<sup>2,3,7</sup>

6 **Author Details:**

7 <sup>1</sup>College of Life Science, Fujian Normal University, Fuzhou 350117, China

8 <sup>2</sup>Queensland Alliance of Agriculture and Food Innovation, The University of Queensland, St  
9 Lucia, QLD, Australia

10 <sup>3</sup>ARC Centre of Excellence for Plant Success in Nature and Agriculture, The University of  
11 Queensland, St Lucia QLD, 4072 Australia

12 <sup>4</sup>Queensland Department of Primary Industries, Brisbane, QLD, 4000, Australia

13 <sup>5</sup>WA State Agricultural Biotechnology Centre, Centre for Crop and Food Innovation, Food  
14 Futures Institute, Murdoch University, Murdoch, WA, 6150, Australia

15 <sup>6</sup>School of Agriculture and Environment, University of Western Australia, 35 Stirling Way,  
16 Crawley WA, 6009, Australia

17 <sup>7</sup>VinUni Big Data Research Institute, VinUniversity, Hanoi, Vietnam

18

19 **Corresponding Author**

20 Robert J. Henry

21 Email: [robert.henry@uq.edu.au](mailto:robert.henry@uq.edu.au)

22 Phone: +61 7 3346 0552

23 Address: Queensland Alliance for Agriculture and Food Innovation (QAAFI), Level 2,  
24 Queensland Biosciences Precinct [#80], the University of Queensland, St Lucia QLD,  
25 Australia, 4072.

26

27 † These authors contributed equally

28

29

30

31

32 © The Author(s) 2026. Published by Oxford University Press on behalf of the Nanjing  
33 Agricultural University. This is an Open Access article distributed under the terms of the  
34 Creative Commons Attribution License <https://creativecommons.org/licenses/by/4.0/>, which  
35 permits unrestricted reuse, distribution, and reproduction in any medium, provided the  
36 original work is properly cited.

37 **Abstract:**

38 Collapsed (haploid) genome assemblies omit large portions of genetic information, especially  
39 in heterozygous, clonally propagated crops such as pineapple. Here, we assembled a telomere-  
40 to-telomere, haplotype-resolved genome for a key pre-Colombian cultivar of pineapple  
41 (*Ananas comosus*) ‘Smooth Cayenne’ (F180) using PacBio Hi-Fi and Hi-C data. The two 25-  
42 chromosome haplotypes span 858 Mb ( $N50 \approx 16.8$  Mb) and are >99% complete, each resolving  
43 all centromeres and 22 of 25 telomeres. Comparison of the phased chromosomes reveals 1.5  
44 million single nucleotide polymorphisms (SNPs) and 1,953 large structural variants (74  
45 inversions, 750 translocations, and 1,129 segmental duplications). This assembly reveals that  
46 inversions have profoundly impacted the ‘Smooth Cayenne’ genome, reshaping ~3-4% of the  
47 total sequence. Structural context dictates the genetic impact of these large inversions, as shown  
48 in recombination landscape analysis of 374 F<sub>1</sub> seedlings, wherein a 1.3 Mb paracentric  
49 inversion on chromosome 20 forms a strict recombination coldspot, whereas a 6 Mb pericentric  
50 inversion on chromosome 24 still permits gene flow likely via short double crossovers, albeit  
51 at lower rates than the rest of the chromosome. Re-anchoring the 11,879 DArTseq markers  
52 from the F<sub>1</sub> seedlings to the phased reference assembly, removes the dense network of spurious  
53 inter-chromosomal linkage seen in the collapsed F153 ‘Smooth Cayenne’ genome, likely  
54 providing markedly cleaner baselines for genome-wide association studies (GWAS) and  
55 genomic prediction. These results establish the new F180 assembly as a very high-quality  
56 reference, illustrate how undetected inversions can silently constrain genetic gain, and  
57 demonstrate the broader value of phased genomes for dissecting heterozygosity, structural  
58 variation and meiotic behaviour in perennial crops.

59 **Keywords**

60 *Ananas comosus*; haplotype-resolved; paracentric inversion; pericentric inversion;  
61 recombination landscape, structural variation

62

63

64

65

66

67

68

## 69 1. Introduction

70 Collapsed assemblies of diploid or polyploid organisms provide only a partial representation  
71 of their genomes, as they capture only one version of each heterozygous region [1]. This  
72 limitation is particularly problematic in highly heterozygous genomes, where valuable genetic  
73 variation is lost. Compounding this issue, early genome assemblies based on short-read  
74 sequencing or low-accuracy long reads often produced fragmented, collapsed assemblies,  
75 especially in heterozygous species [2,3]. The availability of highly accurate long-read and high-  
76 throughput chromosome conformation capture (Hi-C) [4], alongside new assemblers [5], has  
77 enabled the sequencing and assembly of very high-quality haplotype-resolved telomere-to-  
78 telomere (T2T) genomes. Genomic fragmentation and the collapsed nature of these earlier  
79 assemblies greatly limit their utility for fundamental genetic studies and breeding programs.  
80 For instance, in blueberry, a phased genome assembly enabled more accurate single nucleotide  
81 polymorphism (SNP) identification and quantitative trait locus (QTL) mapping, resulting in  
82 improved genome-wide association study (GWAS) outcomes and selection of informative  
83 markers for genomic selection [6]. This demonstrates how haplotype-resolved assemblies  
84 outperform collapsed ones by providing more accurate identification of SNPs and QTLs for  
85 breeding and fundamental plant genetics.

86 Unlike SNPs and small InDels, large structural variants (SVs), particularly inversions, can  
87 profoundly impact genomic architecture by forming haploblocks that suppress recombination  
88 and alter the evolutionary and breeding trajectories of linked loci [7-11]. In heterozygotes,  
89 single crossovers (SCOs) within an inversion generate unbalanced chromatids, which in turn  
90 create unviable gametes, suppress recombination and lead to large blocks of alleles within a  
91 chromosomal region being inherited as a block [12]. However, this suppression is not absolute.  
92 Analysis of insect [13,14] and, more recently, plant genomes [8,15] shows that recombination  
93 can occasionally leak through via double crossovers (DCOs), inversion toggling or short gene  
94 conversion tracts, creating a spectrum from “tight” to “leaky” inversions. Differences in  
95 inversion size, position relative to the centromere, local sequence context or age govern this  
96 spectrum, although the contribution of these remains poorly resolved in most crops [16,17].  
97 Resolving both haplotypes in a T2T assembly is therefore a prerequisite to discovering,  
98 genotyping and ultimately exploiting polymorphic inversions in breeding programmes. Recent  
99 phased-assembly work in the clonally propagated perennial fruit, mango (*Mangifera indica*),  
100 demonstrated that both pericentric and paracentric inversions can harbour breeding-relevant

101 loci and show heterogeneous patterns of recombination suppression [11]. However, equivalent  
102 high-resolution analyses are still lacking in many other clonally propagated fruit crops.

103 Several economically important crops, particularly those propagated vegetatively for centuries  
104 or millennia, exhibit elevated heterozygosity [18]. This diversity, preserved through clonal  
105 propagation, may contribute to traits such as hybrid vigour or resilience to environmental  
106 stresses, thereby sustaining certain genotypes despite minimal genetic improvement [19]. Such  
107 heterozygosity can be linked to a variety of factors, including self-incompatibility and  
108 vegetative propagation over thousands of years leading to the accumulation of different somatic  
109 mutations in each haplotype. This high heterozygosity may also originate from wide  
110 hybridisations in their evolutionary past i.e. crosses between genetically distinct lineages [20],  
111 maintained via vegetative reproduction. However, direct evidence for such hybrid origins in  
112 pineapple remains limited due to the ancient nature of its domestication history. Understanding  
113 haplotype-specific variation in these crops is essential for deciphering the genomic basis of  
114 their adaptability and enhancing modern breeding strategies [21].

115 Pineapple (*Ananas comosus var. comosus*), domesticated from its wild progenitor *Ananas*  
116 *comosus var. microstachys* [22], exemplifies this paradigm. Cultivated through a combination  
117 of sexual recombination, self-incompatibility and clonal propagation since its ancient  
118 domestication, pineapple maintains a highly heterozygous genome [23]. As a globally  
119 important tropical fruit, pineapple production is constrained by a narrow genetic base [24],  
120 exacerbated by reliance on a limited number of cultivars, particularly ‘Smooth Cayenne’,  
121 making it vulnerable to pathogens and abiotic stresses [25,26]. In this study, we report a  
122 haplotype-resolved T2T genome assembly of ‘Smooth Cayenne’ (F180 selection), a key pre-  
123 Columbian variety of pineapple. This haplotype-resolved genome assembly represents a  
124 substantial improvement over previous pineapple genome assemblies [23,27] and provides a  
125 foundational reference for understanding heterozygosity, structural variants, and recombination  
126 dynamics in clonally propagated crops.

127

128

129

130

## 131 2. Results

### 132 2.1 *De novo haplotype-resolved genome assembly of Smooth Cayenne*

133 The ‘Smooth Cayenne’ (F180) genome was assembled using PacBio high-fidelity (HiFi) and  
134 Hi-C sequencing data. In total, we obtained 46.33 Gb ( 47×) of HiFi reads with an N50 of 27  
135 kb and 39.86 Gb ( 40×) of Hi-C data (**Figure S1, Table S1**). GenomeScope analysis estimated  
136 a heterozygosity of 1.57% based on a  $k$ -mer of 21 (**Figure S2**). To achieve the T2T haplotype-  
137 resolved genome assembly, we adopted a haplotype-aware pipeline integrating HiFi and Hi-C  
138 reads to assemble the phased genome (**Figure S3**). Chromosomes of both haplotypes were  
139 manually ordered by size with short arms forward, and their correspondence to F153 [23] is  
140 provided in **Table S2**. After iterative manual refinement, the Hi-C contact maps presented high  
141 consistency across all chromosomes, with strong diagonal signals and clear T2T interactions,  
142 supporting the accuracy of chromosome ordering and orientation (**Figure 1a, b**). Following  
143 removal of organelle-derived and redundant contigs, a nearly complete phased reference  
144 genome was generated, consisting of 50 pseudo-chromosomes (25 per haplotype) totalling  
145 858.47 Mb (433.27 Mb and 425.20 Mb, respectively) (**Figure 1c, Table 1**). The chromosome  
146 lengths for haplotype A (HA) ranged from 12.49 Mb to 21.07 Mb (N50 = 16.52 Mb), while the  
147 lengths for haplotype B (HB) ranged from 12.09 Mb to 20.48 Mb (N50 = 16.86 Mb) (**Table**  
148 **S3**). Telomeric repeat motifs were detected at both ends of 22 chromosomes in each haplotype,  
149 and all 50 centromeres were identified (**Figure S4-S5, Table S3-S4**). A completely collapsed  
150 genome, where two haplotypes were merged and parental alleles randomly switched within a  
151 contig, could not be assembled due to the formation of chimeric chromosomes (**Figure S6**). Of  
152 the 25 chromosomes, eight were assembled as chimeric, retaining large haplotype-specific  
153 blocks, rather than creating a coherent collapsed assembly. This issue stems from significant  
154 structural discrepancies and low sequence collinearity between the haplotypes, which led the  
155 software to incorrectly merge the haplotypes instead of generating a collapsed genome.

156 Different approaches were used to evaluate the quality of the assembly and phasing of HA and  
157 HB. The overall mapping rates for HiFi, stLFR, and RNA-seq reads from various tissues were  
158 99.62%, 97.86%, and 83.11-95.29%, respectively (**Table S5**). We estimated the base-level  
159 accuracy using Merqury based on  $k$ -mer spectra, resulting in a quality value (QV) of 65.21 for  
160 HA and 64.94 for HB, respectively, with a base error rate of  $< 10^{-6}$ . BUSCO assessment  
161 revealed that HA, HB, and the whole genome of F180 contained 99.2%, 99.4%, and 99.5%  
162 complete core orthologous genes, respectively, based on the highly conserved core

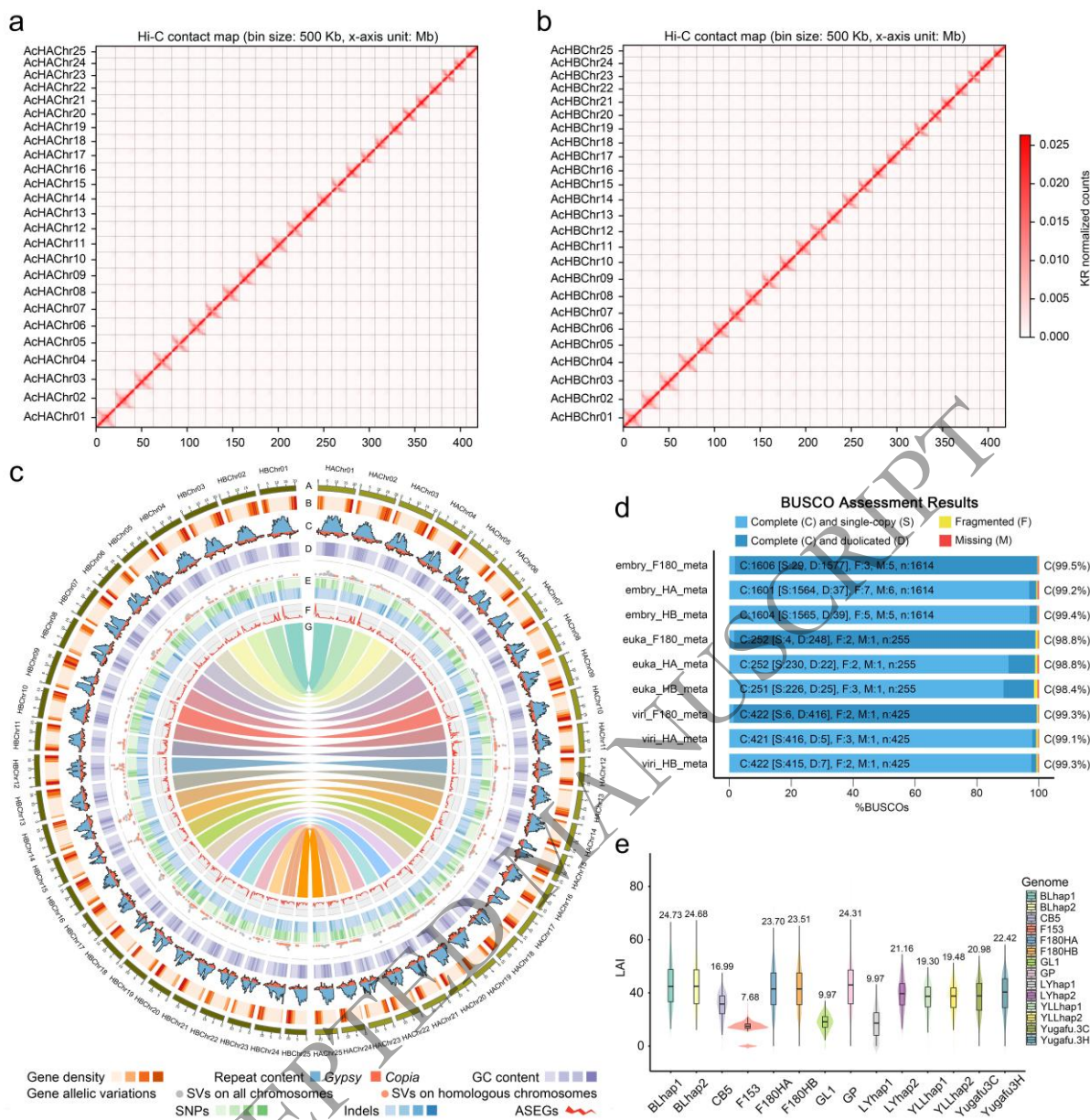
163 embryophyte gene set ( $N = 1,614$ ) and the default gene predictor MetaEuk (**Figure 1d, Figure**  
 164 **S7, Table S6**), ranking F180 among the top three pineapple genomes (**Figure S8**). Particularly,  
 165 the BUSCO value for F180 HB improved from 98.3% to 99.4% following several rounds of  
 166 manual curation. The Long Terminal Repeats (LTR) Assembly Indices (LAI) were evaluated  
 167 at 24.82 for HA and 23.23 for HB, both surpassing the gold quality assembly threshold of “LAI >  
 168 20” [28] and are superior to those of most other pineapple genomes (**Figure 1e**). Together,  
 169 these quality assessments validate that we have assembled a nearly complete, T2T haplotype-  
 170 resolved genome for F180 with high accuracy, contiguity, and fully resolved centromeres.

171 **Table 1. Summary statistics of genome assembly and annotation for *Ananas comosus* var.**  
 172 ***comosus* cultivar ‘Smooth Cayenne’ F180.** Assembly statistics are reported for haplotypes A (HA)  
 173 and B (HB) separately, and for the combined diploid genome (50 chromosomes). Values include  
 174 cumulative assembly size, contiguity (N50), GC content, scaffold number, BUSCO completeness  
 175 (embryophyta\_odb10,  $n = 1,614$ ), LTR Assembly Index (LAI), quality value (QV), gap counts, and  
 176 detection of telomeres and centromeres. Annotation statistics include the number of protein-coding  
 177 genes, gene length metrics, and BUSCO completeness for predicted proteins (embryophyta\_odb10,  $n =$   
 178  $1,614$ ). Detailed mapping statistics, transcriptome alignment rates, and functional annotation  
 179 breakdowns are provided in Supplementary Table S5.

Category	HA (25 chr)	HB (25 chr)	Whole genome (50 chr)
<b>Assembly</b>			
Cumulative size (bp)	414.01 Mb	412.80 Mb	826.81 Mb
N50 (bp)	16.52 Mb	16.86 Mb	16.83 Mb
GC content (%)	39.15	39.06	39.10
Hi-C scaffolds	25	25	50
Complete BUSCO (%) (embryophyta_odb10, $n = 1,614$ )	99.20	99.40	99.50
LAI	24.82	23.23	–
QV	65.21	64.94	65.07
Gaps (count)	65	77	142
Telomeres detected	47/50	47/50	94/100
Centromeres detected	25/25	25/25	50/50
<b>Annotation</b>			
Protein-coding genes	20,280	20,248	40,528
Mean gene length (bp)	4,593	4,598	–
Mean CDS length (bp)	1,341	1,349	–
Average exon number per gene	6.0	6.0	–
Complete BUSCO for proteins (%) (embryophyta_odb10, $n = 1,614$ )	99.10	99.00	–

180

181



182

183 **Figure 1 An overview of the 'Smooth Cayenne' F180 reference genome. a-b.** Hi-C  
 184 interaction heatmaps represent contact matrices of 25 chromosomes ( $2n = 50$ ) in F180; **a.** HA  
 185 genome and **b.** HB genome at 500 kb resolution; **c.** The circos plot of the F180 genome at 500  
 186 kb intervals. The tracks represent the following elements (from outer to inner): (A)  
 187 chromosomes (in Mb), (B) gene density, (C) repeat content: LTR/*Gypsy* and *Copia*, (D) GC  
 188 content, (E) gene allelic variations: SV on all and homologous chromosomes, SNPs, and Indels,  
 189 (F) allele-specific expression genes (ASEGs); **d.** The completeness of the F180 genome,  
 190 including HA, HB and the whole genome assessed in BUSCO (v5.8.0) with MetaEuk tool,  
 191 based on the latest embryophyta\_odb10, eukaryota\_odb10, and viridiplantae\_odb10 database;  
 192 **e.** The violin plot of the longest terminal repeat (LTR) index (LAI) among nine pineapple  
 193 genome assemblies at the chromosome level. The horizontal axis represents different  
 194 accessions, while the number in each box indicates the LAI value.

## 195 2.2 Repetitive sequence and gene annotation

196 Genome-wide repeat annotation revealed that interspersed elements dominate the F180  
197 genome (**Table S7**), comprising 266.49 Mb (61.51%) in HA and 264.46 Mb (62.20%) in HB,  
198 with a large proportion remaining unclassified (39.38% and 40.54%, respectively). Long  
199 terminal repeat retrotransposons (LTR-RTs), the major category of retroelements (Class I),  
200 were the prevalent interspersed repeats, contributing 20.27% of HA and 19.25% of HB. Within  
201 this category, Gypsy/DIRS1 were the most predominant type (15.43-16.18%), followed by  
202 Ty1/Copia (3.66-3.71%) (**Figure 1c, Table S7**). In contrast, DNA transposons (Class II)  
203 represented only a minor fraction (1.05% in HA and 1.36% in HB). Aside from interspersed  
204 repeats, small RNA repeats, simple repeats, and low-complexity repeats constituted only 0.28-  
205 1.89% of the genome. Overall, a cumulative length of 281.50 Mb (64.97%) in HA and 277.76  
206 Mb (65.32%) in HB were masked by RepeatMasker prior to gene annotation.

207 We predicted 20,280 and 20,248 high-confidence protein-coding genes (PCGs), corresponding  
208 to 24,280 and 24,283 transcripts or coding sequences (CDS) for HA and HB, respectively  
209 (**Table S5**). The PCGs spanned 93.16 Mb for HA and 93.12 Mb for HB, with average lengths  
210 of 4,593 bp and 4,598 bp for PCGs, and 1,341 bp and 1,349 bp for CDSs. Each gene harboured  
211 an average of 6.0 exons, with exon lengths averaging 222 bp for HA and 223 bp for HB.  
212 Functional annotation against the Nr, Gene Ontology (GO), EggNOG, KEGG, and SwissProt  
213 databases assigned 97.86% of PCGs to at least one database, with ~26% annotated by all five  
214 databases (**Table S5, Figure S9**). The complete BUSCO values for PCGs were 99.1%, 99.0%  
215 and 99.5% for HA, HB and the whole genome, respectively, indicating high gene set  
216 completeness. The F180 HA and HB assemblies showed strong gene-level collinearity with  
217 F153 and the gap-free pineapple reference genome GP [27] (**Figure S10**).

218 A total of 2,919 and 2,962 ribosomal RNAs (rRNAs) were identified for HA and HB, with  
219 combined lengths of 3.49 Mb and 2.51 Mb, respectively (**Table S8**). At the chromosome level,  
220 5S rRNA was the most abundant, distributed across 13 chromosomes with a large cluster on  
221 chromosome 18 (Chr18), while 28S and 5.8S rRNAs were the least common, present only on  
222 Chr25 in HA and on Chr13 and Chr25 in HB. Additionally, 1,032 and 870 transfer RNAs  
223 (tRNAs) were predicted in HA and HB, with total lengths of 75.82 kb and 63.71 kb,  
224 respectively (**Table S9**).

### 225 2.3 Structural variations and three large inversions

226 Synteny analysis revealed 489 syntenic regions covering ~310 Mb of each haplotype, with a  
227 total of 8,851 SVs detected between the two haplotypes, including 74 inversions (0.8%), 750  
228 translocations (8.5%), and 1,129 duplications (12.8%), and 6,898 presence/absence variations  
229 (PAVs; 77.9%) (**Figure 2a** and **Table S10**). The PAVs comprised 3,391 large insertions (1.64  
230 Mb) and 3,507 large deletions (2.06 Mb), which represent a major source of sequence  
231 divergence between haplotypes. Most syntenic blocks exhibited high collinearity between  
232 haplotypes, though 1,435 HA and 2,092 HB regions were unaligned. Allelic sequence  
233 comparisons between haplotypes revealed a multitude of sequence variations, comprising  
234 1,538,444 SNPs, 112,393 small insertions, 113,639 small deletions, 93 copy gains, 106 copy  
235 losses, 18,083 highly diverged sequences, and 20 tandem repeats.

236 Regarding inversions between the two haplotypes, the combined length of all 74 inversions  
237 was 12.27 Mb (2.96%) in HA and 15.66 Mb (3.79%) in HB, with average lengths of 165.8 kb  
238 and 211.7 kb, respectively. The inversion sizes varied among chromosomes, ranging from 219  
239 bp to 6 Mb in HA and 219 bp to 7 Mb in HB (**Tables S11-S12, Figure S11**). The three largest  
240 inversions (> 1 Mb), Inv560, Inv522, and Inv553, were located on Chr24, Chr11, and Chr20,  
241 respectively, with size ranges of 6.05-7.71 Mb, 1.51-2.48 Mb, and 1.30-1.32 Mb (HA vs HB).  
242 The Hi-C contact matrices confirmed the presence of genuine inversions between these  
243 homologous chromosomes (**Figure 2b, Figure S12-S13**). To assess how these structural  
244 rearrangements affected gene organization, we further examined synteny at the protein-coding  
245 gene level. Much higher syntenic relationships among genes between the two haplotypes were  
246 observed compared to the nucleotide level (**Figure 2c, Figure S14**). Genes on most  
247 chromosomes except Chr11, Chr20, and Chr24 were largely conserved (**Figure 2d, Table S12-**  
248 **S13**). Although the Chr24 inversion is the largest structurally, it contained fewer genes (97 in  
249 HA and 98 in HB) than the Chr20 inversion (174 in HA and 180 in HB), which spanned only  
250 one-fifth of its physical length. Likewise, the Chr11 inversion, despite being the second largest  
251 by nucleotide length, encompassed only 13-14 genes and was therefore barely visible at the  
252 gene-collinearity scale.

253 To further validate the authenticity of the largest inversion (Inv560) on Chr24, the entire dataset  
254 of F180 PacBio long reads was aligned to the F180 HB assembly (**Figure 2e**). Two breakpoints  
255 (at 303,529 bp and 8,009,074 bp) of Inv560 in HB were identified. PacBio reads from the HB  
256 haplotype fully spanned both breakpoints, covering the entire inversion and its flanking

257 sequences. In contrast, no HA reads spanned both breakpoints, with alignments exclusively  
258 detected up to 928 bp upstream of the left border and extending to 1,837 bp downstream of the  
259 right border in the HB reference. These border regions were annotated as nonfunctional (gene  
260 desert) areas, indicating no gene disruption by the breakpoints on Chr24 (**Table S13**). Notably,  
261 Inv560 encapsulated the entire centromere adjacent to a breakpoint in both haplotypes,  
262 rendering Chr24 in HB the most acrocentric chromosome (**Figure S5**). Similarly, the presence  
263 of Inv553 on Chr20 was confirmed by PacBio long-read alignments (**Figure S15**).

264 Representatives of all available pineapple genomes (**Table S14**) were aligned against HA  
265 chromosomes 1 to 25 (**Figure S16**). Dot plots of these alignments revealed that the 1.3 Mb  
266 inversion (Inv553) in Chr20 of HB was unique compared to the other genomes. However,  
267 haplotype 1 (H1) of the domesticated Queen variety ‘Ba Li’ (‘BL’) potentially contains this  
268 inversion nested inside another inversion (**Figure S17**). Obvious counterparts of the 6 Mb  
269 inversion (Inv560) of Chr24 were not detected in any other pineapple genomes (**Figure S16**).  
270 However, several rearrangements and smaller inversions were detected in its vicinity (**Figure**  
271 **S17**). Notably, a large deletion in ‘BL’ H1 included the centromere (**Figure S18**), thereby  
272 indicating incompleteness of its assembly. A counterpart of the second largest inversion in  
273 Chr11 HA was only found in ‘BL’ H2 (**Figure S16**).

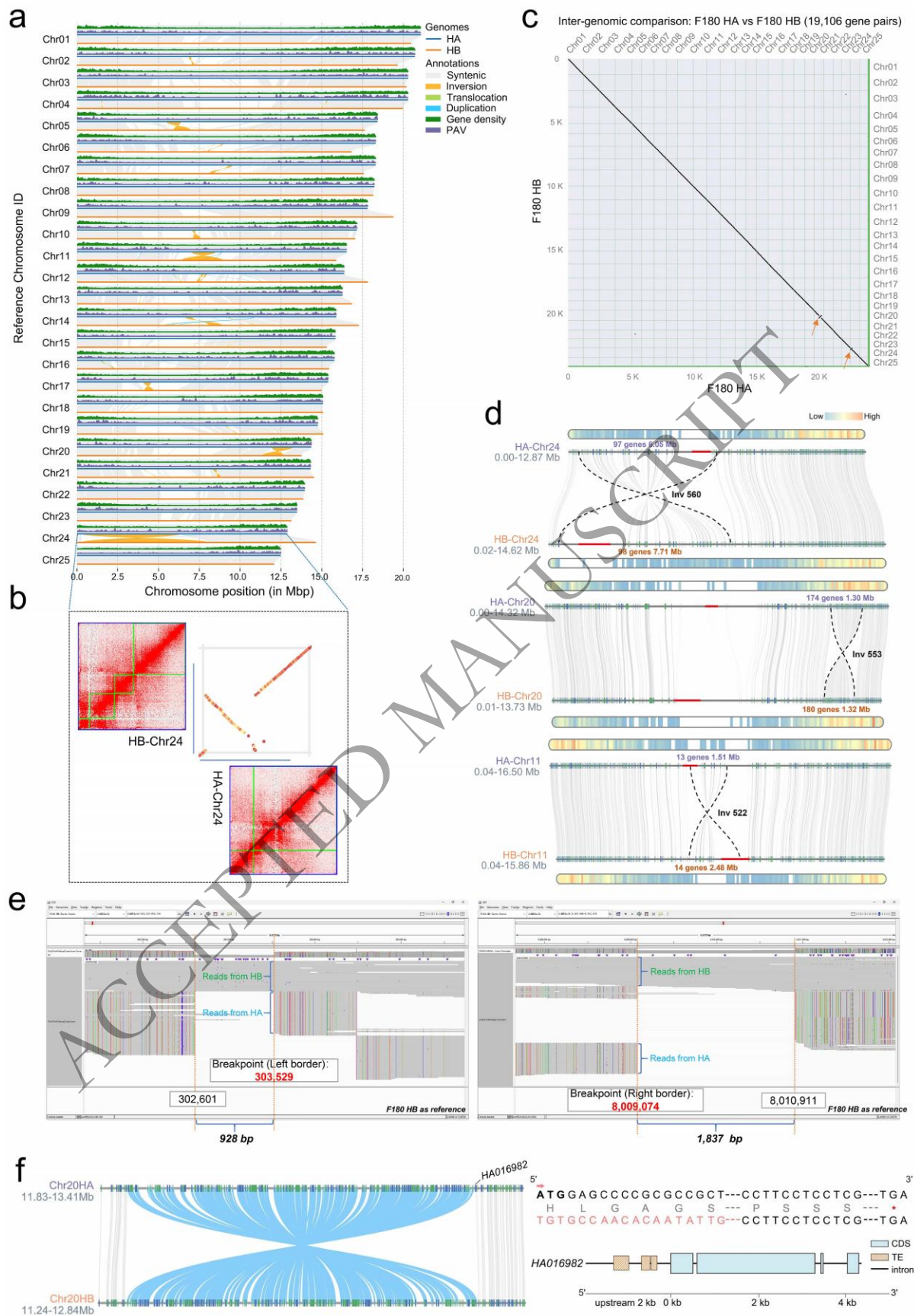
274 Compared to the F180 Chr17 haplotypes, ‘BL’ appeared to have a large chromosome arm  
275 deletion and duplications in haplotypes H1 and H2, respectively (**Figure S16**). A large  
276 chromosomal arm inversion was also obvious in both haplotypes of the phased ‘MD-2’ genome  
277 compared to the F180 Chr25 haplotypes, among other notable differences that similarly  
278 occurred in the original collapsed F153 genome. Levels of collinearity between F180 and the  
279 haplotype-resolved chromosomes of ‘Yugafu’ were generally acceptable, with exception of  
280 Chr5 and Chr25 and several heterozygous inversions. The collinearity of F180 with the  
281 haplotype-resolved genomes of *erectifolius* (‘LY’) and *microstachys* (‘YLL’) varieties was  
282 generally high, with some obvious regions of low homology. In contrast, available *bracteatus*  
283 collapsed assemblies (‘CB5’, ‘GL1’) appeared to possess numerous structural rearrangements  
284 that did not have obvious counterparts in other genomes.

285

286

287

288



289

290

291 **Figure 2 Whole genome syntenic comparisons within ‘Smooth Cayenne’ F180. a.** PAV  
292 landscape and nucleotide-level collinearity between F180 HA and HB. Gray, orange, green,  
293 and blue lines represent synteny, inversion, translocation, and duplication, respectively; The  
294 dark green and purple tracks represent gene density and PAV distribution, calculated using a  
295 50-kb sliding window; **b.** The HiC interaction signal along Chr24 for HA and HB presents high  
296 consistency between neighbouring blocks. Each green box represents a supercontig, while each  
297 royal blue box represents a chromosome; **c.** Synteny dot plot of whole-genome protein-coding  
298 genes between HA and HB generated by MCScanX [29]. Two inversions located at Chr20 and  
299 Chr24 were visible; **d.** Synteny of protein-coding genes of Chr11, Chr20 and Chr24 between  
300 two subgenomes. The blue and green rectangle frames indicate genes transcribed in clockwise  
301 and counterclockwise directions, while the grey bars denote the intergenic regions. The  
302 centromere positions are marked by red bars on the axes; **e.** PacBio read mapping validates the  
303 largest insertion (6.0~7.7 Mb) on Chr24. Reads from HB, but not HA, span the two breakpoints  
304 (at 303,529 bp and 8,009,074 bp) of this Chr24 inversion on the F180 HB reference. Notably,  
305 no HA reads aligned to the 928 bp upstream of the left border or 1,837 bp downstream of the  
306 right border in the HB reference, both of which were annotated as nonfunctional regions; **f.** The  
307 hemizygous gene *HA016982.1* spans the inversion breakpoint on Chr20. The gene retains an  
308 intact coding sequence in HA, while its HB allele is truncated by an ATG deletion. Its 2-kb  
309 promoter region harbours three TE insertions.

310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323  
324  
325  
326  
327

ACCEPTED MANUSCRIPT

#### 328 2.4 Allelic gene variations and inversion-specific comparison

329 In total, 19,058 allelic gene pairs (19,153 one-to-one gene pairs) were identified between HA  
330 and HB, comprising 19,084 genes from HA and 19,065 from HB (**Tables S15-S16, Figure**  
331 **S19**). Most homologous alleles exhibited a one-to-one correspondence, while a few displayed  
332 one-to-multiple or multiple-to-multiple relationships. Additionally, 1,196 and 1,183 genes were  
333 identified as unique to HA and HB, respectively. The  $Ka/Ks$  ratios for allelic gene pairs were  
334 generally less than one (95.79%; mean = 0.32), whereas positively selected genes ( $Ka/Ks > 1$ )  
335 were significantly enriched in defence-related biological processes, particularly responses to  
336 fungi and other organisms (**Figure S20**). RNA-seq analysis of HA-HB allelic genes revealed  
337 clear tissue-specific expression, characterized by strong within-tissue clustering (**Figure 3a, b**).  
338 Overall, these genes showed higher expression in roots and flowers/fruits than in leaves, and  
339 among the three distinct leaf tissues, the meristem exhibited the highest average expression  
340 (**Figure 3c**).

341 We examined allelic gene presence/absence variations within and around inversions, as well as  
342 any resulting positional shifts. Overall, a total of 330 and 337 genes were located within HA  
343 and HB inversions, respectively (**Table S13**). Specifically, 303 allelic gene pairs were present  
344 within inversions between haplotypes, and 21 HA-specific and 29 HB-specific genes were  
345 found. Moreover, within the inversions of HA and HB, six and five genes, respectively, were  
346 identified with their allelic counterparts located outside the inversion regions.

347 Regarding inversion breakpoints, 11 and 10 genes were found to span the breakpoints in HA  
348 and HB, respectively. This set included eight allelic gene pairs and two haplotype-specific  
349 breakpoint-spanning genes whose allelic counterparts did not cross the breakpoints (**Tables**  
350 **S13 and S17**). No genes were detected near the breakpoints of the largest and second-largest  
351 inversions on Chr24 and Chr11, likely due to both inversions overlapping the centromeric  
352 regions (**Figure 2d**). By contrast, a single gene in HA (*HA016982.1*), which spans a breakpoint  
353 on Chr20 and encodes a putative protein (1,184 aa) with 97.5% identity to SMG8 (a nonsense-  
354 mediated mRNA decay factor) [30], was disrupted by the inversion breakpoint on Chr20 in HB  
355 (**Figure 2f**). Further manual checks and read mapping indicated that the HA allele contains a  
356 complete coding sequence with intact ATG start and TGA stop codons, while the corresponding  
357 sequence in HB was incomplete due to the absence of the ATG start codon. This hemizygous  
358 gene exhibited reduced expression, with three transposable element (TE) insertions identified  
359 within its 2 kb upstream region.

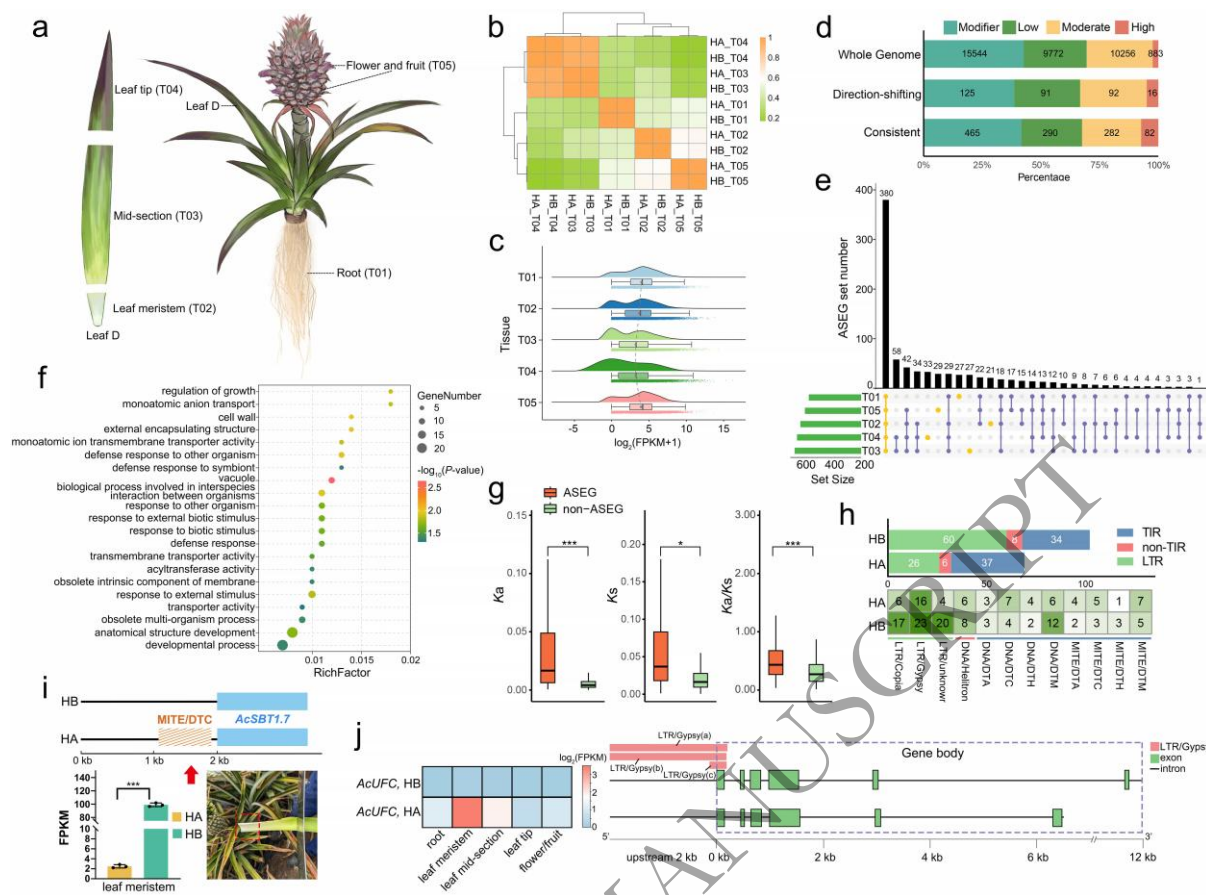
## 360 2.5 *Allele-specific expression and inversion-region expression*

361 We comprehensively identified allele-specific expression genes (ASEGs) from HA and HB  
362 alleles using RNA-seq data collected from five tissues per sample. A comparison of 19,153  
363 allele-defined gene pairs (one-to-one) revealed that 863 pairs (4.28%) showed significant ASE  
364 ( $|\log_2\text{Foldchange}| > 1$  and  $P\text{-value} < 0.05$ ) in at least one tissue (**Table S18**). These ASEGs  
365 were further classified into 713 consistent and 150 direction-shifting ASEGs (**Tables S19-S20**  
366 **and Figures S21-S22**). The percentage of high-impact variants in consistent ASEGs was much  
367 larger than in direction-shifting ASEGs or the genome-wide baseline (**Figure 3d**). Tissue-  
368 specific ASE analysis showed that the majority of ASEGs (44.03%) were expressed across all  
369 five tissues, while smaller proportions were specific to leaf tips (3.82%), flowers/fruits (3.36%),  
370 leaf mid-sections (3.12%), roots (3.12%), and leaf meristems (2.43%) (**Figure 3e**). The ASEGs  
371 expressed in all tissues were mainly enriched in biological processes related to substance  
372 transport and energy metabolism (**Figure S23**), while tissue-specific ASEGs were primarily  
373 associated with defence responses and tissue development (**Figure 3f**). ASEGs also exhibited  
374 significantly higher  $K_a$  and  $K_a/K_s$  ratios than non-ASEGs, indicating their rapid evolutionary  
375 rate (**Figure 3g**). Despite extensive inversions in the pineapple genome, ASEGs within inverted  
376 regions exhibited no significant expression divergence from collinear regions across tissues ( $t$ -  
377 test,  $P > 0.05$ ) (**Figure S24**).

378 TE insertions were detected within the 2 kb upstream promoter regions of 86 out of the 863  
379 ASEG pairs analysed (HA:  $n = 45$ ; HB:  $n = 53$ ; both:  $n = 12$ ) (**Table S21**). The abundance and  
380 types of TEs differed between haplotypes (**Figure 3h**). A potential association between  
381 haplotype-specific TE insertions and allelic expression bias was observed at several loci. A  
382 representative example involves the *AcSBT1.7* gene, which encodes a subtilisin-like protease  
383 known to regulate plant developmental processes, defence, and stress responses [31]. In HA,  
384 an intact MITE/DTC insertion within the 2 kb upstream promoter region coincided with a ~30-  
385 fold suppression of *AcSBT1.7* transcript levels in leaf meristem tissue relative to HB (**Figure**  
386 **3i**). Similarly, the *AcUFC* gene, a homolog of Arabidopsis UPSTREAM OF FLOWERING  
387 LOCUS C (*FLC*), exhibited HA-specific allelic dominance in leaf tissues (**Figure 3j**). Notably,  
388 this regulatory locus interacts with *FLC* and *DFC* (DOWNSTREAM OF *FLC*) to mediate  
389 vernalisation-related signalling pathways in Arabidopsis [32]. Three LTR/Gypsy  
390 retrotransposons were found exclusively in the promoter and coding regions of the HB allele.

391

392



393  
 394 **Figure 3 Functional and evolutionary dynamics of allele-specific expressed genes (ASEGs)**  
 395 **in pineapple subgenomes.** **a.** Five tissues include roots (RNA-seq T01), meristem (T02), mid-  
 396 range (T03), leaf tip of mature leaf (T04), and combined flower and fruit (T05); **b.** Correlation  
 397 analysis and clustering of the expression of alleles in five tissues. Each data point represents  
 398 the mean of three biological replicates for HA or HB in a given tissue; **c.** Raincloud plot of  
 399  $\log_2(\text{FPKM}+1)$  values of genes across five tissues. The average expression level of each tissue  
 400 is connected by a grey dashed line. Similar expression trends were observed for subgenomes;  
 401 **d.** The percentage stacked bar chart of pineapple ASEGs with high, moderate, low, and  
 402 modifier impact variants. Numbers denote genes with  $\geq 1$  variant, and genes may occur in  
 403 multiple categories. “Whole genome” denotes the genome-wide background for comparison.  
 404 “Direction-shifting” ASEGs are defined as those showing a reversal of allelic imbalance  
 405 between tissues. “Consistent” ASEGs maintain stable imbalance toward the same allele across  
 406 all five tissues; **e.** The tissue-specific distribution of ASEGs detected in five tissues in the  
 407 pineapple subgenomes. Yellow dots represent ASEGs that are allele-specific in only one tissue  
 408 or in all tissues; **f.** GO enrichment analysis of tissue-specific ASEGs across the five sampled  
 409 tissues. The y-axis shows significantly enriched GO terms while the x-axis represents the  
 410 RichFactor. Dot size indicates the number of ASEGs associated with each GO term, and color  
 411 corresponds to the statistical significance ( $-\log_{10}$  P-value); **g.** Comparison of  $K_a$ ,  $K_s$ , and  $K_a/K_s$   
 412 between ASEGs and non-ASEGs by independent  $t$ -test.  $*P < 0.5$ ;  $***P < 0.001$ ; **h.** The  
 413 distribution of TE numbers and types between HA and HB. Darker green indicates higher TE  
 414 abundance; **i.** The ASEG (*AcSBT1.7*) with a specific intact MITE/DTC insertion in the  
 415 upstream 2 kb region. The bar plot shows the FPKM of the gene in leaf meristem; **j.** An example  
 416 of an ASEG (*AcUFC*) with LTR/Gypsy insertions in promoter and gene regions. Left:  
 417 Expression heatmap of *AcUFC* in HA and HB across five tissues. Right: Schematic showing  
 418 the LTR/Gypsy insertions.

419 **2.6 Impacts of reference quality and inversions on inter-chromosomal linkage and**  
420 **recombination patterns**

421 Pairs of DArTseq<sup>TM</sup> markers ( $n = 11,879$ ) whose alleles are highly correlated ( $R^2 \geq 0.90$ ) and  
422 lie on different chromosomes are plotted in **Figure 4a**. When the markers were aligned to the  
423 short-read, collapsed ‘Smooth Cayenne’ F153 assembly [23], large numbers of inter-  
424 chromosomal linkages (241 unique marker pairs) were apparent. Realigning the identical  
425 marker set to the long-read, haplotype-resolved F180 assembly (HA and HB) significantly  
426 reduced the number of inter-chromosomal linkages to only 81 and 97 pairs, respectively. F153  
427 contained  $\sim 1.94$ -fold more duplicated marker binding sites than F180, a discrepancy most  
428 likely attributable to assembly artefacts (**Table S22**). It is also noteworthy that whole-genome  
429 alignments with ‘MD-2’ required extensive filtering to remove spurious alignments and local  
430 alignments with DArT markers revealed levels of sequence duplication almost 6-fold higher in  
431 ‘MD-2’ than in the F180.

432 Using 669 high-quality DArTseq<sup>TM</sup> markers common to the F180 haplotypes HA and HB,  
433 recombination events in F180 were scored across 374 F<sub>1</sub> progeny. Across the 25 chromosomes,  
434 SCO frequencies averaged three events per chromosome in HA and HB, while DCO  
435 frequencies averaged one per chromosome for both haplotypes (**Table 2**). Chr3 and Chr4  
436 exhibited the highest SCO activity in HA, whilst Chr3 had the highest activity in the HB  
437 haplotype. DCOs were most frequent on Chr3 and Chr13 (three DCOs each) in HA and HB.  
438 To further examine recombination patterns, we plotted marker density and SCO counts in 50  
439 kb bins across chromosomes 20 and 24 for both haplotypes (**Figure 4b**). In each case, the  
440 positions of inversions are marked by red bars along the x-axis. In HA, Inv553 on Chr20  
441 displayed only one SCO although it appears just within a boundary and is likely the result of  
442 minor positional inaccuracy. The size of the inversion and the lack of crossovers suggests  
443 strong recombination suppression (**Table 3**). A chi-squared pairwise comparison of proportions  
444 using SCO counts within the inverted region of the HA assembly and that of comparable length  
445 regions either side of the inversion, demonstrated a significant difference between the inverted  
446 region and the preceding region (Bonferroni-adj  $P = 0.012$ ), and the inverted region and the  
447 following region (Bonferroni-adj  $P = 4.5e^{-5}$ ) (**Table S23**). The HB assembly showed a similar  
448 difference between SCO counts in the inversion and regions preceding and following with P-  
449 values of 0.013 and 0.0023. By contrast, Inv560 on Chr24 of HA contained 12 SCOs,  
450 suggesting that recombination can still proceed within this region, although the SCO density  
451 in the inverted region was markedly lower than the rest of the chromosome. There was no

452 significant difference in SCO counts between the inversion and that of a similarly sized region  
 453 at the opposite end of the chromosome. The HB assembly did, however, show a significant  
 454 difference (Bonferroni-adj  $P = 0.0058$ ) with only three SCOs within the Inv560. For  
 455 comparison, full marker counts before recombination analysis demonstrated a moderate marker  
 456 density across both chromosomes, including within inverted segments (**Figures 4b**). **Figure 5**  
 457 models the meiotic consequences of a SCO event in inversion heterozygotes of F180, including  
 458 inversion loop formation and the generation of non-viable or imbalanced gametes.

459 **Table 2. Average numbers of single-crossover (SCO) and double-crossover (DCO) events**  
 460 **per chromosome per individual, calculated from 374 F<sub>1</sub> progeny of reciprocal ‘Smooth**  
 461 **Cayenne-F180’ × ‘MD-2’ crosses.** Recombination was inferred independently for markers  
 462 mapped to each phased F180 haplotype, yielding per-chromosome HA and HB estimates.  
 463 Values are means per progeny per chromosome; haplotype-specific chromosome sizes (Mb)  
 464 are provided for context.

Chromosome	Av SCO Counts		Av DCO Counts		Chromosome size			
	HA	HB	HA	HB	HA	HB		
					cM	Mb	cM	Mb
1	3	4	1	1	74	21.1	71	20.5
2	2	2	1	0	69	20.7	65	19.7
3	5	6	3	3	65	20.3	64	20.2
4	5	4	2	2	67	20.3	66	20.0
5	4	4	2	2	62	18.4	57	17.6
6	3	2	1	0	62	18.3	56	16.8
7	2	2	0	0	58	18.3	56	17.6
8	3	3	1	1	65	18.2	64	18.2
9	4	4	2	1	60	17.8	66	19.4
10	3	3	1	1	57	17.2	56	17.1
11	3	3	1	1	53	16.5	51	15.9
12	3	3	1	1	51	16.4	56	17.8
13	4	4	3	3	53	16.3	55	16.9
14	3	2	1	0	52	15.9	57	17.3
15	3	2	1	0	53	15.8	50	15.3
16	3	4	1	2	52	15.8	50	15.5
17	2	2	1	0	50	15.4	49	15.0
18	2	2	1	0	50	15.1	54	15.1
19	2	2	0	0	48	14.8	49	15.1
20	2	2	0	0	48	14.4	50	13.8
21	2	2	1	0	47	14.3	48	14.5
22	2	2	0	0	47	14.0	48	13.9
23	1	2	0	0	41	13.5	40	13.1
24	3	2	1	1	42	12.9	48	14.6
25	1	1	0	0	17	12.5	14	12.1
<i>Av/ Total</i>	3	3	1	1	1,344	414.0	1,337	412.8

466 **Table 3. Double crossover (DCO) events comprised of SCOs detected within or near the**  
 467 **large inversions on chromosomes 20 (Inv553) and 24 (Inv560) in haplotype A (HA) of**  
 468 **‘Smooth Cayenne’ F180.** Positions are given in base pairs, with genetic distances in  
 469 centimorgans (cM). “In” and “Out” indicate whether markers fall inside or outside the inversion  
 470 boundaries. Recombination in F180 was tracked using the haploid QTL2 model with F180 (A)  
 471 vs. ‘MD-2’ (B) allele coding.

Haplotype	Chr	Position 1		Position 2		Distance (cM)
<i>HA: Inv553: 20</i>		<i>11,928,447</i>		<i>13,232,671</i>		
HA	20	10,350,897	Out	13,296,283	Out	9.8
		10,350,897	Out	13,138,073	Out	9.3
		11,843,718	Out	13,846,945	Out	6.7
<i>HA: Inv560: 24</i>		<i>298,388</i>		<i>6,348,311</i>		
HA	24	323,091	In	6,935,708	Out	12.4
		450,224	In	3,488,511	In	10.1
		3,488,511	In	7,272,635	Out	12.6
		3,876,767	In	7,272,635	Out	11.4
		6,270,022	In	9,423,969	Out	10.5
		6,270,022	In	9,533,575	Out	10.9

472

473

474

475

476

477

478

479

480

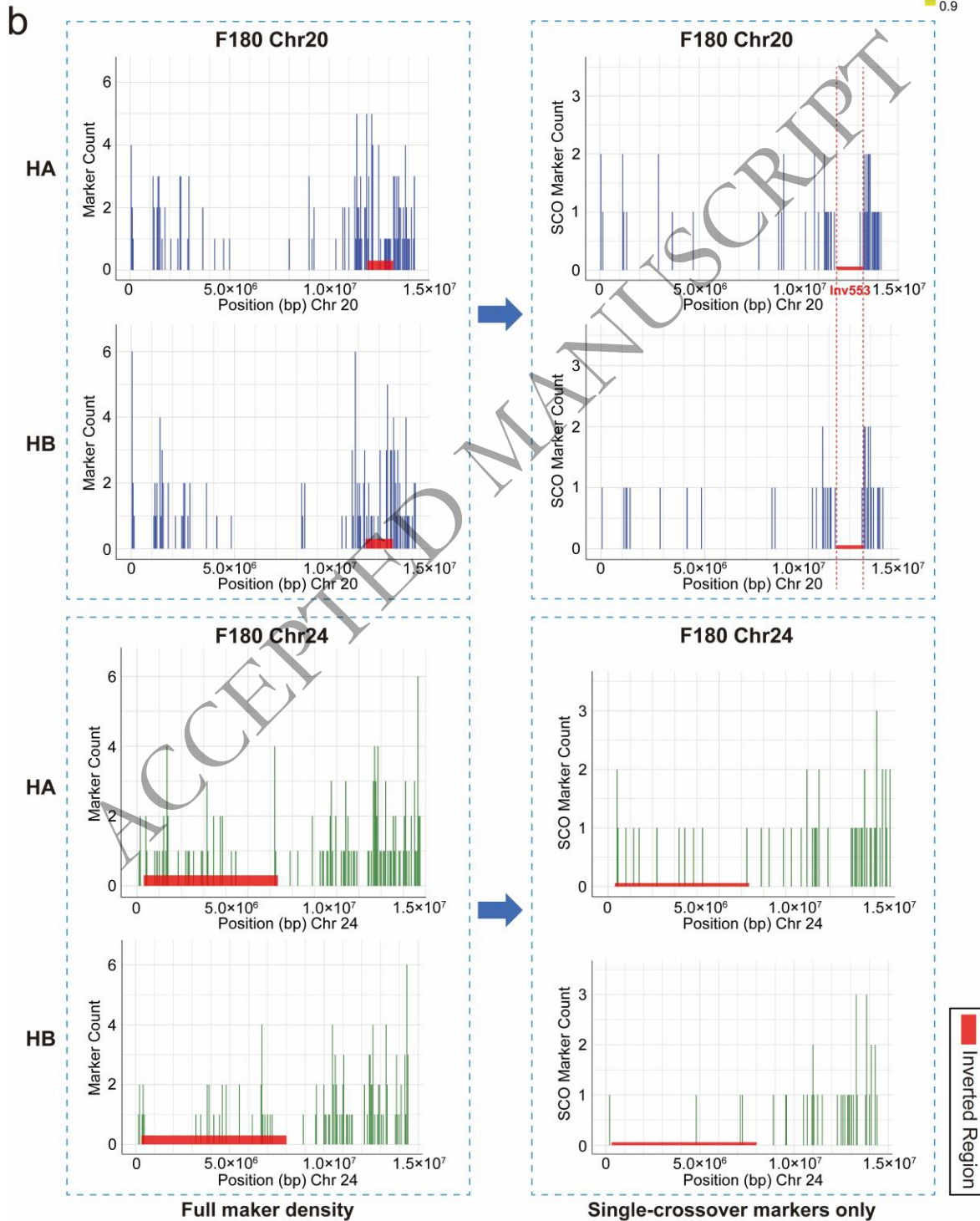
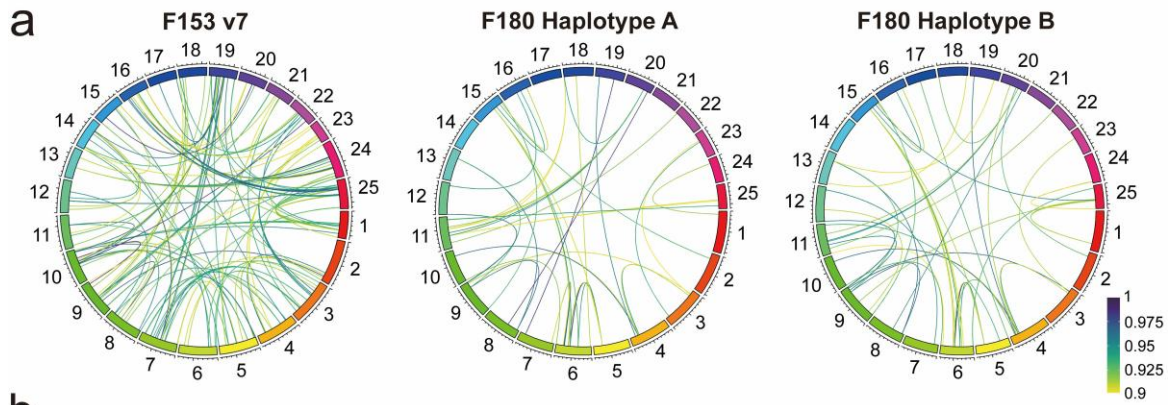
481

482

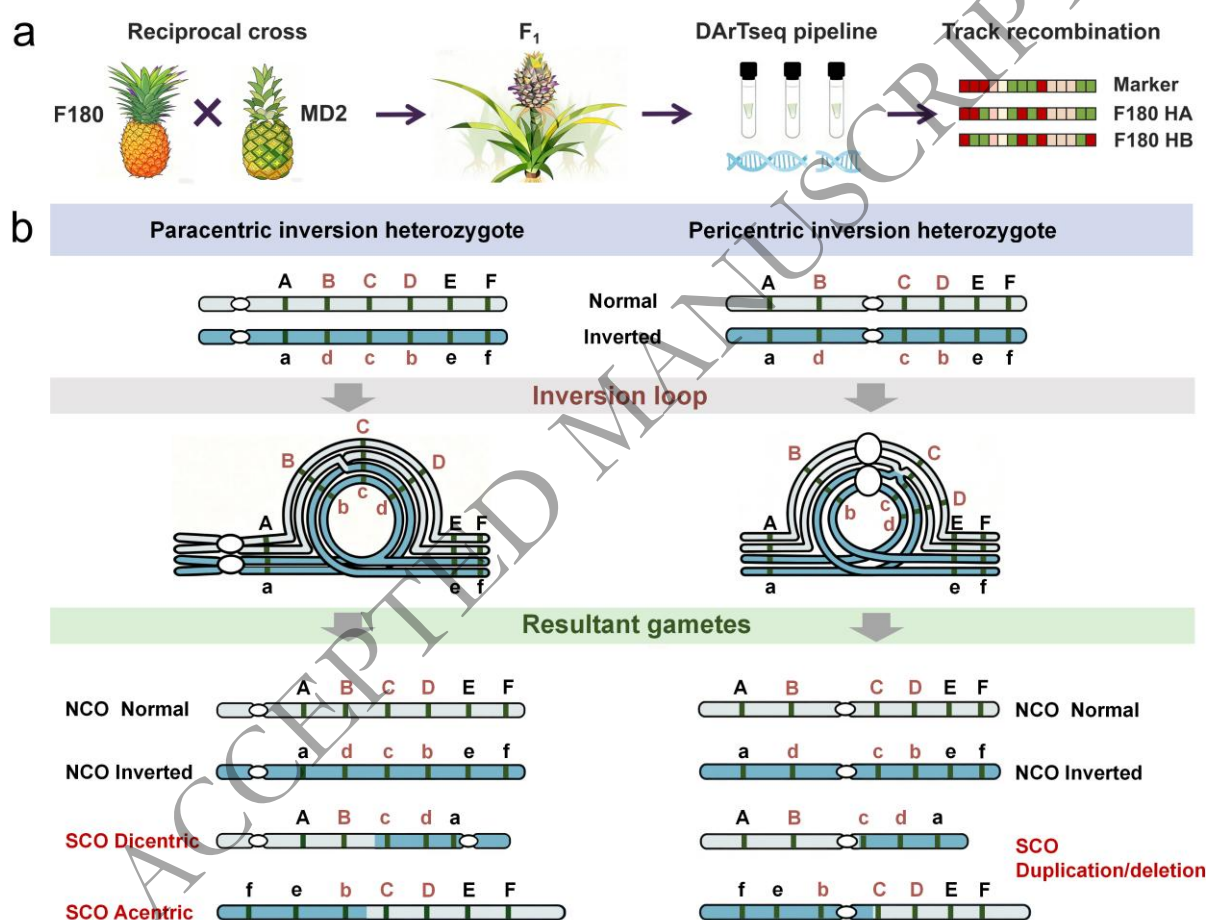
483

484

485



487 **Figure 4. Genome assembly quality and inversion-associated recombination patterns in**  
 488 **‘Smooth Cayenne’ F180. a.** Inter-chromosomal linkage disequilibrium (LD;  $R^2 \geq 0.90$ ) for  
 489 11,879 DArTseq™ loci mapped to three references: the collapsed ‘Smooth Cayenne’ F153  
 490 genome (left) and the two phased F180 haplotypes (HA and HB). Each line joins a marker pair  
 491 located on different chromosomes; line thickness reflects the number of high-LD links. **b.** Full  
 492 (left panels) vs. SCO (right panels) marker distributions in 50 kb bins across Chr20 and Chr24  
 493 for HA and HB. Red horizontal bars indicate two large inversions: paracentric Inv553 on Chr20  
 494 and pericentric Inv560 on Chr24. Blue and green bars represent marker densities on Chr20 and  
 495 Chr24, respectively. All panels show data (full makers or SCO events) pooled from all 374  
 496 progeny. Near absence of SCO markers within Inv553 indicates strong recombination  
 497 suppression, whereas the reduced SCO signals within Inv560 suggest partial suppression. A  
 498 single SCO near the Inv553 boundary represents a minor positional outlier and does not alter  
 499 the absence of crossovers within the inversion.



500

501

502 **Figure 5. Mechanistic model of how inversions suppress recombination. a.** Strategy for  
 503 tracking recombination with DArTseq genotyping in an F<sub>1</sub> population derived from reciprocal  
 504 cross between F180 and ‘MD-2’; **b.** Meiotic consequences of paracentric (left) and pericentric  
 505 (right) inversions in ‘Smooth Cayenne’ F180. In paracentric inversions (left, Chr20), single  
 506 crossovers (SCOs) produce dicentric bridges and acentric fragments, which cause gamete  
 507 inviability and establish a strict recombination coldspot. In pericentric inversions (right, Chr24),  
 508 although SCOs generate unbalanced gametes carrying segmental duplications or deletions,  
 509 short-range double crossovers (DCOs, not depicted) can restore euploidy, permitting limited  
 510 gene flow and attenuating the coldspot. NCO, non-crossover.

## 511 2.7 Haplotype-dynamic NLR repertoires in F180

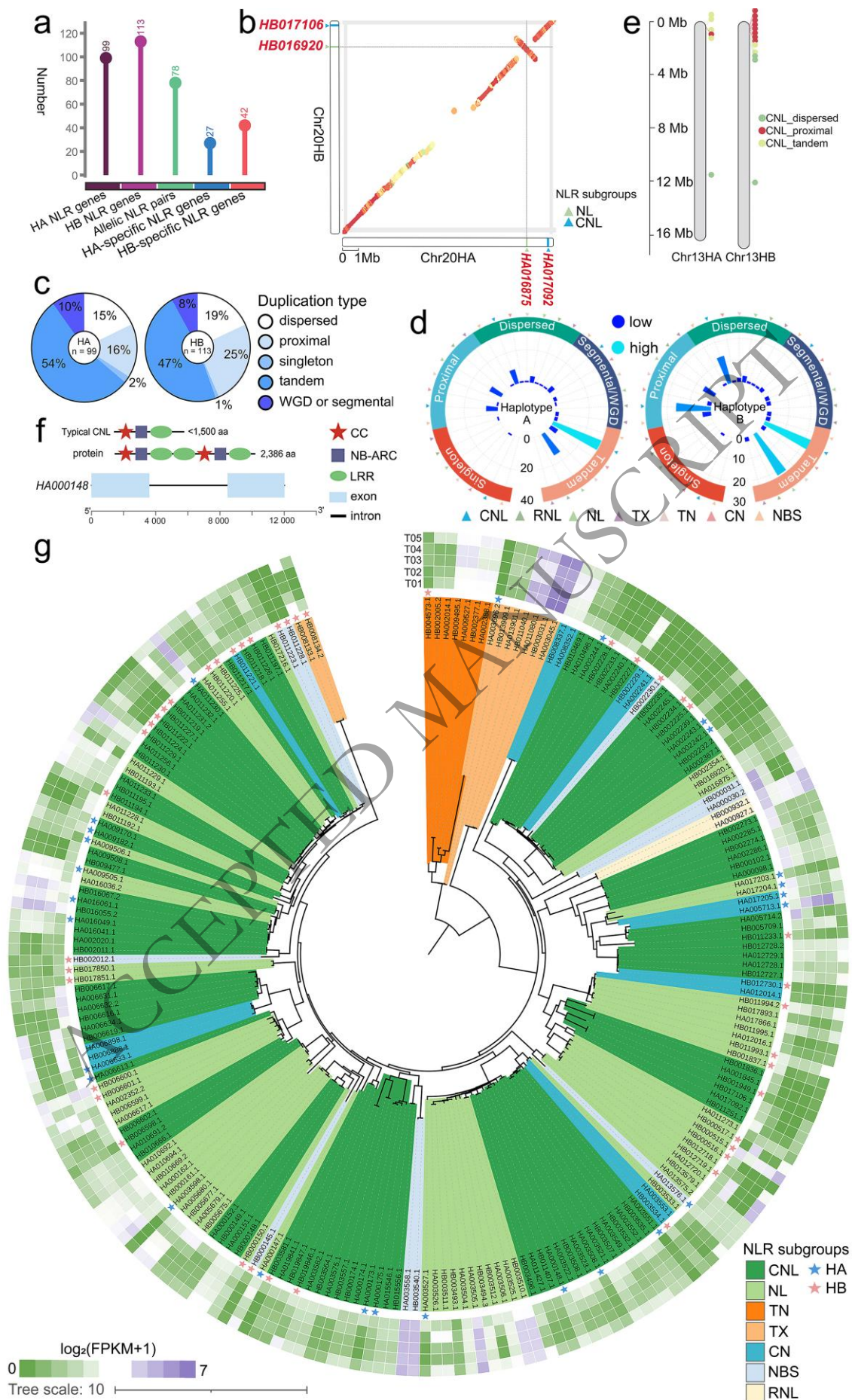
512 To explore the potential disease resistance-related NLR genes in F180, we identified 99 and  
513 113 NLR genes in HA and HB, respectively (**Tables S24-S25**). These NLRs were classified  
514 into seven subgroups, which were distinguished based on the presence or absence of TIR, CC,  
515 NBS, and RPW8 domains. On average, CNL was the most prevalent NLR subgroup, followed  
516 by NL, while TX and RNL were the least common in both haplotypes (**Figure S25 and Table**  
517 **S26**). NLR genes were unevenly distributed across haplotypes, with Chr7, Chr17, Chr23, and  
518 Chr24 lacking NLRs (**Figure S26**). The number of NLR genes on each chromosome varied by  
519 up to 12 genes (on Chr13) between haplotypes. HA uniquely harboured two NLRs on Chr10,  
520 while HB retained one on Chr5. Most of the pineapple NLR genes resided in clusters and were  
521 located near telomeric regions (**Figure S27**). A total of 78 NLR allelic gene pairs were  
522 identified, with 27 and 42 NLR genes specifically detected in the HA and HB haplotypes,  
523 respectively (**Figure 6a and Tables S27-S28**), shedding light on the dynamic nature of R genes  
524 [33]. An NLR allelic gene pair (*HA016875-HB016920*) was exclusively identified within the  
525 Inv553 on Chr20 (**Figure 6b**).

526 Analysis of genomic duplication dynamics delineated distinct evolutionary trajectories of NLR  
527 expansion mechanisms in both haplotypes, with tandem, proximal, and dispersed duplications  
528 serving as principal drivers (**Figure 6c**). Furthermore, divergence in NLR subgroup expansion  
529 mechanisms was associated with distinct duplication modes (**Figure 6d**). CNL and NL  
530 subgroups mainly expanded through tandem and proximal duplication events, respectively.  
531 Interestingly, NLR duplication patterns exhibited marked heterogeneity between haplotypes.  
532 For instance, proximal duplications in HB generated twice as many CNL subgroups as those  
533 in HA, while a parallel trend was observed for the NL subgroup amplified through tandem  
534 duplication. Specifically, proximal duplication on Chr13 HB accounted for the haplotype-  
535 specific CNL expansion (**Figure 6e**). Diverging from canonical CNL architecture, the  
536 atypically large CNL gene *HA000148* (~12 kb) on Chr1 HA possessed a dual CC-NB-ARC  
537 domain (**Figure 6f**), potentially enabling multi-pathogen effector surveillance [34]. No  
538 structural orthologues were identified in HB.

539 Consistent with the genome-wide expression pattern of HA-HB allelic gene pairs, NLR genes  
540 exhibited elevated transcriptional activity in roots and reproductive organs (flowers/fruits)  
541 relative to leaves across both subgenomes (**Figure S28**). The phylogenetic analysis revealed  
542 complex evolutionary diversity and variation in subfamily expression levels (**Figure 6g**).

543 Allelic NLR gene pairs generally clustered together and displayed similar expression patterns.  
544 For example, TX-type genes (e.g., *HA003045-HB003031*) were highly expressed in all five  
545 tissues. No significant expression divergence was observed between HA and HB-specific NLR  
546 genes ( $P > 0.05$ ) (**Figure S29**). Collectively, the majority of HA/HB-specific NLR genes  
547 exhibited low transcriptional activity (FPKM  $< 10$ ; **Figure S30**). However, notable exceptions  
548 included *HA005713* and *HA011230*, which exhibited exceptionally high expression in leaves,  
549 and *HB000517* and *HB012719*, with elevated expression in roots and leaf tips. These genes  
550 warrant further functional characterization to dissect their roles in tissue-specific immune  
551 responses.

ACCEPTED MANUSCRIPT



553 **Figure 6. Genomic architecture and functional diversification of NLR genes in pineapple**  
554 **subgenomes. a.** The lollipop plot showing the distribution of NLR genes between two  
555 haplotypes; **b.** Chromosomal locations of NLR genes on Chr20, indicated by coloured triangles  
556 and dash lines; **c.** Pie charts of the proportions of different duplication types in subgenomes; **d.**  
557 Distribution of NLR genes across duplication events. The outer ring represents five duplication  
558 types, while the inner ring indicates the abundance of NLR subgroups (coloured triangles); **e.**  
559 The haplotype-specific CNL expansion on Chr13HB. Three duplication events are displayed:  
560 dispersed (green), proximal (red), and tandem (light green); **f.** Schematic representation of a  
561 representative CNL protein and its gene structure. The scale bar indicates genomic positions in  
562 base pairs. CC: Coiled-Coil; NB-ARC: Nucleotide-Binding adaptor shared by APAF-1, R  
563 proteins, and CED-4; LRR: Leucine-Rich Repeat; **g.** Raincloud plot showing  $\log_2(\text{FPKM}+1)$   
564 values of NLR genes in HA (left) and HB (right) across five tissues; **h.** Phylogenetic tree of  
565 212 NLR protein sequences from the two subgenomes, constructed using IQ-TREE with 1,000  
566 bootstrap replicates. NLR protein IDs are colour-coded by subgroup. The heatmap shows  
567 relative expression levels across five tissues. Blue and pink stars indicate HA- and HB-specific  
568 NLR genes, respectively.

### 569 3. Discussion

570 The haplotype-resolved genome assembly of the pre-Columbian pineapple cultivar ‘Smooth  
571 Cayenne’ presented here marks a significant advancement over previous assemblies [23],  
572 addressing the limitations of haploid representations in highly heterozygous genomes [1]. The  
573 results from our *de novo* assembly (Section 2.1) demonstrate the power of integrating PacBio  
574 HiFi long-read sequencing (47× coverage) and Hi-C long-range data (40× coverage) [4],  
575 producing a high-quality phased genome consisting of 50 pseudo-chromosomes (433.3 Mb for  
576 HA, 425.2 Mb for HB) with exceptional contiguity (N50 of 16.38–16.86 Mb) and completeness  
577 (BUSCO > 99%, LAI > 23). Here, we were able to construct complete haplotype-resolved  
578 assemblies for all 25 chromosomes, including the identification of all 50 centromeres and  
579 telomeric repeats on 22 out of 25 chromosomes per haplotype, further underscoring the  
580 robustness of this resource. Repetitive elements, predominantly LTR retrotransposons (20.27%  
581 in HA, 19.25% in HB), constitute ~62% of the genome, substantially more than the collapsed  
582 F153 assembly (44%) [23]. The prediction of ~20,248–20,280 (HA and HB, respectively) high-  
583 confidence PCGs, with 97.86% functionally annotated and BUSCO completeness of 99.0–  
584 99.5%, provides a robust foundation for genetic studies and breeding applications. The high  
585 collinearity with the recently released gap-free pineapple reference genome composite [27]  
586 validates the accuracy of our gene annotations, enabling the precise identification of trait-  
587 associated loci to address challenges like flowering asynchrony and pathogen susceptibility  
588 [35–37]. The presence of 2,919–2,962 rRNAs and 870–1,032 tRNAs, with distinct chromosomal  
589 distributions, offers additional targets for studying non-coding RNA functions in pineapple  
590 adaptation, fruit ripening, and development [38,39].

591 Attempts to generate a collapsed genome failed due to significant haplotype divergence and  
592 low sequence collinearity, underscoring the necessity of phased assemblies for capturing the  
593 full complexity of clonally propagated, highly heterozygous crops like pineapple. Our phased  
594 assembly resolved critical structural variations, including the large inversions on chromosomes  
595 20 and 24, which would otherwise remain obscured in collapsed genomes through chimeric  
596 contig formation [18]. Allelic variations and inversion-specific comparisons reveal the extent  
597 of haplotype-specific diversity in ‘Smooth Cayenne’, contributing to the high genome-wide  
598 heterozygosity (1.57%) observed in this accession. Allelic diversity between the two  
599 haplotypes was substantial. For example, 1.5 million SNPs were identified between haplotypes.  
600 However, it is important to note that despite the clonal nature of pineapple cultivation likely  
601 occurring for thousands of years, in turn preserving heterozygosity, the initial domestication  
602 process of these varieties has greatly reduced the relative diversity of common varieties, i.e.  
603 ‘Queen’ and ‘Smooth Cayenne’, in comparison to related *Ananas* cultivars and species [22].  
604 The identification of 19,058 shared allelic gene pairs, alongside 1,196 HA-specific and 1,183  
605 HB-specific genes, highlights the genetic divergence preserved by clonal propagation. *Ka/Ks*  
606 analysis of these haplotype-specific genes displays the predominance of purifying selection  
607 ( $Ka/Ks \leq 1$ ), ensuring genomic stability. Genes under positive selection ( $Ka/Ks > 1$ ), i.e., genes  
608 that display heightened ratios of non-synonymous mutations, show clear selective trends  
609 toward photosynthetic and organelle-related processes and defence-related genes. The presence  
610 of a variety of defence-related genes suggests adaptive evolution against pathogens,  
611 particularly fungi-related genes, as fungi are the major pathogen type of cultivated pineapple  
612 [26,40,41]. There were also ~45 additional genes associated with defence response and the  
613 regulation of biological activity in HA, suggesting genes key to both defence and the regulation  
614 of biological processes, indicating the importance of haplotype-specific contributions to the  
615 resilience and agronomic performance of ‘Smooth Cayenne’. Haplotype-specific NLR gene  
616 repertoires, comprising 99 and 113 genes in HA and HB, respectively, exemplify this diversity,  
617 offering a dynamic resource for combating pineapple diseases.

618 Linkage analysis of the ‘Smooth Cayenne’ × ‘MD-2’ progeny reveals a consistent pattern in  
619 which crossovers cluster near telomeres and are scarce in centromeric regions [11]. Across 374  
620 F<sub>1</sub> progeny, we observed a genome-wide mean of crossover (CO) events that falls squarely  
621 within the one-to-five CO per chromosome range reported for many crops [42] and illustrates  
622 classical CO interference, whereby the formation of one exchange suppresses nearby events  
623 [43]. A total of 1,953 SVs (including 74 inversions, 750 translocations, and 1,129 duplications)

624 and extensive sequence variations (e.g., 1,538,444 SNPs), underscore the haplotype-specific  
625 differences between HA and HB. The three largest inversions, found on Chr24 (6.05–7.71 Mb),  
626 Chr11 (1.51–2.48 Mb), and Chr20 (1.30–1.32 Mb), which were validated by Hi-C matrices and  
627 PacBio read mapping, are particularly significant. These inversions, encompassing large  
628 segments of the F180 genome, have the potential to suppress recombination, preserving  
629 haplotype-specific gene blocks as single inheritance units [44,45]. This has profound  
630 implications for breeding strategies, as it complicates the introgression of valuable traits. The  
631 two largest inversions in ‘Smooth Cayenne’ display contrasting crossover behaviour that does  
632 not fully align with expectations from their structural classifications. Inversion 553 on Chr20,  
633 which is paracentric, shows an absence of recombination. Instead, a higher density of SCOs  
634 occurs at either end (**Figure 4b**). This suggests that recombination suppression within the  
635 inversion may reduce crossover interference and promote exchanges in adjacent regions, as has  
636 been observed in *Drosophila* [46]. In contrast, inversion 560 on Chr24, which is pericentric,  
637 shows limited but detectable recombination, with several SCOs and DCOs occurring within  
638 one half of the inversion. The relatively large size of this pericentric inversion may permit  
639 partial alignment between homologous chromosomes, allowing occasional recombination to  
640 occur despite centromere involvement [47]. Strong suppression is expected in both para- and  
641 pericentric inversions, with additional centromere associated suppression in the latter.  
642 However, the extent of suppression varies with inversion size, breakpoint locations, and  
643 genomic context [12,48], highlighting the complexity of inversion-specific recombination  
644 dynamics. From a breeding perspective, these patterns suggest that Inv553 on Chr20 may act  
645 as a true recombination coldspot, locking ~170 genes into a single haploblock that can be  
646 inherited intact, whereas Inv560 on Chr24 may permit limited gene flux through recombination  
647 or gene conversion. Here, we present a mechanistic model of how inversions suppress  
648 recombination in F180 (**Figure 5**). In the paracentric inversion heterozygotes on Chr20 (1.3  
649 Mb), SCOs produce abnormal dicentric and acentric chromosomes, leading to inviable  
650 gametes, effectively suppressing recombination and creating a strict recombination coldspot in  
651 the region. In contrast, although SCOs in the pericentric inversions on Chr24 (6.7 Mb) can  
652 yield unbalanced gametes, DCOs may restore euploidy, allowing limited recombination and  
653 weakening the coldspot. Such differences are critical for understanding how inversions shape  
654 haplotype inheritance, mutational load, and the potential for linkage drag versus allele fixation  
655 in breeding populations. The absence of these two major inversions in key commercial  
656 pineapple genomes, including ‘MD-2’ and ‘Queen’, as seen in **Figure S16**, may stem from true  
657 varietal differences in ‘Queen’. Whereas for ‘MD-2’, may reflect a somaclonal origin or

658 differing parental contributions. Other explanations could be linked to ‘Smooth Cayenne’  
659 domestication from *Ananas comosus* var. *microstachys* or ancient hybridisation events with  
660 *Ananas comosus* var. *bracteatus* and others [22,23]. The collapsed ‘Smooth Cayenne’ F153  
661 assembly likely masks these inversions by failing to fully resolve haplotypes, highlighting a  
662 technical limitation rather than a true difference from F180 and further demonstrating the  
663 necessity of haplotype-resolved assembly. Evidence also suggests the improper assembly of  
664 the phased ‘MD-2’ genome, which contains numerous artificial duplications. Conserved ASE  
665 across inverted and collinear regions suggests large-scale inversions in the pineapple genome  
666 may preserve local regulatory landscapes, potentially through intact chromatin topology or  
667 compensatory epigenetic mechanisms [49], allowing ASEGs within inverted regions to  
668 maintain expression stability across tissues. By determining the orientation and allele content  
669 of such regions in diverse parents and progeny, strategies can be devised to manipulate genetic  
670 linkage for important agronomic traits. Generating improved assemblies of other commercial  
671 and wild varieties and systematically identifying their structural variants will be critical for  
672 enabling such approaches.

673 When the same 11,879 DArTseq loci as were used in the linkage analysis were mapped to the  
674 earlier collapsed F153 assembly [23], over 240 high-LD links joined non-homologous  
675 chromosomes, likely an artefact of misplaced or merged contigs. After realignment to the  
676 phased HA and HB genomes, those spurious links dropped to less than 100, leaving a much-  
677 reduced set of more reliable long-range associations. For breeders, this reduction in “false LD”  
678 translates directly into cleaner GWAS signals and more accurate genomic-prediction models,  
679 especially where LD decay and haploblock analysis is conducted on a chromosome-by-  
680 chromosome basis [50]. The linkage results reinforce the practical value of using the haplotype-  
681 resolved F180 reference for downstream genetics, instead of non-T2T references. Taken  
682 together with the larger gene set (~20,250 per haplotype) and the ability to trace haplotype-  
683 specific disruptions such as the inversion-breakpoint gene on chromosome 20, the phased F180  
684 resource offers a markedly higher quality resource than previous ‘Smooth Cayenne’  
685 assemblies. QTL mapping for targets such as flowering asynchrony [35,36], pathogen  
686 susceptibility [25,40], and climatic resilience [51] should proceed with higher resolution and  
687 lower false-positive rates, accelerating marker-assisted and genomic-selection pipelines in  
688 pineapple breeding programmes [52].

689 Allele-specific expression plays a critical role in heterosis and adaptive trait regulation [53],  
690 which has led to its extensive study in multiple fruit crops [33,54-56]. Our study presents the  
691 first comprehensive analysis of ASE in pineapple across five distinct tissues (Section 2.5). The  
692 low prevalence of ASE in pineapple (4.28% of genes in at least one tissue) suggests that  
693 homologous alleles are predominantly expressed in a coordinated manner, with transcriptional  
694 balance maintained across most heterozygous loci. ASEGs can be grouped into two primary  
695 patterns: consistent ASEGs (stable allelic bias across tissues) and direction-shifting ASEGs  
696 (tissue-dependent reversal of allelic bias). These patterns align with the dominance and  
697 overdominance hypotheses in the genetic basis of heterosis, respectively [57-59]. Underlying  
698 this divergence, genomic variation annotation highlights that consistent ASEGs harboured a  
699 greater proportion of high-impact variants compared to direction-shifting ASEGs. The  
700 consistent allelic bias towards HA or HB probably reflects compensatory buffering against  
701 deleterious effects of high-impact variants (e.g., protein-truncating mutations) in the alternative  
702 allele [57,60]. Additionally, the number of consistent ASEGs is much higher than that of  
703 direction-shifting ASEGs, mirroring dominance-biased allelic regulation, consistent with  
704 findings in tea plant [58] and bamboo [61]. Both conserved and tissue-specific ASEGs are  
705 involved in a range of physiological processes. These findings imply that ASE may orchestrate  
706 a trade-off between core metabolic maintenance and adaptive plasticity. This may confer  
707 evolutionary advantages by balancing growth-defence resource allocation in diploid pineapple.  
708 Heterozygous and homozygous polymorphic TE insertions were detected in the promoter  
709 regions of ASEGs across pineapple haplotypes. In this case, specific TE insertions (e.g.,  
710 MITE/DTC in *AcSBTL7* and LTR/Gypsy in *AcUFC*) in the upstream regulatory regions are  
711 associated with significantly reduced expression of the corresponding alleles. These insertions  
712 likely act as potent episodic epimutagens, disrupting transcription factor binding or chromatin  
713 accessibility via DNA hypermethylation, thereby enforcing allelic expression imbalance  
714 [62,63]. Similar TE-mediated allelic suppression has been reported in carnation [64] and citrus  
715 [56]. Conversely, TE insertions may also enhance gene expression, as revealed in red-skinned  
716 apples, where a retrotransposon upregulated *MdMYB1* to drive anthocyanin accumulation [65].

717 NLRs play essential roles in plant immunity by detecting pathogen effectors and show  
718 extensive structural diversification among plant species [66,67]. Our newly assembled T2T and  
719 fully phased pineapple genome enables high-resolution analysis of NLR polymorphisms  
720 between haplotypes (Section 2.7). NLR genes typically cluster near the telomeres, where they  
721 physically interact to optimise immune responses [68]. Proximal duplication likely drove the

722 expansion of CNL genes specific to HB on Chr13. This expansion may be favoured by  
723 pathogen-driven selection, thereby increasing immune gene copy numbers in HB (HA: 99 vs.  
724 HB: 113). We identified 78 allelic NLR pairs, with 27 and 42 NLR genes specifically found in  
725 HA and HB, respectively. This uneven inter-haplotype gene distribution likely reflects an  
726 evolutionary trade-off between persistent pathogen pressure and the fitness costs of expanding  
727 NLR repertoires [69,70]. HA- and HB-specific NLR genes (e.g., *HA005713* and *HB000517*)  
728 exhibited high expression levels (FPKM > 100) in certain tissues, suggesting localised immune  
729 functions. Their pathogen-response specificity warrants further investigation.

#### 730 **4. Conclusion**

731 Beyond pineapple, this study underscores the broader value of haplotype-resolved genomes for  
732 understanding heterozygous, clonally propagated crops. The high contiguity (N50 of 16.38–  
733 16.86 Mb), completeness (BUSCO > 99%), and resolution of almost all 50 telomeres and  
734 centromeres position the F180 assembly as a high-quality reference, surpassing most existing  
735 pineapple genomes. A comparison of the two haplotypes displays the large differences inherent  
736 in the ‘Smooth Cayenne’ pineapple genome, wherein 1.5 million SNPs and 1,953 large  
737 structural variants, including 74 inversions, 750 translocations, and 1,129 segmental  
738 duplications were identified. Inversions were found to be a particularly important process,  
739 collectively reshaping 3–4% of the genome. Among these, a 1.3 Mb paracentric inversion on  
740 chromosome 20 acts as a true recombination coldspot, whereas a 6 Mb pericentric inversion  
741 on chromosome 24 still supports gene flow, likely via short double crossovers. This contrast  
742 demonstrates that chromosomal context, not size alone, dictates whether inversions hinder or  
743 permit allelic shuffling, a principle likely to hold in many perennial crops where undetected  
744 inversions can stall genetic gain. These insights extend to other perennial crops where large  
745 inversions may silently constrain genetic gain. The resource therefore enables fundamental  
746 studies of how heterozygosity, structural variation and allele-specific expression interact in  
747 adaptation and domestication, while guiding both conventional and genome-editing strategies  
748 to broaden pineapple’s genetic base and mitigate vulnerability to biotic and abiotic stress.

## 749 5. Materials and Methods

### 750 5.1 Sample collection, nucleic acid extraction and sequencing, and Hi-C library 751 preparation and sequencing

752 For long-read DNA sequencing, meristematic tissue taken from the D leaf of a mature ‘Smooth  
753 Cayenne’ clone F180 pineapple plant was collected from a commercial farm, ‘Sandy Creek  
754 Pineapples’ located in the Glasshouse Mountains, QLD (-26.872209, 152.962088). The sample  
755 was immediately snap-frozen using liquid nitrogen, and then stored at -80 °C. The Tissue Lyser  
756 II instrument under cryogenic conditions was used to homogenise the sample. The CTAB  
757 (Cetyltrimethylammonium bromide) method, as described in Furtado (2014) [71], was used to  
758 extract whole genomic DNA. Long-read sequencing was completed at the Institute for  
759 Molecular Bioscience, University of Queensland, using two PacBio HiFi SMRT Cells on the  
760 Sequel II instrument. For the annotation, total RNA was extracted from five different tissues,  
761 including three leaves, one root, and a combined fruit and flower tissue (**Figure 3a**), using the  
762 NucleoSpin RNA Plant Mini Kit from Macherey-Nagel. Total RNA was sequenced using the  
763 MGI T7 DNBSEQ platform. RNA-seq library construction followed standard protocols,  
764 beginning with RNA quality assessment using an Agilent 4200 Tape Station to ensure RNA  
765 integrity (RIN  $\geq$  6). Ribosomal RNA (rRNA) was depleted, then the remaining RNA  
766 fragmented, followed by first- and second-strand cDNA synthesis, end repair, A-tailing, and  
767 adapter ligation using the MGIEasy RNA library prep kit. Following ligation, libraries were  
768 PCR amplified and size-selected (~300 bp). The final step involved circularisation and rolling  
769 circle amplification to generate DNA nanoballs (DNBs), which were then loaded onto patterned  
770 flow cells for sequencing. Immature leaves attached to the shoot apex were collected from F180  
771 for Hi-C library preparation and Illumina sequencing, which was carried out by DNA Zoo at  
772 the University of Western Australia. Library preparation and analysis were conducted  
773 following the *in situ* Hi-C protocol and data pipeline described by Rao et al. (2014) [72].

### 774 5.2 Genome assembly

775 Prior to genome assembly, the *k*-mers (21-mers) of PacBio HiFi subreads were counted using  
776 Jellyfish (v2.2.10) [73], and then GenomeScope2 [74] was applied to estimate the genome sizes,  
777 heterozygosity, and repeat contents based on the 21-mers count histograms. We adopted a  
778 custom assembly workflow to assemble the haplotype-resolved genome (**Figure S3**). First,  
779 based on HiFi and Hi-C reads, Hifiasm (v0.19.5-r587) [5] was applied to build a phased diploid  
780 assembly, generating primary contigs of 603.89 Mb and 424.73 Mb, with N50 values of 6.83

781 Mb and 6.90 Mb, respectively. Then, with the aid of Hi-C data, HapHiC (v1.0.5) [75] was used  
782 to *de novo* cluster, reassign, order, and orient the first (HA, 603.89 Mb) and second (HB, 424.73  
783 Mb) haplotype contigs. Initially, the HA contigs were anchored into 25 groups with sizes  
784 ranging from 9.48 Mb to 38.56 Mb, while the HB contigs were clustered into 25 groups with  
785 sizes ranging from 5.25 Mb to 20.95 Mb. After constructing the final scaffolds, HapHiC  
786 automatically generated a shell script for visualisation and curation in JuiceBox. Then,  
787 JuiceBox [76] was used for visualising Hi-C interactions and performing manual modification  
788 to obtain 25 pseudo-chromosomes in each haplotype assembly (HA and HB). Telomere regions  
789 were identified using quarTeT (v1.2.1) [77] based on the telomeric seven-base repeat pattern  
790 (CCCTAAA/TTTAGGG at the 5'/3' end) in the plant genome. In the manual modification  
791 process, contigs containing telomeres were manually examined to ensure their proper  
792 placement at the chromosome ends. Some unplaced contigs were also rescued according to the  
793 strong Hi-C interaction signals. Finally, the BAM file and agp file containing the contigs from  
794 both haplotype assemblies were regenerated and JuiceBox was used to visualise the combined  
795 Hi-C contact matrix. Some redundant or excess contigs incorrectly phased by Hifiasm were  
796 manually swapped between HA and HB according to Hi-C signals. In addition to the telomeres  
797 at the chromosomal level, we also identified seven small unanchored contigs (ranging from  
798 21.5 kb to 96.6 kb) with telomeres at one end. The initial genome assembly was searched using  
799 BLASTN (2.15.0+) [78] against the chloroplast and mitochondrial genome sequences of  
800 *Ananas comosus* downloaded from the NCBI GenBank to remove the organelle-derived  
801 contigs. Next, genome self-alignment was performed using Minimap2 (v2.28) [79] to eliminate  
802 redundant sequences, yielding the final F180 genome assembly. Only three of the seven small  
803 unanchored contigs with telomeres remained in the final assembly. BLASTN was further used  
804 to determine whether the telomeres on these unanchored contigs match the telomere-lacking  
805 ends of the chromosomes.

### 806 **5.3 Assessment of genome assembly quality**

807 Multiple approaches were involved to assess the quality of genome assembly and phasing.  
808 JuiceBox was used for visualising the Hi-C interaction signals. QuarTeT [77] was used to  
809 perform telomere identification and centromere candidate prediction. The assembly contiguity  
810 was assessed using QCAST (v5.0.2) [80]. The completeness of the genome and gene model  
811 was assessed using the latest version of BUSCO (v5.8.0) [81] (Benchmarking Universal  
812 Single-Copy Orthologs) with the Miniprot [82] and MetaEuk [83] tool, based on 1,614 highly

813 conserved plant core genes from the “embryophyta\_odb10” dataset. The base-level accuracy  
814 and completeness of the genome assembly were evaluated with Merquy (v1.3) [84], using  
815 consensus quality value (QV) scores estimated from the *k*-mer spectrum of the PacBio HiFi  
816 read set with a 19-mer. The LTR Assembly Index (LAI), calculated using the LTR\_retriever  
817 pipeline (v2.9.0), was also used to assess assembly quality. The PacBio HiFi reads, stLFR  
818 paired-end short reads, and RNA-seq reads were mapped to the genome assembly to assess  
819 assembly accuracy using Minimap2 [79], BWA (v0.7.17-r1188) [85] and HISAT2 (v2.2.1) [86],  
820 respectively. SAMtools (v1.9) [87] subcommands ‘depth’ and ‘coverage’ were used to obtain  
821 base depth and coverage metrics for each chromosome.

#### 822 5.4 Genome repeats and gene annotation

823 Before annotation, we manually curated the chromosome order and orientation for the two  
824 haplotypes based on the Hi-C interaction heatmap, ensuring chromosomes were arranged from  
825 largest to smallest, with the short arms in front and long arms at the rear. Genomic transposable  
826 and repeat elements (TEs) were *de novo* identified by RepeatModeler2 (v2.0.4) [88] with  
827 default parameters. The output file “genome-families.fa” was subsequently used to mask the  
828 repeat regions in the genome with RepeatMasker (v4.1.5) [89] using the “-xsmall” soft-  
829 masking option to enable prediction of PCGs containing repetitive sequences [90]. The quality-  
830 controlled and trimmed paired-end RNA-seq reads were aligned to the soft-masked genomes  
831 using HISAT2 (v2.2.1) [86]. PCGs were annotated using the BRAKER3 annotation pipeline  
832 [91], which combines GeneMark-ETP, AUGUSTUS, and the TSEBRA combiner,  
833 incorporating both transcriptome evidence and a large protein database, with statistical models  
834 that are iteratively trained and specifically tailored to the target genome. The RNA-seq  
835 transcripts were obtained from five pineapple tissue samples: root, leaf meristem, leaf mid-  
836 section, leaf tip, and a combined flower and fruit sample. The gene annotation completeness of  
837 predicted PCGs in HA and HB was evaluated using BUSCO (v5.8.0). Functional annotation of  
838 PCGs was performed through: (1) BLASTP searches (E-value  $\leq 1e-5$ ) against NCBI non-  
839 redundant (Nr; *Viridiplantae* taxonomy) and Swiss-Prot protein databases; (2) functional  
840 assignment using Gene Ontology (GO), Kyoto Encyclopedia of Genes and Genomes (KEGG),  
841 and InterProScan; and (3) orthologous group classification by Clusters of Orthologous Groups  
842 (COGs) and EggNOG. Barrnap (v0.9) (<https://github.com/tseemann/barrnap>) was used to  
843 predict rRNAs, including 5S, 5.8S, 18S, and 28S rRNA genes, in the F180 genome. The tRNA  
844 genes were predicted by tRNAscan-SE (v2.0.12) [92] with default parameters.

845 **5.5 Identification of structural variations, allelic gene pairs and allele-specific expression**

846 Genome collinearity and structural variations, including presence/absence variants (PAVs; >50  
847 bp insertions/deletions), between F180 subgenomes HA and HB were identified using Synteny  
848 and Rearrangement Identifier (SyRi v1.7.0) [93] and a custom Python script. Single nucleotide  
849 polymorphisms (SNPs) were called using GATK (v4.2.6.1)[94] through: (1) variant calling  
850 (SNPs and InDels) with HaplotypeCaller (--stand-call-conf 30); (2) filtering with  
851 VariantFiltration using the parameters: “DP < 1 || DP > 400 || QD < 2.0 || FS > 60.0 || MQ <  
852 40.0 || MQRankSum < -12.5 || ReadPosRankSum < -8.0”; and (3) final selection of high-  
853 confidence SNPs with SelectVariants. The SnpEff [95] (v5.2) software was used to annotate  
854 the potential effects of variations within the pineapple ASEGs based on the gene annotations  
855 of the reference genome (HA) with default parameters. InDels were identified using  
856 Assemblytics (v1.2.1) [96] based on whole-genome alignments generated by NUCmer from  
857 MUMmer [97] (v4).

858 Whole-genome nucleotide collinearity was assessed using Minimap2 [79], NUCmer [97], and  
859 Mauve (v2.40) [98] under default parameters. Dot plots were generated from Minimap2 PAF  
860 files using a modified paf2dotplot script (<https://github.com/zengxiaofei/paf2dotplot>), and  
861 from NUCmer .nc.coords files using a custom Python script. Whole-genome dot plot  
862 alignments were performed to compare the two haplotypes and to assess F180 haplotypes  
863 against published assemblies of pineapple cultivars and related varieties, confirming the  
864 presence or absence of inversions 553 and 560 across lineages. Sequences used for the dot plots  
865 are detailed in **Table S14**. Synteny analysis at the PCG level was performed to assess the impact  
866 of structural rearrangements on gene organization. All annotated protein sequences from HA  
867 and HB were compared using all-against-all BLASTP searches ( $E\text{-value} \leq 1e\text{-}5$ ), and the  
868 resulting homologous pairs were fed into MCScanX [29] to identify collinear gene blocks. A  
869 synteny-based toolkit SynGap (v1.2.5) [99] *genepair* module was applied to profile integrative  
870 gene synteny between homologous chromosomes, using the HA and HB genome sequences  
871 and the GFF3 annotation files of the longest transcripts as input. The output file  
872 (“\*.final.genepair”), including the syntenic and the best two-way BLAST gene pairs, was  
873 considered as the allele table.

874 An RNA-seq analysis pipeline was employed to examine the differential expression of each  
875 allele pair in the allele table. In detail, for triplicate RNA-seq reads from five distinct tissues,  
876 roots, leaf meristems, leaf mid-sections, leaf tips, and flowers/fruits, we first used fastp (v0.23.4)

877 to remove adapters and filter out low-quality bases/reads with parameters “-q 20 -l 75 -w 8”.  
878 The trimmed reads were then aligned to the combined F180 haplotype assembly using HISAT2  
879 (v2.2.1). The output SAM (Sequence Alignment/Map) files were filtered to retain uniquely  
880 mapped reads (NH:i:1). The filtered SAM files were then converted to BAM (Binary  
881 Alignment/Map) format, sorted, and indexed by SAMtools, and subsequently fed into  
882 FeatureCounts (v1.6.2) for quantification of FPKM (Fragments Per Kilo base of exon model  
883 per Million mapped reads) values. The expression count matrix for all allele pairs was  
884 constructed, and differential expression was assessed using DESeq2 (v1.46.0) in R. The  
885 correlation analysis of the allelic expression was conducted using corplot (v0.84) package in  
886 R software.

887 Allele-specific expression (ASE) was defined by an absolute  $\log_2$  fold change  $> 1$  with a P-  
888 value  $< 0.05$ . ASE patterns were classified into two categories: consistent ASE and direction-  
889 shifting ASE [57]. ASE patterns were classified as consistent if they exhibited stable allelic  
890 imbalance toward one allele across all five tissues (fixed  $\log_2$ FoldChange sign), or as direction-  
891 shifting if allelic imbalance direction showed inter-tissue reversal (sign flip). The  $\log_2$  fold  
892 change heatmaps were generated using the R packages ggplot2 (v3.5.1), tidyr (v1.3.1), and  
893 dplyr (v1.1.4). The upset plot was created using UpSetR package (v1.4.0). For the ASEG sets  
894 derived from different tissues, GO enrichment analysis was performed using the “GO  
895 enrichment” function in TBtools (v1.113) [100]. The box plots of  $K_a$ ,  $K_s$ , and  $K_a/K_s$  between  
896 ASEGs and non-ASEGs were visualised by ggplot2 and ggsignif (v0.6.4). Expression patterns  
897 of ASEGs were created using TBtools with the “Log<sub>2</sub> Scale” option. A range of genomic  
898 features, including chromosome length, gene density, repetitive elements, GC content, SVs,  
899 SNPs, Indels, and allele-specific expression genes (ASEGs) were visualised on each pseudo-  
900 chromosome using Circos (v0.69.8) [101]. Transposable elements were also independently  
901 identified in each haplotype assembly using EDTA [102] (v2.2.2) with the “-sensitive 1 -anno  
902 1” parameter settings. The TEanno.gff3 files were used to analyse the correlation with gene  
903 expression.

#### 904 ***5.6 Analysis of recombination modulation by haplotype-specific inversions***

905 A total of 953 pineapple plants were generated for linkage analyses, including 512 seedlings  
906 spanning 17 biparental families plus 53 parents, and 374 seedlings from reciprocal crosses  
907 between ‘Smooth Cayenne’  $\times$  ‘MD-2’. Freeze-dried leaf bases (meristem) were sent to  
908 Diversity Arrays Pty Ltd for DArTseq™ genotyping. Genomic DNA was digested with

909 PstI/MseI, ligated to barcoded adapters, and SNPs were called with the DArTsoft14 pipeline  
910 [103,104]. A total of 32,141 marker sequences were aligned to the collapsed ‘Smooth Cayenne’  
911 F153 v7 assembly and to both phased haplotypes of the F180 v4 T2T genome. Aligned genomic  
912 regions with an E-value  $< 1e-7$  were retained and the proportion of markers with  $> 1$  aligned  
913 region calculated. After filtering for call rate ( $> 0.75$ ), minor allele frequency ( $> 0.025$ ), and  
914 unique chromosomal position, 11,879 high-quality SNPs were retained. Missing genotypes  
915 were imputed in TASSEL v5.2.44 with the LD-KNNi algorithm at 96% masking accuracy  
916 [105,106], and this data was used to calculate inter-chromosomal linkage disequilibrium (LD)  
917 ( $R^2 \geq 0.90$ ) and visualised as Circos plots with the GWLD package [107]. To track  
918 recombination derived from F180, a subset of this data comprising the parent cultivars F180  
919 and ‘MD-2’ and 374 of their progeny was filtered to retain only markers that were heterozygous  
920 in F180 and homozygous in ‘MD-2’. This ensured that segregation patterns in the progeny  
921 reflected recombination events occurring exclusively in F180, while ‘MD-2’ contributed a  
922 constant allele background. Progeny SNP calls were recoded under the A/B system, where “A”  
923 represented the first F180 allele together with the ‘MD-2’ allele, and “B” represented the  
924 second F180 allele together with the ‘MD-2’ allele. The marker file was analysed in R using  
925 the QTL2 package with the ‘haploid’ model [108]. The Kosambi mapping function (scale factor  
926  $2^7$ ) was applied to convert physical distances into genetic distances, with recombination  
927 frequencies capped at 0.5, producing map lengths of up to  $\sim 130$  cM for the largest  
928 chromosomes. The function `locate_xo()` in the Qtl2 package was used to identify crossover  
929 events. This function uses the genotype data and the genetic map to locate crossover positions  
930 represented in each individual. Genetic probabilities with an assumed error rate of 0.01 were  
931 used to account for genotyping errors. Additionally, the `maxmarg()` function was used to filter  
932 genotypes based on a minimum posterior probability (`minprob = 0.8`), ensuring that only high-  
933 confidence genotype calls were retained. Single crossovers (SCOs) were defined as intervals  $\geq$   
934 0.25 cM. The individual seedling data were then collated into one dataset for F180 using a  
935 sliding window of 5 cM to avoid counting crossover events more than once during collation,  
936 and double crossovers (DCOs) were identified when two consecutive SCOs were separated by  
937  $\leq 20$  cM.

### 938 **5.7 Identification of putative pineapple NLR genes**

939 The RGAugury [109] pipeline was used to predict NLR genes in each haplotype genome (HA  
940 and HB). The putative NBS domain-encoding genes were further classified into different

941 subgroups based on their domain structures, including CN (coiled-coil (CC) and NB-ARC),  
942 CNL (CC, NB-ARC and leucine-rich repeat (LRR)), NBS (NB-ARC), NL (NB-ARC and LRR),  
943 RNL (RPW8, NB-ARC and LRR), TN (Toll/interleukin-1 receptor (TIR) and NB-ARC), TX  
944 (TIR and unknown domain). The distribution of NLR genes on chromosomes was drawn using  
945 RIdeogram (v0.2.2) [110] in R. NLR protein sequences were aligned using MUSCLE (v3.8.1)  
946 [111] with default settings and then subjected to IQ-TREE [112] (v2.0.3) for phylogenetic tree  
947 construction with 1,000 bootstrap replicates. The “intersect” function embedded in BEDTools  
948 (v2.30.0) [113] was used to assess the effects of inversion on NLR genes.

949

## 950 **Acknowledgements**

951 This work was supported by Hort Innovation, “Building an Advanced Genomics Platform for  
952 Australian Horticulture (AS21006)”. Pineapple material was provided by Sam Pike (Sandy  
953 Creek Pineapples) at multiple growth stages. We thank the University of Queensland (UQ)’s  
954 Institute for Molecular Bioscience for the generation of the Pac-Bio data, and the Queensland  
955 Department of Primary Industries for funding. We also thank DNA Zoo for the generation of  
956 Hi-C libraries and sequencing, and the Rajeev Varshney sequencing lab at Murdoch University  
957 for the library preparation and generation of RNA-seq data. Analyses were performed on the  
958 Bunya high-performance computing cluster hosted by the Research Computing Centre at UQ.  
959 We acknowledge the China Scholarship Council (Grant No. 201908350014) for supporting  
960 JF’s research visit to UQ. The funder had no role in study design, data collection and analysis,  
961 decision to publish, or preparation of the manuscript.

## 962 **Author Contributions**

963 RV and RH conceived the study, secured funding, and supervised the project. PJM, JF, and GS  
964 collected plant material and oversaw laboratory workflows, including DNA extraction and QC.  
965 RV, VG, PK, PJM, RH, RGRC, and AC generated long-read, short-read, and Hi-C data. JF  
966 performed genome assembly, polishing, scaffolding, and quality assessment. JF, GS, MW, LZ,  
967 JW, and YZ conducted repeat and gene annotation, comparative genomics and structural  
968 variation analyses, including inversion discovery and recombination analysis. JF, GS, MW, and  
969 PJM prepared figures and curated data for public deposition. PJM, JF, GS, and MW drafted the  
970 manuscript; all authors contributed to review and editing, approved the final version, and  
971 agreed to be accountable for the work.

972 **Data Availability**

973 All Pac-Bio CCS, Hi-C, and RNA-seq data has been deposited in the Genome Sequence  
974 Archive (GSA) in National Genomics Data Center (NGDC, <https://ngdc.cncb.ac.cn/>) under  
975 project ID PRJCA041511. Whole genome sequence data have been deposited in the Genome  
976 Warehouse (GWH) at NGDC under accession number GWHGEOA00000000.1 (Haplotype A)  
977 and GWHGEOB00000000.1 (Haplotype B).

978 **Conflict of Interest Statement**

979 The authors have no conflicts of interest to declare. All coauthors have seen and agree with the  
980 contents of the manuscript and there are no financial interests to report. We certify that the  
981 submission is original work and is not under review at any other publication.

982 **Supplementary Data**

983 Supplementary data is available at *Horticulture Research* online.

ACCEPTED MANUSCRIPT

- 986 1. Guiglielmoni, N., Houtain, A., Derzelle, A. *et al.* Overcoming uncollapsed haplotypes in long-  
987 read assemblies of non-model organisms. *BMC Bioinformatics*. 2021;**22**: 303
- 988 2. Zhang, G., Fang, X., Guo, X. *et al.* The oyster genome reveals stress adaptation and  
989 complexity of shell formation. *Nature*. 2012;**490**: 49-54
- 990 3. Penaloza, C., Gutierrez, A.P., Eöry, L. *et al.* A chromosome-level genome assembly for the  
991 Pacific oyster *Crassostrea gigas*. *GigaScience*. 2021;**10**: giab020
- 992 4. Belton, J.M., McCord, R.P., Gibcus, J.H. *et al.* Hi-C: a comprehensive technique to capture  
993 the conformation of genomes. *Methods*. 2012;**58**: 268-76
- 994 5. Cheng, H., Concepcion, G.T., Feng, X. *et al.* Haplotype-resolved de novo assembly using  
995 phased assembly graphs with hifiasm. *Nature Methods*. 2021;**18**: 170-5
- 996 6. Benevenuto, J., Ferrão, L.F.V., Amadeu, R.R. *et al.* How can a high-quality genome assembly  
997 help plant breeders? *GigaScience*. 2019;**8**: giz068
- 998 7. Liu, Y., Du, H., Li, P. *et al.* Pan-genome of wild and cultivated soybeans. *Cell*. 2020;**182**: 162-  
999 76. e13
- 1000 8. Huang, K., Andrew, R.L., Owens, G.L. *et al.* Multiple chromosomal inversions contribute to  
1001 adaptive divergence of a dune sunflower ecotype. *Molecular Ecology*. 2020;**29**: 2535-49
- 1002 9. Zhang, W., Tan, C., Hu, H. *et al.* Genome architecture and diverged selection shaping pattern  
1003 of genomic differentiation in wild barley. *Plant Biotechnology Journal*. 2023;**21**: 46-62
- 1004 10. Boideau, F., Richard, G., Coriton, O. *et al.* Epigenomic and structural events preclude  
1005 recombination in *Brassica napus*. *New Phytologist*. 2022;**234**: 545-59
- 1006 11. Wilkinson, M.J., McLay, K., Kainer, D. *et al.* Centromeres are hotspots for chromosomal  
1007 inversions and breeding traits in mango. *New Phytologist*. 2025;**245**: 899-913
- 1008 12. Kirkpatrick, M. How and why chromosome inversions evolve. *PLoS Biology*. 2010;**8**:  
1009 e1000501
- 1010 13. Navarro, A., Betrán, E., Barbadilla, A. *et al.* Recombination and gene flux caused by gene  
1011 conversion and crossing over in inversion heterokaryotypes. *Genetics*. 1997;**146**: 695-709
- 1012 14. De-Kayne, R., Gordon, I.J., Terblanche, R.F. *et al.* Incomplete recombination suppression  
1013 fuels extensive haplotype diversity in a butterfly colour pattern supergene. *PLoS Biology*.  
1014 2025;**23**: e3003043
- 1015 15. Schwartz, C., Lenderts, B., Feigenbutz, L. *et al.* CRISPR-Cas9-mediated 75.5-Mb inversion  
1016 in maize. *Nature plants*. 2020;**6**: 1427-31
- 1017 16. Wellenreuther, M. & Bernatchez, L. Eco-evolutionary genomics of chromosomal inversions.  
1018 *Trends in Ecology & Evolution*. 2018;**33**: 427-40
- 1019 17. Huang, K., Ostevik, K.L., Elphinstone, C. *et al.* Mutation load in sunflower inversions is  
1020 negatively correlated with inversion heterozygosity. *Molecular Biology and Evolution*.  
1021 2022;**39**: msac101
- 1022 18. Gemenet, D.C. & Khan, A. Opportunities and challenges to implementing genomic selection  
1023 in clonally propagated crops. *Genomic Selection for Crop Improvement: New Molecular  
1024 Breeding Strategies for Crop Improvement*. 2017;185-98
- 1025 19. Bisognin, D.A. Breeding vegetatively propagated horticultural crops. *Crop Breeding and  
1026 Applied Biotechnology*. 2011;**11**: 35-43
- 1027 20. Anushma, P., Dhanyasree, K. & Rafeekher, M. Wide hybridization for fruit crop  
1028 improvement: a review. *International Journal of Chemical Studies*. 2021;**9**: 769-73
- 1029 21. Bhat, J.A., Yu, D., Bohra, A. *et al.* Features and applications of haplotypes in crop breeding.  
1030 *Communications Biology*. 2021;**4**: 1266
- 1031 22. Chen, L., VanBuren, R., Paris, M. *et al.* The bracteatus pineapple genome and domestication  
1032 of clonally propagated crops. *Nature Genetics*. 2019;**51**: 1549-58
- 1033 23. Ming, R., VanBuren, R., Wai, C.M. *et al.* The pineapple genome and the evolution of CAM  
1034 photosynthesis. *Nature Genetics*. 2015;**47**: 1435-42
- 1035 24. He, Y., Luan, A., Wu, J. *et al.* Overcoming key technical challenges in the genetic  
1036 transformation of pineapple. *Tropical Plants*. 2023;**2**: 6

- 1037 25. Shu, H., Sun, W., Li, K. *et al.* The cause for water-heart fruit of pineapple and protective  
1038 measurements. *American Journal of Plant Sciences*. 2019;**10**: 885-92
- 1039 26. Sapak, Z. & Nusaibah, S.A. Common Diseases in Pineapple and Their Management in  
1040 Advances in Tropical Crop Protection (ed Wong, M.-Y.) 85-104 (Springer Nature  
1041 Switzerland, 2024).
- 1042 27. Feng, J., Zhang, W., Chen, C. *et al.* The pineapple reference genome: Telomere-to-telomere  
1043 assembly, manually curated annotation, and comparative analysis. *Journal of Integrative*  
1044 *Plant Biology*. 2024;**66**: 2208-25
- 1045 28. Ou, S., Chen, J. & Jiang, N. Assessing genome assembly quality using the LTR Assembly  
1046 Index (LAI). *Nucleic Acids Research*. 2018;**46**: e126
- 1047 29. Wang, Y., Tang, H., Wang, X. *et al.* Detection of colinear blocks and synteny and evolutionary  
1048 analyses based on utilization of MCScanX. *Nature Protocols*. 2024;**19**: 2206-29
- 1049 30. Yamashita, A., Izumi, N., Kashima, I. *et al.* SMG-8 and SMG-9, two novel subunits of the  
1050 SMG-1 complex, regulate remodeling of the mRNA surveillance complex during nonsense-  
1051 mediated mRNA decay. *Genes & Development*. 2009;**23**: 1091-105
- 1052 31. Jin, X.Y., Liu, Y.H., Hou, Z.M. *et al.* Genome-Wide Investigation of *SBT* Family Genes in  
1053 Pineapple and Functional Analysis of AcoSBT1.12 in Floral Transition. *Frontiers in Genetics*.  
1054 2021;**12**: 15
- 1055 32. Andersson, C.R., Helliwell, C.A., Bagnall, D.J. *et al.* The *FLX* Gene of *Arabidopsis* is  
1056 required for *FRI*-dependent activation of *FLC* expression. *Plant and Cell Physiology*.  
1057 2008;**49**: 191-200
- 1058 33. Luo, Y.Y., Liu, Z.Y., Jin, Z.X. *et al.* Phased T2T genome assemblies facilitate the mining of  
1059 disease-resistance genes in *Vitis davidii*. *Horticulture Research*. 2025;**12**: 13
- 1060 34. Wang, W.D., Chen, L.Y., Fengler, K. *et al.* A giant NLR gene confers broad-spectrum  
1061 resistance to *Phytophthora sojae* in soybean. *Nature Communications*. 2021;**12**: 8
- 1062 35. Wang, R., Hsu, Y., Bartholomew, D.P. *et al.* Delaying natural flowering in pineapple through  
1063 foliar application of aviglycine, an inhibitor of ethylene biosynthesis. *HortScience*. 2007;**42**:  
1064 1188-91
- 1065 36. Kuan, C., Yu, C., Lin, M. *et al.* Foliar application of aviglycine reduces natural flowering in  
1066 pineapple. *HortScience*. 2005;**40**: 123-6
- 1067 37. Sanewski, G., Ko, L., DeFaveri, J. *et al.* Genetic resistance to the root rot pathogen  
1068 *Phytophthora cinnamomi* in Ananas. *Acta Horti*. 2016;**1111**: 281-6
- 1069 38. Kumar, V. & Majee, A. Long noncoding RNAs in fruit crops. *Annu Rev Plant Biol*. 2021;**72**:  
1070 245-71
- 1071 39. Song, L., Fang, Y., Chen, L. *et al.* Role of non-coding RNAs in plant immunity. *Plant*  
1072 *Communications*. 2021;**2**: 100180
- 1073 40. Young, A.J., Pathania, N., Manners, A. *et al.* Heart rot of Australian pineapples caused by  
1074 *Dickeya zaeae*. *Australasian Plant Pathology*. 2022;**51**: 525-33
- 1075 41. Kuruppu, M., Siddiqui, Y. & Khalil, H.B. Comprehensive analysis of causal pathogens and  
1076 determinants influencing black rot disease development in MD2 pineapples. *Frontiers in*  
1077 *Microbiology*. 2025;**15**: 1514235
- 1078 42. Taagen, E., Bogdanove, A.J. & Sorrells, M.E. Counting on crossovers: controlled  
1079 recombination for plant breeding. *Trends in Plant Science*. 2020;**25**: 455-65
- 1080 43. Payero, L. & Alani, E. Crossover recombination between homologous chromosomes in  
1081 meiosis: recent progress and remaining mysteries. *Trends in Genetics*. 2025;**41**: 47-59
- 1082 44. Korunes, K.L. & Noor, M.A. Pervasive gene conversion in chromosomal inversion  
1083 heterozygotes. *Molecular Ecology*. 2019;**28**: 1302-15
- 1084 45. Berdan, E.L., Barton, N.H., Butlin, R. *et al.* How chromosomal inversions reorient the  
1085 evolutionary process. *Journal of Evolutionary Biology*. 2023;**36**: 1761-82
- 1086 46. Kulathinal, R.J., Bennett, S.M., Fitzpatrick, C.L. *et al.* Fine-scale mapping of recombination  
1087 rate in *Drosophila* refines its correlation to diversity and divergence. *Proceedings of the*  
1088 *National Academy of Sciences*. 2008;**105**: 10051-6
- 1089 47. Stevison, L.S., Hoehn, K.B. & Noor, M.A. Effects of inversions on within-and between-  
1090 species recombination and divergence. *Genome Biology and Evolution*. 2011;**3**: 830-41

- 1091 48. Termolino, P., Falque, M., Aiese Cigliano, R. *et al.* Recombination suppression in  
1092 heterozygotes for a pericentric inversion induces the interchromosomal effect on crossovers in  
1093 Arabidopsis. *The Plant Journal*. 2019;**100**: 1163-75
- 1094 49. Dixon, J.R., Selvaraj, S., Yue, F. *et al.* Topological domains in mammalian genomes identified  
1095 by analysis of chromatin interactions. *Nature*. 2012;**485**: 376-80
- 1096 50. Jiao, Y., Peluso, P., Shi, J. *et al.* Improved maize reference genome with single-molecule  
1097 technologies. *Nature*. 2017;**546**: 524-7
- 1098 51. Sinaga, A.O.Y. & Marpaung, D.S.S. Abiotic stress-induced gene expression in pineapple as a  
1099 potential genetic marker. *Advanced Agrochem*. 2024;**3**: 133-42
- 1100 52. Hasan, N., Choudhary, S., Naaz, N. *et al.* Recent advancements in molecular marker-assisted  
1101 selection and applications in plant breeding programmes. *Journal of Genetic Engineering and*  
1102 *Biotechnology*. 2021;**19**: 128
- 1103 53. Liu, J., Li, M.J., Zhang, Q. *et al.* Exploring the molecular basis of heterosis for plant breeding.  
1104 *Journal of Integrative Plant Biology*. 2020;**62**: 287-98
- 1105 54. Sun, X.P., Jiao, C., Schwaninger, H. *et al.* Phased diploid genome assemblies and pan-  
1106 genomes provide insights into the genetic history of apple domestication. *Nature Genetics*.  
1107 2020;**52**: 1423-32
- 1108 55. Li, Q.H., Qiao, X., Li, L.Q. *et al.* Haplotype-resolved T2T genome assemblies and  
1109 pangenome graph of pear reveal diverse patterns of allele-specific expression and the  
1110 genomic basis of fruit quality traits. *Plant Communications*. 2024;**5**: 21
- 1111 56. Yin, M.Q., Song, X.C., He, C. *et al.* The haplotype-resolved genome assembly of an ancient  
1112 citrus variety provides insights into the domestication history and fruit trait formation of  
1113 loose-skin mandarins. *Genome Biology*. 2025;**26**: 61
- 1114 57. Shao, L., Xing, F., Xu, C.H. *et al.* Patterns of genome-wide allele-specific expression in  
1115 hybrid rice and the implications on the genetic basis of heterosis. *Proceedings of the National*  
1116 *Academy of Sciences of the United States of America*. 2019;**116**: 5653-8
- 1117 58. Zhang, X.T., Chen, S., Shi, L.Q. *et al.* Haplotype-resolved genome assembly provides insights  
1118 into evolutionary history of the tea plant *Camellia sinensis*. *Nature Genetics*. 2021;**53**: 1250-9
- 1119 59. Usai, G., Giordani, T., Vangelisti, A. *et al.* Haplotype-resolved genome assembly of *Ficus*  
1120 *carica* L. reveals allele-specific expression in the fruit. *Plant Journal*. 2025;**121**: 16
- 1121 60. Zhan, W.M., Cui, L.H., Yang, S.L. *et al.* Natural variations of heterosis-related allele-specific  
1122 expression genes in promoter regions lead to allele-specific expression in maize. *BMC*  
1123 *Genomics*. 2024;**25**: 11
- 1124 61. Hou, Y.G., Gan, J.W., Fan, Z.Y. *et al.* Haplotype-based pangenomes reveal genetic variations  
1125 and climate adaptations in moso bamboo populations. *Nature Communications*. 2024;**15**: 15
- 1126 62. Quadrana, L., Etcheverry, M., Gilly, A. *et al.* Transposition favors the generation of large  
1127 effect mutations that may facilitate rapid adaption. *Nature Communications*. 2019;**10**: 10
- 1128 63. Zhang, L., Shi, Y., Gong, W.F. *et al.* The tetraploid *Camellia oleifera* genome provides  
1129 insights into evolution, agronomic traits, and genetic architecture of oil *Camellia* plants. *Cell*  
1130 *Reports*. 2024;**43**: 24
- 1131 64. Lan, L., Leng, L.H., Liu, W.C. *et al.* The haplotype-resolved telomere-to-telomere carnation  
1132 (*Dianthus caryophyllus*) genome reveals the correlation between genome architecture and  
1133 gene expression. *Horticulture Research*. 2024;**11**: 13
- 1134 65. Zhang, L.Y., Hu, J., Han, X.L. *et al.* A high-quality apple genome assembly reveals the  
1135 association of a retrotransposon and red fruit colour. *Nature Communications*. 2019;**10**: 13
- 1136 66. Ngou, B.P.M., Heal, R., Wyler, M. *et al.* Concerted expansion and contraction of immune  
1137 receptor gene repertoires in plant genomes. *Nature Plants*. 2022;**8**: 1146-52
- 1138 67. Jones, J.D.G., Vance, R.E. & Dangl, J.L. Intracellular innate immune surveillance devices in  
1139 plants and animals. *Science*. 2016;**354**: 6
- 1140 68. Adachi, H., Derevnina, L. & Kamoun, S. NLR singletons, pairs, and networks: evolution,  
1141 assembly, and regulation of the intracellular immunoreceptor circuitry of plants. *Current*  
1142 *Opinion in Plant Biology*. 2019;**50**: 121-31
- 1143 69. Chen, S.H., Martino, A.M., Luo, Z.Y. *et al.* A high-quality pseudo-phased genome for  
1144 *Melaleuca quinquenervia* shows allelic diversity of NLR-type resistance genes. *Gigascience*.  
1145 2023;**12**: 18

- 1146 70. Tian, D., Traw, M.B., Chen, J.Q. *et al.* Fitness costs of R-gene-mediated resistance in  
1147 *Arabidopsis thaliana*. *Nature*. 2003;**423**: 74-7
- 1148 71. Furtado, A. DNA extraction from vegetative tissue for next-generation sequencing in Cereal  
1149 Genomics: Methods and Protocols 2013/11/19 edn, Vol. 1099 (eds Henry, R.J. & Furtado,  
1150 A.) 1-5 (Humana Press, 2014).
- 1151 72. Rao, S.S., Huntley, M.H., Durand, N.C. *et al.* A 3D map of the human genome at kilobase  
1152 resolution reveals principles of chromatin looping. *Cell*. 2014;**159**: 1665-80
- 1153 73. Marçais, G. & Kingsford, C. A fast, lock-free approach for efficient parallel counting of  
1154 occurrences of k-mers. *Bioinformatics*. 2011;**27**: 764-70
- 1155 74. Ranallo-Benavidez, T.R., Jaron, K.S. & Schatz, M.C. GenomeScope 2.0 and Smudgeplot for  
1156 reference-free profiling of polyploid genomes. *Nature Communications*. 2020;**11**: 1432
- 1157 75. Zeng, X., Yi, Z., Zhang, X. *et al.* Chromosome-level scaffolding of haplotype-resolved  
1158 assemblies using Hi-C data without reference genomes. *Nature Plants*. 2024;**10**: 1184–200
- 1159 76. Durand, N.C., Robinson, J.T., Shamim, M.S. *et al.* Juicebox Provides a Visualization System  
1160 for Hi-C Contact Maps with Unlimited Zoom. *Cell Systems*. 2016;**3**: 99-101
- 1161 77. Lin, Y., Ye, C., Li, X. *et al.* quarTeT: a telomere-to-telomere toolkit for gap-free genome  
1162 assembly and centromeric repeat identification. *Horticulture Research*. 2023;**10**: uhad127
- 1163 78. Altschul, S.F., Gish, W., Miller, W. *et al.* Basic local alignment search tool. *Journal of*  
1164 *Molecular Biology*. 1990;**215**: 403-10
- 1165 79. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*. 2018;**34**:  
1166 3094-100
- 1167 80. Gurevich, A., Saveliev, V., Vyahhi, N. *et al.* QUAST: quality assessment tool for genome  
1168 assemblies. *Bioinformatics*. 2013;**29**: 1072-5
- 1169 81. Tegenfeldt, F., Kuznetsov, D., Manni, M. *et al.* OrthoDB and BUSCO update: annotation of  
1170 orthologs with wider sampling of genomes. *Nucleic Acids Research*. 2025;**53**: D516-D22
- 1171 82. Li, H. Protein-to-genome alignment with minimap2. *Bioinformatics*. 2023;**39**: btad014
- 1172 83. Levy Karin, E., Mirdita, M. & Söding, J. MetaEuk—sensitive, high-throughput gene  
1173 discovery, and annotation for large-scale eukaryotic metagenomics. *Microbiome*. 2020;**8**: 48
- 1174 84. Rhie, A., Walenz, B.P., Koren, S. *et al.* Merqury: reference-free quality, completeness, and  
1175 phasing assessment for genome assemblies. *Genome Biology*. 2020;**21**: 245
- 1176 85. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform.  
1177 *Bioinformatics*. 2009;**25**: 1754-60
- 1178 86. Kim, D., Paggi, J.M., Park, C. *et al.* Graph-based genome alignment and genotyping with  
1179 HISAT2 and HISAT-genotype. *Nature Biotechnology*. 2019;**37**: 907-15
- 1180 87. Li, H., Handsaker, B., Wysoker, A. *et al.* The sequence alignment/map format and SAMtools.  
1181 *Bioinformatics*. 2009;**25**: 2078-9
- 1182 88. Flynn, J.M., Hubley, R., Goubert, C. *et al.* RepeatModeler2 for automated genomic discovery  
1183 of transposable element families. *Proceedings of the National Academy of Sciences*.  
1184 2020;**117**: 9451-7
- 1185 89. Tarailo-Graovac, M. & Chen, N. Using RepeatMasker to identify repetitive elements in  
1186 genomic sequences. *Current Protocols in Bioinformatics*. 2009;**25**: 4.10.1-4.4
- 1187 90. Nakandala, U., Masouleh, A.K., Smith, M.W. *et al.* Haplotype resolved chromosome level  
1188 genome assembly of *Citrus australis* reveals disease resistance and other citrus specific genes.  
1189 *Horticulture Research*. 2023;**10**: uhad058
- 1190 91. Gabriel, L., Brúna, T., Hoff, K.J. *et al.* BRAKER3: Fully automated genome annotation using  
1191 RNA-seq and protein evidence with GeneMark-ETP, AUGUSTUS, and TSEBRA. *Genome*  
1192 *Research*. 2024;**34**: 769-77
- 1193 92. Chan, P.P., Lin, B.Y., Mak, A.J. *et al.* tRNAscan-SE 2.0: improved detection and functional  
1194 classification of transfer RNA genes. *Nucleic Acids Research*. 2021;**49**: 9077-96
- 1195 93. Goel, M., Sun, H., Jiao, W.B. *et al.* SyRI: finding genomic rearrangements and local sequence  
1196 differences from whole-genome assemblies. *Genome Biology*. 2019;**20**: 277
- 1197 94. DePristo, M.A., Banks, E., Poplin, R. *et al.* A framework for variation discovery and  
1198 genotyping using next-generation DNA sequencing data. *Nature Genetics*. 2011;**43**: 491-8

- 1199 95. Cingolani, P., Platts, A., Wang, L.L. *et al.* A program for annotating and predicting the effects  
1200 of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila*  
1201 *melanogaster* strain *w<sup>1118</sup>*; *iso-2*; *iso-3*. *Fly*. 2012;**6**: 80-92
- 1202 96. Nattestad, M. & Schatz, M.C. Assemblytics: a web analytics tool for the detection of variants  
1203 from an assembly. *Bioinformatics*. 2016;**32**: 3021-3
- 1204 97. Marçais, G., Delcher, A.L., Phillippy, A.M. *et al.* MUMmer4: A fast and versatile genome  
1205 alignment system. *Plos Computational Biology*. 2018;**14**: 14
- 1206 98. Darling, A.E., Mau, B. & Perna, N.T. progressiveMauve: multiple genome alignment with  
1207 gene gain, loss and rearrangement. *PloS One*. 2010;**5**: e11147
- 1208 99. Wu, F., Mai, Y., Chen, C. *et al.* SynGAP: a synteny-based toolkit for gene structure annotation  
1209 polishing. *Genome Biology*. 2024;**25**: 218
- 1210 100. Chen, C., Chen, H., Zhang, Y. *et al.* TBtools: An Integrative Toolkit Developed for Interactive  
1211 Analyses of Big Biological Data. *Molecular Plant*. 2020;**13**: 1194-202
- 1212 101. Krzywinski, M., Schein, J., Birol, I. *et al.* Circos: an information aesthetic for comparative  
1213 genomics. *Genome Research*. 2009;**19**: 1639-45
- 1214 102. Ou, S.J., Su, W.J., Liao, Y. *et al.* Benchmarking transposable element annotation methods for  
1215 creation of a streamlined, comprehensive pipeline. *Genome Biology*. 2019;**20**: 18
- 1216 103. Kilian, A., Wenzl, P., Huttner, E. *et al.* Diversity arrays technology: a generic genome  
1217 profiling technology on open platforms. *Methods Mol Biol*. 2012;**888**: 67-89
- 1218 104. Kilian, A., Sanewski, G. & Ko, L. The application of DArTseq technology to pineapple. *Acta*  
1219 *Horticulturae*. 2016;**1111**: 181-8
- 1220 105. Bradbury, P.J., Zhang, Z., Kroon, D.E. *et al.* TASSEL: software for association mapping of  
1221 complex traits in diverse samples. *Bioinformatics*. 2007;**23**: 2633-5
- 1222 106. Money, D., Gardner, K., Migicovsky, Z. *et al.* LinkImpute: fast and accurate genotype  
1223 imputation for nonmodel organisms. *G3: Genes, Genomes, Genetics*. 2015;**5**: 2383-90
- 1224 107. Zhang, R., Wu, H., Li, Y. *et al.* GWLD: an R package for genome-wide linkage  
1225 disequilibrium analysis. *G3: Genes, Genomes, Genetics*. 2023;**13**: jkad154
- 1226 108. Broman, K.W., Gatti, D.M., Simecek, P. *et al.* R/qtl2: software for mapping quantitative trait  
1227 loci with high-dimensional data and multiparent populations. *Genetics*. 2019;**211**: 495-502
- 1228 109. Li, P.C., Quan, X.D., Jia, G.F. *et al.* RGAugury: a pipeline for genome-wide prediction of  
1229 resistance gene analogs (RGAs) in plants. *BMC Genomics*. 2016;**17**: 852
- 1230 110. Hao, Z.D., Lv, D.K., Ge, Y. *et al.* RIdiogram: drawing SVG graphics to visualize and map  
1231 genome-wide data on the idiograms. *PeerJ Computer Science*. 2020;**6**: e251
- 1232 111. Edgar, R.C. MUSCLE: multiple sequence alignment with high accuracy and high throughput.  
1233 *Nucleic Acids Research*. 2004;**32**: 1792-7
- 1234 112. Nguyen, L.T., Schmidt, H.A., von Haeseler, A. *et al.* IQ-TREE: A Fast and Effective  
1235 Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Molecular Biology*  
1236 *and Evolution*. 2015;**32**: 268-74
- 1237 113. Quinlan, A.R. & Hall, I.M. BEDTools: a flexible suite of utilities for comparing genomic  
1238 features. *Bioinformatics*. 2010;**26**: 841-2