# The genome and population genomics of allopolyploid *Coffea arabica* reveal the diversification history of modern coffee cultivars

Jarkko Salojärvi[1,2,3,*], Aditi Rambani[4,†], Zhe Yu[5,†], Romain Guyot[6,7,†], Susan Strickler[4,†], Maud Lepelley[8], Cui Wang[2], Sitaram Rajaraman[2], Pasi Rastas[9], Chunfang Zheng[5], Daniella Santos Muñoz[5], João Meidanis[10], Alexandre Rossi Paschoal[11], Yves Bawin[12], Trevor Krabbenhoft[13], Zhen Qin Wang[13], Steven Fleck[13], Rudy Aussel[8,14], Laurence Bellanger[8], Aline Charpagne[15], Coralie Fournier[15], Mohamed Kassam[15], Gregory Lefebvre[15], Sylviane Métairon[15], Déborah Moine[15], Michel Rigoreau[8], Jens Stolte[15], Perla Hamon[6], Emmanuel Couturon[6], Christine Tranchant-Dubreuil[6], Minakshi Mukherjee[13], Tianying Lan[13], Jan Engelhardt[16], Peter Stadler[17], Samara Mireza Correia De Lemos[18], Suzana Ivamoto Suzuki[19], Ucu Sumirat[20], Wai Ching Man[21], Nicolas Dauchot[22], Simon Orozco-Arias[7], Andrea Garavito[23], Catherine Kiwuka[24], Pascal Musoli[24], Anne Nalukenge[24], Erwan Guichoux[25], Havinga Reinout[26], Martin Smit[26], Lorenzo Carretero-Paulet[27], Oliveiro Guerreiro Filho[28], Masako Toma Braghini[28], Lilian Padilha[29], Gustavo Hiroshi Sera[30], Tom Ruttink[12,33], Robert Henry[31], Pierre Marraccini[32], Yves Van de Peer[33,34,35,40], Alan Andrade[36], Douglas Domingues[18], Giovanni Giuliano[37], Lukas Mueller[4], Luiz Filipe Pereira[38], Stephane Plaisance[39], Valerie Poncet[6], Stephane Rombauts[33,40], David Sankoff[5], Victor A. Albert[13,*], Dominique Crouzillat[8,*], Alexandre de Kochko[6,*], Patrick Descombes[15,*]

[1] School of Biological Sciences, Nanyang Technological University, Singapore 637551, Singapore
[2] Organismal and Evolutionary Biology Research Programme, University of Helsinki, 00014 Helsinki, Finland
[3] Singapore Centre for Environmental Life Sciences Engineering, Nanyang Technological University, Singapore 637551, Singapore
[4] Boyce Thompson Institute, University of Cornell, Ithaca NY 14853, US
[5] Department of Mathematics and Statistics, University of Ottawa, Ottawa, Canada K1N 6N5
[6] Institut de Recherche pour le Développement (IRD), Université de Montpellier, 34394 Montpellier, France
[7] Department of Electronics and Automation, Universidad Autónoma de Manizales, Manizales 170002, Colombia
[8] Société des Produits Nestlé SA, Nestlé Research, 37097 Tours CEDEX 2, France
[9] Institute of Biotechnology, University of Helsinki, 00014 Helsinki, Finland
[10] Institute of Computing, University of Campinas, 13083-852 Campinas, Sao Paulo, Brazil
[11] Department of Computer Science, The Federal University of Technology – Paraná (UTFPR), 86300-000, Cornélio Procópio, Brazil
[12] Plant Sciences Unit, Flanders research Institute for Agriculture, Fisheries and Food (ILVO), 9090 Melle, Belgium
[13] Department of Biological Sciences, University at Buffalo, New York, USA
[14] Centre d'Immunologie de Marseille-Luminy, Aix Marseille Université, France
[15] Société des Produits Nestlé SA, Nestlé Research, 1015 Lausanne, Switzerland
[16] Department of Computer Science, University of Leipzig, 04107 Leipzig, Germany
[17] Department of Computer Science and Interdisciplinary Center for Bioinformatics, University of Leipzig, 04107 Leipzig, Germany
[18] Group of Genomics and Transcriptomes in Plants, São Paulo State University, UNESP, Rio Claro, SP, Brazil, 13506-900
[19] Centro de Ciências Agrárias, Universidade Estadual de Londrina, 86057-970 Londrina, Brazil
[20] Indonesian Coffee and Cocoa Research Institute (ICCRI), Jember 68118 Indonesia
[21] Texas A&M University, Urbana, Illinois 61801, USA
[22] Research Unit in Plant Cellular and Molecular Biology, University of Namur, Namur 5000, Belgium
[23] Departamento de Ciencias biológicas, Facultad de Ciencias Exactas y Naturales, Universidad de Caldas, Manizales, Colombia

1

56    [24] National Agricultural Research Organization (NARO), Uganda
57    [25] Biodiversité Gènes & Communautés, INRA, 33610 CESTAS, France
58    [26] Hortus Botanicus Amsterdam, 1018 DD Amsterdam, Netherlands
59    [27] Departamento de Biología y Geología, Universidad de Almería, Almería, Spain
60    [28] Instituto Agronômico (IAC) Centro de Café 'Alcides Carvalho', Fazenda Santa Elisa, Caixa
61    Postal 28, Campinas (SP), Brasil 13012 – 970
62    [29] Embrapa Café / Instituto Agronômico (IAC) Centro de Café 'Alcides Carvalho', Fazenda
63    Santa Elisa, Caixa Postal 28, Campinas (SP), Brasil 13012 - 970
64    [30] Instituto de Desenvolvimento Rural do Paraná- IAPAR, 86047-902 Londrina, Brasil
65    [31] Queensland Alliance for Agriculture and Food Innovation, University of Queensland,
66    Brisbane 4072, Australia
67    [32] CIRAD - UMR DIADE (IRD-CIRAD-Université de Montpellier) BP 64501, F-34394
68    Montpellier Cedex 5, France
69    [33] Department of Plant Biotechnology and Bioinformatics, Ghent University
70    [34] Department of Biochemistry, Genetics and Microbiology, University of Pretoria, Pretoria
71    0028, South Africa
72    [35] College of Horticulture, Academy for Advanced Interdisciplinary Studies, Nanjing
73    Agricultural University, Nanjing, 210095, China
74    [36] Embrapa Café/Inovacafé Laboratory of Molecular Genetics Campus da UFLA-MG, 37200-
75    900 Lavras-MG, Brazil
76    [37] Italian National Agency for New technologies, Energy and Sustainable Economic
77    Development (ENEA), Casaccia Res. Ctr., 00123 Roma, Italy
78    [38] Embrapa Café / Lab. Biotecnologia, Área de Melhoramento Genético, Londrina – PR, Brasil
79    - 86047-902
80    [39] VIB Nucleomics Core, B-3000 Leuven, Belgium
81    [40] Center for Plant Systems Biology, VIB, 9052 Ghent, Belgium
82
83
84    *Correspondence should be addressed to: Patrick Descombes
85    (patrick.descombes@rd.nestle.com), Alexandre de Kochko (alexandre.dekochko@ird.fr),
86    Jarkko Salojärvi (jarkko@ntu.edu.sg), Victor A. Albert (vaalbert@buffalo.edu), or Dominique
87    Crouzillat (dcrouzillat@gmail.com).
88
89    [†]These authors contributed equally to this work
90

91    **Abstract**

92    *Coffea arabica,* an allotetraploid hybrid of *C. eugenioides* and *C. canephora,* is the source of
93    approximately 60% of coffee products worldwide, and its cultivated accessions have undergone
94    several population bottlenecks. We present chromosome-level assemblies of a di-haploid *C.*
95    *arabica* accession and modern representatives of its diploid progenitors, *C. eugenioides* and *C.*
96    *canephora*. The three species exhibit largely conserved genome structures between diploid
97    parents and descendant subgenomes, with no obvious global subgenome dominance. We find
98    evidence for a founding polyploidy event 350,000-610,000 years ago, followed by several pre-
99    domestication bottlenecks, resulting in narrow genetic variation. A split between wild
100   accessions and cultivar progenitors occurred ~30.5 kya, followed by a period of migration
101   between the two populations. Analysis of modern varieties, including lines historically
102   introgressed with *C. canephora,* highlights their breeding histories and loci that may contribute
103   to pathogen resistance, laying the groundwork for future genomics-based breeding of *C.*
104   *arabica*.
105

106   **Introduction**

107   Polyploidy is a powerful evolutionary force that has shaped genome evolution across many
108   eukaryotic lineages, possibly offering adaptive advantages in times of global change[1,2]. Such

2

109  whole genome duplications (WGDs) are particularly characteristic of plants[3], and a great
110  proportion of crop species are polyploid[4-11]. Our understanding of genome evolution following
111  WGD is still incomplete, but outcomes can result in genomic shock, in terms of activation of
112  cryptic transposable elements, subgenome-partitioned gene regulation or fractionation,
113  homoeologous exchange, meiotic instability, and even karyotype variation [8, 12-16]. Alternatively,
114  few or none of the above phenomena can materialize, and the two genomes can coexist
115  harmonically, gradually adapting to new ploidy levels[17]. Regardless, the most common fate
116  of polyploids appears to be fractionation and eventual reversion to the diploid state[18].
117
118  With an estimated production of 10 million metric tons per year, coffee is one of the most traded
119  commodities in the world.  The most broadly appreciated coffee is produced from the
120  allotetraploid species *Coffea arabica*, especially from cultivars belonging to the Bourbon or
121  Typica lineages and their hybrids[19]. *C. arabica* (2n = 4x = 44 chromosomes) resulted from a
122  natural hybridization event between the ancestors of present-day *C. canephora* (Robusta coffee,
123  subgenome CC) and *C. eugenioides* (subgenome EE; each with 2n = 2x = 22). The founding
124  WGD has been dated to between 10,000 to one million years ago[20-23], with the Robusta-derived
125  subgenome of *C. arabica* closest related to *C. canephora* accessions from northern Uganda[24].
126  Arabica cultivation was initiated in 15th -16th century Yemen. (**Ext. data Fig. 1**). Around 1600,
127  "seven seeds" were smuggled out of Yemen[25], establishing Indian *C. arabica* cultivar lineages.
128  A century later, the Dutch began cultivating Arabica in Southeast Asia – thus setting up the
129  founders of the contemporary Typica group. One plant, shipped to Amsterdam in 1706, was
130  used to establish Arabica cultivation in the Caribbean in 1723. Independently, the French
131  cultivated Arabica on the island of Bourbon (presently Réunion)[26], and the descendants of a
132  single plant that survived by 1720 form the contemporary Bourbon group. Contemporary
133  Arabica cultivars descend from these Typica or Bourbon lineages, except for a few wild
134  ecotypes with origins in natural forests in Ethiopia. Due to its recent allotetraploid origin and
135  strong bottlenecks during its history, cultivated *C. arabica* harbors a particularly low genetic
136  diversity[20] and is susceptible to many plant pests and diseases, such as coffee leaf rust (*Hemileia*
137  *vastatrix*). As a result, the classic Bourbon-Typica lineages can only be cultivated successfully
138  in a few regions around the world. Fortunately, a spontaneous *C. canephora* x *C. arabica* hybrid
139  resistant to *H. vastatrix* was identified on the island of Timor[27] in 1927. Many modern Arabicas
140  contain *C. canephora* introgressions derived from this hybrid, ensuring rust resistance, but
141  having also unwanted side effects, such as decreased beverage quality[28].
142
143  Modern genomic tools and a detailed understanding of the origin and breeding history of
144  contemporary varieties are vital to developing new Arabica cultivars, better adapted to climate
145  change and agricultural practices[29-31]. Here, we present chromosome-level assemblies of *C.*
146  *arabica* and representatives of its progenitor species, *C. canephora* (Robusta) and *C.*
147  *eugenioides* (hereafter Eugenioides). Whole-genome resequencing data of 41 wild and
148  cultivated accessions facilitated in-depth analysis of Arabica history and dissemination routes,
149  as well as the identification of candidate genomic regions associated with pathogen resistance.
150
151  **Results and Discussion**

152  **Chromosome-level assemblies and annotations, genome fractionation, and subgenome**
153  **dominance**
154  As reference individuals, we chose the di-haploid Arabica line ET-39[32], a previously sequenced
155  doubled haploid Robusta[33], and the wild Eugenioides accession Bu-A, respectively. Long and
156  short-read-based hybrid assemblies were obtained (Online methods and **Supplementary**
157  **sections 2.1-2.2**), spanning 672 Mb (Robusta), 645 Mb (Eugenioides) and 1,088 Mb (Arabica),
158  respectively. Upon HiC scaffolding, the Robusta and Arabica assemblies consisted of 11 and
159  22 pseudochromosomes, and spanned 82.7% and 62.5%, respectively, of the projected genome
160  sizes (**Table 1**). To improve the Arabica assembly, we  generated a second assembly using
161  PacBio HiFi technology followed by Hi-C scaffolding (Online methods and **Supplementary**
162  **section 2.3**). This assembly was 1,198 Mb long, of which 1,192 Mb (93.1% of the predicted

163 genome size based on cytological evidence[34]) was anchored to pseudochromosomes (**Table 1**).
164 Gene space completeness, assessed using Benchmarking Universal Single-Copy Orthologs
165 (BUSCO)[35], was >96% for all assemblies. Importantly, 93.2% of the BUSCO genes were
166 duplicated in the HiFi assembly (**Table 1**), indicating that most of the gene duplicates from the
167 allopolyploidy event were retained.
168
169 The Robusta and Eugenioides genomes contained, respectively, 67.5% and 59.7% transposable
170 elements (TEs) (**Supplementary section 3.2**), with Gypsy LTR retrotransposons accounting
171 for most of the difference between the two species. This difference was greatly reduced (63.1%
172 and 63.8%) in the two Arabica subgenomes (subCC and subEE, stemming from Robusta and
173 Eugenioides ancestors, respectively), possibly indicating TE transfer via homoeologous
174 exchange. Robusta contained considerably more recent LTR TE insertion elements than
175 Eugenioides. Again, the two Arabica subgenomes showed greater similarity to each other in
176 recent LTR TE insertions than the two progenitor genomes. No major evidence was found for
177 LTR TE mobilization following Arabica allopolyploidization, in contrast to what has been
178 observed in tobacco[36], but similar to *Brassica* synthetic allotetraploids[37]. Arabica genome
179 evolution observed instead more closely follows the "harmonious coexistence" pattern[38] seen
180 in *Arabidopsis* hybrids[17,39].
181
182 High-quality gene annotations, followed by manual curation of specific gene families
183 (**Supplementary Sections 3.1-3.4**), resulted in 28,857, 33,505, 56,670, and 69,314 gene
184 models for the Robusta, Eugenioides, PacBio Arabica, and Arabica HiFi assemblies,
185 respectively (**Table 1**). Altogether ~97% of Robusta and 99.6% of Arabica HiFi gene models
186 were placed on the pseudochromosomes, with 33,618 and 35,449, respectively, to subgenomes
187 subCC and subEE (**Table 1**). Annotation completeness from BUSCO was ≥95% for
188 Eugenioides and Robusta, and reached 97.3% for Arabica HiFi.
189
190 Comparison of Arabica subCC and subEE against their Robusta and Eugenioides counterparts
191 revealed high conservation in terms of chromosome number, centromere position and numbers
192 of genes per chromosome (**Fig. 1, Supplementary section 4**). Patterns of gene loss following
193 the *gamma* paleohexaploidy event displayed high structural conservation between Robusta and
194 Eugenioides during the 4-6 million years (My) since their initial species split[22,23]
195 (**Supplementary section 4**). Likewise, the structures of the two Arabica subgenomes were
196 highly conserved between each other, with, since the Arabica-founding allotetraploidy event,
197 only ~5% of BUSCO genes having reverted to the diploid state (**Fig. 1A; Table 1**). Syntenic
198 comparisons revealed that genomic excision events, removing one or several genes at a time in
199 similar proportions across the two subgenomes, have been the main driving force in genome
200 fragmentation both before and after the polyploidy event (**Fig. 1B, Supplementary section 4**).
201 Fractionation occurred mostly in pericentromeric regions, whereas chromosome arms showed
202 more moderate paralogous gene deletion (**Fig 1C, Supplementary section 4**). The Arabica
203 allopolyploidy event seemingly did not affect the rate of genome fractionation, which remained
204 roughly constant when comparing deletions in progenitor species versus Arabica subgenomes
205 after the event. In support of the dosage-balance hypothesis[40], subgenomic regions with high
206 duplicate retention rates were significantly enriched for genes that originated from the Arabica
207 WGD (Fisher exact test, $p<2.2e-16$). In contrast, low duplicate retention rate regions
208 significantly overlapped with genes originating from small-scale (tandem) duplications (**Table
209 S1**). Genes with high retention rates were enriched in Gene Ontology (GO) categories such as
210 "cellular component organization or biogenesis", "primary metabolic process", "developmental
211 process" and "regulation of cellular process", while low retention rate genes were enriched in
212 categories such as "RNA-dependent DNA biosynthetic process" and "defense response" (in
213 both subgenomes) and "spermidine hydroxycinnamate conjugate biosynthetic process"
214 (involved in plant defense[41]) and "plant-type hypersensitive response" (in subEE) (**Tables S2-
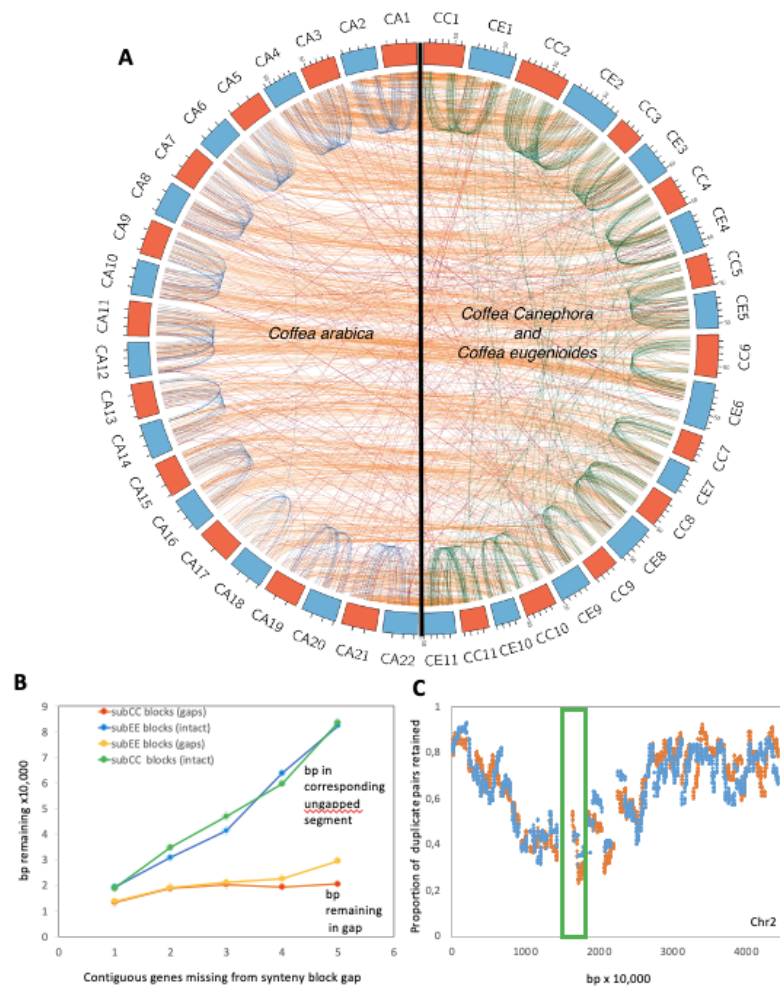215 S5**).
216
217

218



219
220
221 ***Figure 1.*** *Patterns of synteny, fractionation and gene loss in* Coffea arabica *(CA) and its*
222 *progenitor species* C. canephora *(CC) and* C. eugenioides *(CE).* ***A.*** *Corresponding syntenic*
223 *blocks between CA subgenomes subCC (orange) and subEE (blue), and with the CC (orange)*
224 *and CE (blue) genomes.* ***B.*** *bp in intergenic DNA in synteny block gaps caused by fractionation*
225 *in a subCC-subEE comparison, compared to numbers of bp in homoeologous unfractionated*
226 *regions, as a function of numbers of consecutive genes deleted.* ***C.*** *Gene retention rates in*
227 *synteny blocks plotted along subCC chromosome 2; subCC is plotted in orange and subEE in*
228 *blue. The green box indicates the pericentromeric region.*

229
230 To study possible expression biases between subgenomes, we identified syntelogous gene pairs
231 and removed the pairs showing homoeologous exchanges in the Arabica subgenomes (see
232 under **Origin and domestication of Arabica coffee**, below)[42] (**Supplementary section 5**).
233 Overall, no significant global subgenome expression dominance was observed (**Tables S6-S7**).
234 However, gene families regularly displayed mosaic patterns of expression, including several
235 encoding enzymes that contribute to cup quality, such as *N*-methyltransferase (*NMT*), terpene
236 synthase (*TPS*), and fatty acid desaturase 2 (*FAD2*) families, all having some genes being more
237 expressed in one of the two subgenomes (**Ext. data Fig. 2**), as per a recent study[43]. Similar
238 gene family-wise patterns occur in other evolutionarily recent polyploids such as rapeseed[10]
239 and cotton[44], which are also at their early stages of transitioning back to a diploid state.

240
241
242

5

243
244
245 **Origin and domestication of Arabica coffee**
246 To obtain a genomic perspective on the evolutionary history of Arabica, we sequenced 46
247 accessions, including three Robusta, two Eugenioides, and 41 Arabica. The latter included an
248 18th-century type specimen, kindly provided by the Linnaean Society of London, 12 cultivars
249 with different breeding histories, the Timor hybrid and five of its backcrosses to Arabica, and
250 17 wild and three wild/cultivated accessions collected from the Eastern and Western sides of
251 the Great Rift Valley[45,46] (**Table S8, Fig. 2A**).
252
253 Homoeologous exchange (HE) between subgenomes has been observed in several recent
254 polyploids[8,42,47]. Arabica generally displays bivalent pairing of homologous chromosomes and
255 disomic inheritance[48], but since the subgenomes share high similarity, occasional
256 homoeologous pairing and exchange may also occur. We therefore explored the extent of HE
257 among Arabica accessions and its possible contribution to genome evolution. Overall, all
258 accessions shared a fixed allele bias toward subEE at one end of chromosome 7, which
259 contained genes enriched for chloroplast-associated functions (**Ext. data Fig. 3A,**
260 **Supplementary section 5**, **Table S9**). Since the Arabica plastid genome is derived from
261 Eugenioides[49], HE in this region was likely selected for, due to compatibility issues between
262 nuclear and chloroplast genes encoding chloroplast-localized proteins[50]. Surprisingly, all but
263 one accession (BMJM) showed significant ($p$-value $<9.8e-37$) 3:1 allelic biases towards subCC.
264 The highly concordant HE patterns, present in both wild and cultivated Arabicas (**Ext. data**
265 **Fig. 4)**, suggested that i) the allelic bias is an adaptive trait not associated with breeding, and
266 ii) it originated in a common ancestor of all sampled accessions, possibly immediately after the
267 founding allopolyploidy event. Some exchanges, shared by only a few accessions, probably
268 originated more recently (**Ext. data Fig. 3B**). More recent HE events were also found in some
269 cultivars and also showed a bias towards subCC, except for BMJM, which showed bias towards
270 subEE due to a single large crossover in chromosome 1 (**Ext. data Fig. 3A**). An interesting
271 hypothesis for future investigation is that in a low-diversity polyploid species such as Arabica,
272 HE could be a major contributor to phenotypic variation observed among closely related
273 accessions[51].
274
275 We next studied population genetic statistics for each of the subgenomes (**Table S10**). The 17
276 wild samples demonstrated low genomic diversities, indicative of small effective population
277 sizes, while negative Tajima's D suggested an expanding population, possibly following one
278 or more population bottlenecks. The cultivars and wild population samples had similar genetic
279 diversities, as demonstrated by low fixation index ($F_{ST}$) values. In cultivars, nucleotide
280 diversities were only slightly lower than in wild populations and Tajima's D scores were less
281 negative, suggesting that only minor bottlenecks and subsequent population expansions
282 occurred during domestication.
283
284 SNP tree estimation and ADMIXTURE analyses (**Fig. 2B**) identified a three-population
285 solution for subCC: Typica-Bourbon cultivars (Population 1), wild accessions (Population 2),
286 and Timor hybrid-derived cultivars (Population 3). The old BMJM and the recently established
287 Geisha cultivars showed admixed states on both subgenomes, similar to about half of the wild
288 accessions. Indian varieties encompassed both Typica and Bourbon variation, in agreement
289 with previous studies[20]. The Linnaean sample grouped with the cultivars, supporting its
290 hypothesized origin from the Dutch East Indies[25]. A complementary analysis using PCA (**Ext.**
291 **data Fig. 5)** was in agreement with ADMIXTURE analysis.
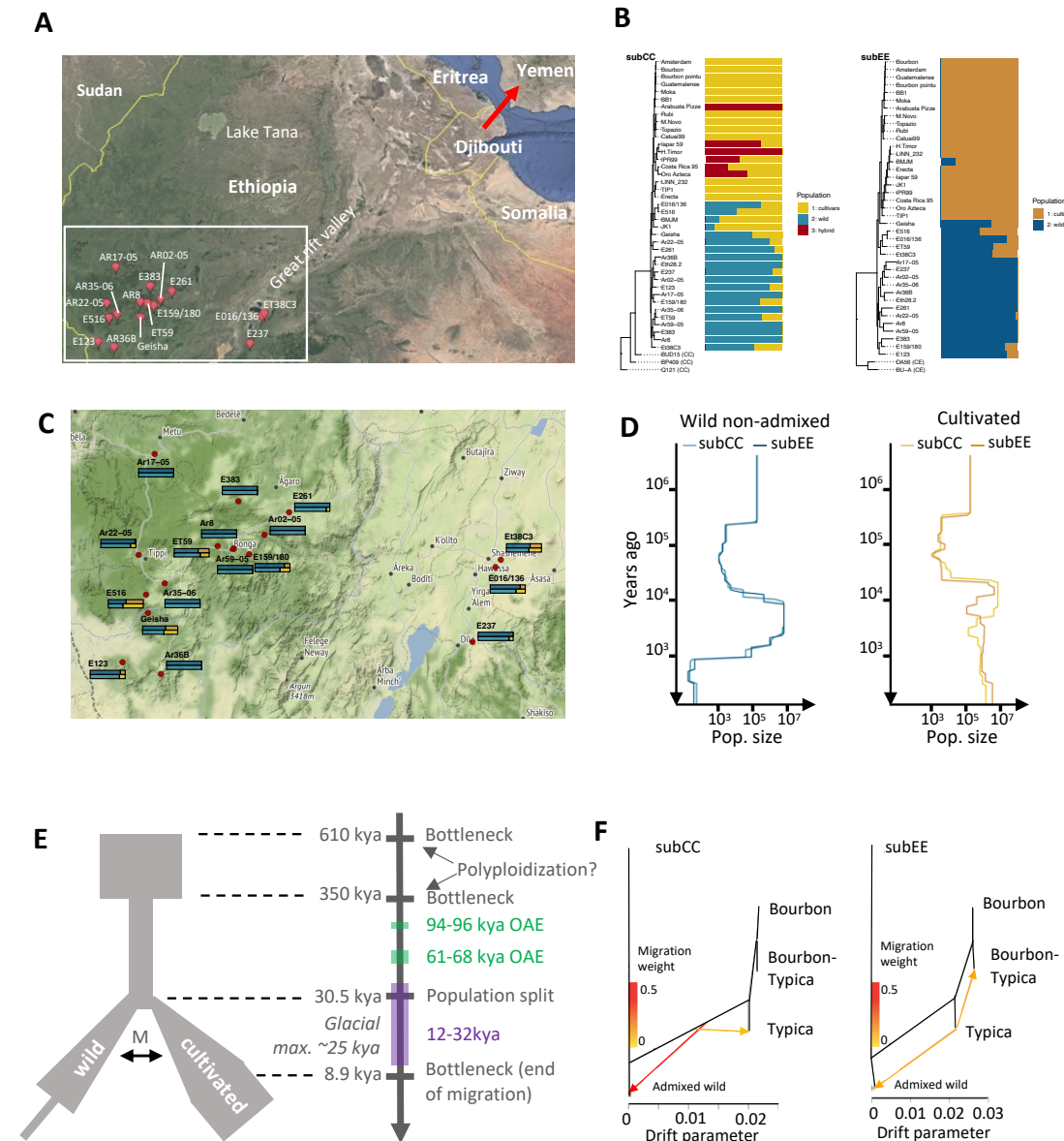292

293
294 **Figure 2.** *Population history of* Coffea arabica. *A. Geographic origin of resequenced wild* C.
295 arabica *accessions. The red arrow indicates the probable route of migration to Yemen in*
296 *historical times.* **B.** *Ancestral population assignments of* C. arabica *accessions for subCC (left)*
297 *and subEE (right). Relationships among individuals are illustrated with phylogenetic trees*
298 *obtained from independent SNPs. For magnified views of the trees, see* **Fig. S45**. *C.*
299 *Magnification of panel A, showing the admixture values for each of the accessions in subCC*
300 *(top) and subEE (bottom); the colors correspond to the analysis in panel B.* **D.** *Population sizes*
301 *of wild and cultivated accessions, inferred using SMC++, suggest genetic bottlenecks at ~350*
302 *and 1 kya (limited to non-admixed wild individuals).* **E.** *FastSimcoal2 output, suggesting a*
303 *population split ~30.5 kya, followed by a period of migration between the populations until*
304 *~8.9 kya. This timing corresponds with increased population diversity in cultivars at a similar*
305 *time, calculated using SMC++. Green rectangles along the timeline show "windows of*
306 *opportunity", times when Yemen was connected with the African continent wherein human*
307 *migrations to the Arabian Peninsula may have occurred. The purple rectangle shows the last*
308 *ice age.* **F.** *Directional gene flow analysis using Orientagraph suggests two hypotheses: gene*
309 *flow from the shared ancestral population of all cultivars to the Ethiopian wild individuals*
310 *(subCC), or gene flow from the Typica lineage to Ethiopia (subEE).*

311
312

313  In wild accessions, both subgenomes concordantly showed two population bottlenecks (**Fig.**
314  **2D**) in the SMC++[52] modeling. Assuming a 21-year generation time[53], the oldest bottleneck
315  initiated abruptly around 350 thousand years ago (kya) and ended around 15 kya, at the start of
316  the African humid period (AHP)[54], when climatic conditions were more favourable for Arabica
317  growth. The more recent bottleneck initiated more gradually around 5 kya and lasts to this day.
318  Cultivated accessions, however, exhibited the oldest, but not the more recent bottleneck. In part
319  due to these differences, we also modeled Arabica population history using FastSimcoal2[55],
320  modeling the wild population and cultivars as two separate lineages. In the best-fitting model
321  (**Fig. 2E**) the wild population was predicted to split from the cultivar founding population 1,450
322  generations ago (~30 kya), i.e., before the last glacial maximum. The original founding event
323  was analyzed using the non-admixed wild individuals, revealing an ancestral population
324  bottleneck at 350 kya (**Ext. data Fig. 6A**). Divergence estimates based on gene fractionation,
325  the distribution of nonsynonymous mutations (**Ext. data Fig. 6B**), and calibrated SNP trees
326  (**Fig. 2B**) suggested the allopolyploid founding event occurred at 610 kya, which is close to
327  previous estimates[22,23]. The 350 kya bottleneck, on the other hand, corresponds to that found in
328  the SMC++ analyses (**Fig. 2D**). We therefore consider 610-350 kya a likely time range for the
329  polyploidization event (**Fig. 2E**). The wild and pre-cultivar lineages maintained some gene flow
330  (in terms of migration) until ~8-9 kya, which may have contributed to the modeled increase in
331  effective population size (**Fig. 2D-E**).
332
333  While these data were not able to identify the precise place of origin of the modern cultivated
334  population (see also the following section), the extended period of migration between wild and
335  cultivated accessions suggests that they were separated only by a relatively small geographic
336  distance, such as along the two sides of the African Great Rift Valley (**Fig. 2A-C**). It is also
337  possible that the cultivated lineage could have extended as far as Yemen, and that the end of
338  migration between the two populations could have been caused by the widening of the Bab al-
339  Mandab strait (separating Yemen and Africa) due to rising sea levels[56] at the end of the AHP.
340  A native Arabica population exists in Yemen[57], which could support this hypothesis. The
341  Linnaean sample, together with the Typica and Bourbon cultivars, originate from this second
342  population that was also used to establish cultivation in Yemen, as suggested by the SNP,
343  ADMIXTURE, and PCA analyses (**Fig. 2B**, **Ext. data Fig. 5**).
344
345  In conclusion, our analyses suggest that the Arabica allopolyploidy event occurred between
346  610 and 350 kya, when considering that inbreeding present in *Coffea* populations would
347  accelerate coalescence estimation[58,59]. Earlier work proposing more recent timings, such as 20
348  kya[20], could be underestimates stemming from confounding effects of population bottlenecks
349  in cultivated and wild lineages.
350
351  **Origin of modern cultivars**
352
353  The known breeding history of several of our Arabica cultivars provided us with a gold standard
354  set for deducing the Arabica pedigree using Kinship-based INference for Gwas (KING)[60] (**Fig.**
355  **3**). The method correctly identified the relationships between Bourbon and Typica group
356  cultivars and the Bourbon-Typica crosses in subCC. In contrast, the subEE pedigree showed
357  lower (2nd) order relationships, possibly due to HE in that subgenome (**Ext. data Fig. 7**). Timor
358  hybrid-derived accessions did not show significant relationships to mainline cultivars in subCC
359  (likely due to Robusta introgressions in this subgenome that broke the haplotype blocks, see
360  below), while subEE showed 2nd degree relationships to both the Typica and Bourbon groups
361  (**Fig. 3; Ext. data Fig. 7**), confirming that subEE has not received substantial introgression.
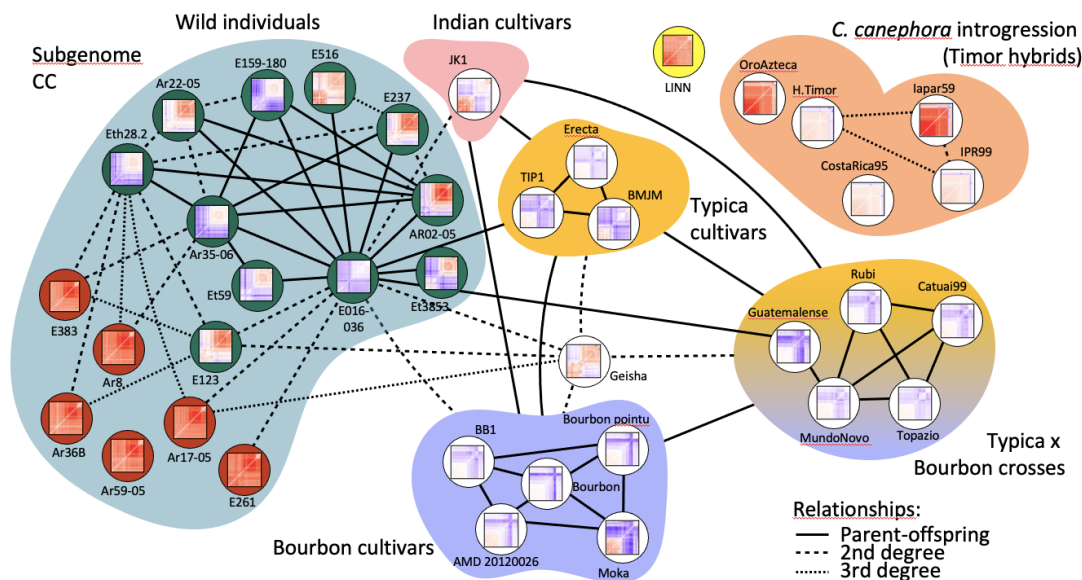
8

**Figure 3.** *Kinship estimation of* C. arabica *accessions, inferred from SNPs in subCC. The degree of relatedness was estimated using Kinship-based INference for GWAS (KING) and describes the number of generations between the related accessions. Thumbnail images show false discovery rate corrected F3 tests of introgression Z-statistics for each of the target individuals. Each cell in the matrix illustrates an F3 test result for the target accession containing introgression from two different sources (x- and y- axis); blue color illustrates significant gene flow (or allele sharing via identity by descent[61]; IBD) from the two source accessions to the target, while red color illustrates lack of gene flow. See **Ext. data Fig. 7** for corresponding analyses in subEE. In the wild accessions, the dark green background highlights the admixed individuals (**Fig. 3B**), while the non-admixed individuals are highlighted with red background. Relationships follow standard nomenclature (e.g., 2nd degree refers to an individual's grandparents, grandchildren etc., whereas 3rd degree refers to great-grandparents, great-grandchildren, etc.)*

Interestingly, Typica, Bourbon and JK1 individuals were also 1st degree related, suggesting direct parent-offspring relationships. Besides confirming their shared Yemeni origins, this finding also underscores the Yemeni germplasm's limited genetic diversity. Further, the old cultivar lines JK1 (Indian), Erecta (Indonesian Typica), BMJM (Caribbean Typica), TIP1 (Brazilian Typica), and BB1 (Brazilian Bourbon) showed 2nd or higher degree relationships with a cluster of closely related wild admixed accessions, centered on E016/136 (**Fig. 2B**). The recently established Geisha cultivar showed similar relationships to the wild admixed individuals and the Bourbon and Typica groups, suggesting common origins. Interestingly, admixed wild accession E016/136 was closely related to both wild and cultivated populations.

In a comparison of geographic origins, wild individuals from the Eastern side of the Great Rift Valley had some levels of admixture and were closely interrelated, while on the Western side, the admixed, related individuals were mostly concentrated around the Gesha region (**Fig. 2C, Fig. 3**). The E016/136 admixed accession, closest to cultivars, demonstrated a first-degree relationship with several wild accessions, of which only Ar35-06 and Eth28.2 were pure representatives of the wild population (**Fig. 2B**). Therefore, these two accessions are genetically closest, in our sample, to the hypothetical true wild parent of cultivated Arabica, with E016/136 representing an intermediate form. Ar35-06 was collected near Gesha mountain, close to the origin of the modern Geisha cultivar. Altogether these data point to the Gesha region as a hotspot of wild accessions amenable to domestication.

399    Admixed wild samples may have originated from a recent hybridization event that occurred
400    before or after their collection from the wild. A third alternative is that the Yemeni population
401    (and hence the cultivars) originated from an admixed population from the Eastern side of the
402    Great Rift Valley or the Gesha region. Analysis of admixture patterns with Orientagraph[62] (**Fig.**
403    **2F**) suggested hybridization with the common ancestor of the Bourbon and Typica lineages in
404    subCC, and of Typica in subEE. In the case of recent hybridization, introduced haplotypes
405    would exist as long contiguous blocks (as in the Timor hybridization, which occurred 100 years
406    ago), while for older events the blocks would be more fragmented due to crossing over.
407    Analysis using the distance fraction ($d_f$) statistic[63] showed the latter to be the case (**Ext. data**
408    **Fig. 8**, indicating that admixture events among wild accessions were not very recent, supporting
409    our third hypothesis.

411    Domestication and cultivation usually involve strong population bottlenecks based on high wild
412    diversity, resulting in reduced genetic diversity in cultivars[64]. However, Arabica nucleotide
413    diversity was already very low in the wild, probably as a result of earlier bottlenecks (**Figs. 2D-**
414    **E**), but only marginally reduced in the pre-cultivated lineage (**Ext. data Fig. 9A**). Bourbon had
415    lower diversity than Typica, probably resulting from the known single-individual bottleneck in
416    this group. Also, the inbreeding coefficients in the wild and cultivated accessions were similar
417    (**Ext. data Fig. 9B**), differing from general expectations for a domesticated species[64].

419    To look for pathways under purifying selection in cultivars, we identified genes with high $F_{ST}$
420    (95 % quantile) between cultivars and wild accessions. This resulted in a set of 1,908 genes that
421    were enriched for the GO categories "cellular response to nitrogen starvation", "regulation of
422    innate immune response" and "regulation of defense response" (**Table S11**), and contained
423    homologs of ammonium transporters *AMT1* and *AMT2*, important for nitrogen uptake in
424    *Coffea*[65], a homolog of the salicylic acid receptor *NONEXPRESSER OF PR GENES 1* (*NPR1*),
425    required in SA signaling and systemic acquired resistance[66], as well as a homolog of the
426    *Arabidopsis LSU2* gene, previously identified as a hub convergently targeted by effectors of
427    pathogens from different kingdoms[67]. A second screen, focused on genes with a large number
428    of high-impact nonsynonymous mutations shared among cultivars (>40% individuals having
429    the mutation), generated a list of 556 genes that were significantly enriched for only one GO
430    category, "defense response" (**Table S12**). From the 22 genes in this category, 16 were NB-
431    ARC domain-containing resistance (R) genes, and two were members of the leucine-rich repeat
432    (LRR) defense gene family. High diversity in immune related responses is one possible
433    pathogen resistance mechanism in plant communities[68], and therefore reduced diversity may
434    have compromised modern Arabica cultivar immunity.

436    The high level of conservation between the Arabica subgenomes and their diploid progenitors
437    may have facilitated spontaneous interspecific hybridization events. This was the case for the
438    Timor hybrid, a spontaneous Robusta x Arabica hybrid resistant to *Hemileia vastatrix*[27]. Our
439    sample set included five descendants of the original Timor hybrid, obtained by backcrossing to
440    Arabica. As expected, the hybridization affected subCC more profoundly, with much higher
441    levels of nucleotide divergence apparent ($F_{ST}$=0.185) than in subEE ($F_{ST}$=0.0897), when
442    comparing cultivars and hybrids. The divergence from wild populations was even greater, with
443    $F_{ST}$=0.254 for subCC and $F_{ST}$=0.138 for subEE, illustrating that introgression occurred almost
444    exclusively within subCC.

446    In the Timor hybrids, the regions found with $d_f$ statistics[63] largely overlapped the introgressed
447    loci identified using $F_{ST}$ scans (**Fig. 4A**) and were found in large blocks, reflecting recent
448    hybridization, and covering 7-11% of the genome (**Fig. 4A, Ext. data Fig. 8**). Transposon
449    Insertion Polymorphisms (TIPs) also overlapped with introgressed regions (Gypsy p=0.0002;
450    Copia p=0.035; Fisher exact test), confirming their recent origin from Robusta (**Fig. 4B**). The
451    introgressed regions overlapped with regions of higher subgenome fractionation (p=0.001873;
452    **Table S13**), possibly due to heterologous recombination between subCC and Robusta, resulting
453    in unequal crossing-over.

454
455 An introgressed region shared by all Timor hybrid lines was evident on chromosome 4 (**Fig.**
456 **4A**). We identified a set of 233 genes shared by all hybrids (**Table S14**). The set contained
457 members of three colocalized tandemly duplicated blocks of resistance-related genes on
458 chromosome 4, subCC, and showed high $F_{ST}$ values between cultivars and introgressed lines.
459 A tandem array of five genes were homologs of *Arabidopsis RPP8*, an NLR resistance locus
460 conferring pleiotropic resistance to several pathogens[69,70]. *RPP8* shows a great amount of
461 variation in *Arabidopsis* alone, where intrachromosomal gene conversion combined with
462 balancing selection contributes to its exceptional diversity[71]. The same subCC region also
463 included a tandem array of ten homologs of *CONSTITUTIVE EXPRESSER OF PR GENES 1*
464 (*CPR1*), a negative regulator of defense response that targets resistance proteins[72,73]. Finally,
465 we identified three duplicates encoding leaf rust 10 disease-resistance locus receptor-like
466 protein kinases (LRK10L). The LRK10L are a gene family that is widespread across plants.
467 First identified as a protein kinase in a locus contributing leaf rust resistance in wheat[74], they
468 found to be upregulated during various biotic and abiotic stresses[75] and confirmed as positive
469 regulators of wheat hypersensitive resistance response to stripe rust fungus[75] and powdery
470 mildew[76].

471
472 The high $F_{ST}$ values between cultivated and introgressed, but not wild individuals (**Fig. 4B**),
473 indicate that the wild population cannot be the source for allelic asymmetries. Nucleotide
474 diversities further illustrate this point; some genes demonstrate lower nucleotide diversity in
475 wild individuals, suggesting these genes to have experienced selective sweeps. To further
476 narrow down candidate genes involved in leaf rust resistance, we reanalyzed comparative gene
477 expression data from susceptible and resistant accessions after *H. vastatrix* inoculation[77]. This
478 analysis identified 723 differentially expressed genes, most of which were associated with
479 defense responses (**Fig. 4B, Tables S14-S14b**). The combination of high $F_{ST}$ values, nucleotide
480 diversities, and differential expression data highlight several strong candidate genes (one *RPP8*,
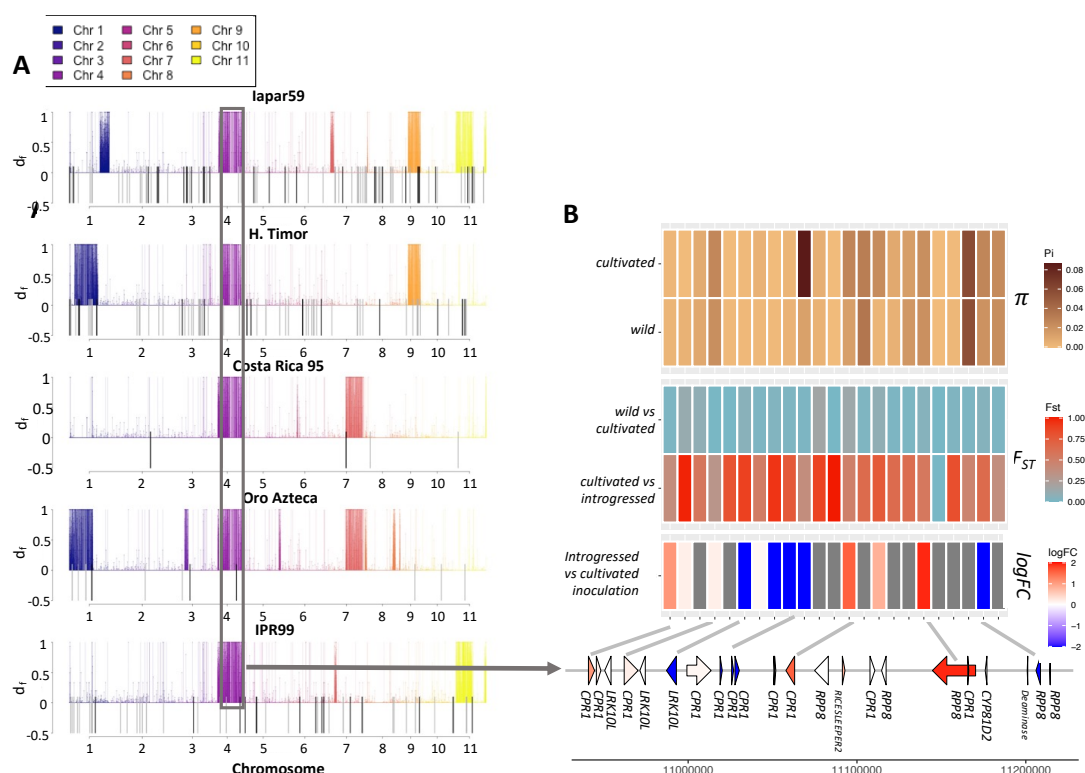481 six *CPR1* and one *LRK10L*) at this locus.

482
483



484
485

486 **Figure 4.** *Introgression of* Coffea canephora *into* H. vastatrix-*resistant* C. arabica *lineages.*
487 *A. Introgression $d_f$ statistic estimated for different Timor hybrid derivatives. Colored lines*
488 *above the axis mark regions of significant introgression in the line under inspection, and are*
489 *colored by chromosome. The shared introgressed region on Chr 4 is colored in purple and*
490 *boxed. Transposon Insertion Polymorphisms are represented as lines below the X axis and*
491 *exhibit overlap with introgressed regions.* ***B.*** *The shared introgressed genomic region on*
492 *subCC chromosome 4 contains a cluster of R genes (RPP8), a cluster of homologs of a*
493 *negative regulator of R genes (CPR1), and a cluster of homologs of leaf rust resistance 10*
494 *kinases (LRK10L) (bottom). The heatmap shows, from the bottom up, (i) log fold change of*
495 *gene expression after* H. vastatrix *inoculation, when comparing resistant Timor hybrid*
496 *lineage against a susceptible cultivar; red color means elevated expression in the hybrid, and*
497 *blue decreased expression. (ii) Fixation index ($F_{ST}$) values for the introgessed lines vs*
498 *cultivars and between cultivars and wild accessions. (iii) Nucleotide diversity for the wild and*
499 *cultivated accessions for each gene coding region, plus the flanking 2kb upstream and*
500 *downstream of the region.*
501
502 **Conclusions and outlook**
503
504 Besides providing genomic resources for molecular breeding of one of the most important
505 agricultural commodities, our Arabica, Robusta and Eugenioides genomes provide a unique
506 window into the genome evolution of a recently formed allopolyploid stemming from two
507 closely related species. Our Arabica data did not suggest a genomic shock induced by
508 allopolyploidy, but instead, only higher LTR transposon turnover rate. Genome fractionation
509 rates remained basically unaltered before and after the allopolyploidy event. Likewise, no
510 global subgenome dominance in gene expression was observed, but rather a mosaic-type
511 pattern as in other recent polyploids[44,47], affecting the expression of individual gene family
512 members. However, similar to octoploid strawberry[8], we detected genome dominance in terms
513 of biased homoeologous exchanges favoring subCC. Since Robusta has one of the widest
514 geographic ranges in the *Coffea* genus, whereas Eugenioides is more range-limited, this biased
515 HE might be adaptive. This hypothesis was supported by the site frequency spectrum of HE
516 loci, showing signs of directional selection (**Ext. data Fig. 3**). Intriguingly, transposable
517 insertion polymorphisms significantly overlapped with tandem gene duplications and
518 biosynthetic gene clusters, hinting at their possible roles in cluster evolution.
519
520 Domestication of perennial species like Arabica coffee differs markedly from that of annual
521 crops, consisting instead of three phases: selection of outstanding genotypes from wild forests,
522 clonal propagation and cultivation, then breeding and diversification[78]. In addition to being a
523 perennial crop, Arabica is also a predominantly autogamous allopolyploid, which puts it in a
524 class of its own. We show here that genetic diversity was already very low among wild
525 accessions, due to multiple pre-domestication bottlenecks, and that the genotypes selected for
526 cultivation by humans (both the ancient cultivated Ethiopian landraces, and the recent Geisha
527 cultivar) already were somewhat admixed between divergent lineages. The resequenced
528 accessions displayed a geographic split along the Eastern versus Western sides of the Great Rift
529 Valley, with cultivated coffee variants all placed with the Eastern population. Such admixture
530 has played a large role in breeding many fruit-bearing crops, the non-polyploid allogamous
531 perennial lychee being one of the most extreme cases[59].
532
533 The prevalent autogamy of Arabica, combined with the multiple genetic bottlenecks it
534 underwent in the wild, may have selectively purged deleterious alleles, explaining the capacity
535 of the species to survive single-plant bottlenecks that occurred during its cultivation. An
536 additional element buffering deleterious alleles was probably Arabica's allopolyploidy itself,
537 which provided some level of heterosis[79]. However, the narrow genetic basis of both cultivated
538 and wild modern Arabica constitutes a major drawback, as well as an obstacle for its breeding
539 using wild genepool diversity. On the other hand, the extensive collinearity of its CC and EE
540 subgenomes with those of its Robusta and Eugenioides progenitors is likely to facilitate

541 introgression of interesting traits from these species, as already happened historically in the
542 Timor spontaneous hybrid. The high-quality genome sequences of the three species provided
543 in this work, together with the identification of the genomic region conferring resistance to
544 coffee leaf rust, constitute a cornerstone for the breeding of novel Arabica varieties with
545 superior adaptability and pathogen resistance.
546
547 **Data availability**
548 Coffee genome assemblies are available at CoGe (https://genomevolution.org/): *C. canephora*:
549 50947, *C. eugenioides*: 60235, and *C. arabica*: 66663 (Pacbio HiFi) and 53628 (Pacbio). All
550 genome information, including the VCF files with SNP information are available at
551 ftp.solgenomics.net; the genome data is also available at ORCAE
552 (https://bioinformatics.psb.ugent.be/orcae/overview/Coara and
553 https://bioinformatics.psb.ugent.be/gdb/coffea_arabica/).
554 The sequencing data have been deposited to NCBI under bioproject ID: PRJNA698600.
555
556 **TABLES**
557
558 **Table 1.** Statistics of the *Coffea* assemblies presented in this paper.

| Assembly | *C. eugenioides* | *C. canephora* | *C. arabica* | *C. arabica* HiFi |
|---|---|---|---|---|
| Projected genome size (Mb)* | 682 | 705 | 1281 | 1281 |
| Total assembly length (Mb) | 661 | 672 | 1,088 | 1,198 |
| % of projected genome | 96.9% | 95.3% | 84.9% | 93.5% |
| N scaffolds | 253 | 3,033 | 8,474 | 132 |
| Scaffold N50 | 61.3 Mb | 50.1 Mb | 32.7 Mb | 53.7 Mb |
| N contigs | 5,736 | 3,755 | 11,863 | 238** |
| Contig N50 (Mb) | 0.40 | 0.76 | 0.23 | 30.0 |
| Pseudochromosomes (Mb) | n.a. | 583 | 801 | 1192 |
| % of projected genome | n.a. | 82.7% | 62.5% | 93.1% |
| N. genes | 33,505 | 28,880 | 56,670 | 69,314 |
| Genes in pseudochromosomes | n.a. | 27,881 | 50,410 | 69,067 |
| % genes in pseudochromosomes | n.a. | 97% | 89% | 99.6% |
| BUSCO genome | | | | |
| complete | 96.7% | 97.4% | 97.6% | 97.9% |
| single | 88.5% | 94.8% | 20.1% | 4.3 % |
| duplicated | 8.2% | 2.6% | 77.5% | 93.6 % |
| fragmented | 1.1% | 0.9% | 0.8% | 0.8 % |
| missing | 2.2% | 1.7% | 1.6% | 1.3 % |
| total | 2,326 | 2,326 | 2,326 | 2,326 |
| BUSCO annotation | | | | |
| complete | 94.9% | 96.2% | 92.1% | 97.3% |
| single | 82.4% | 92.8% | 33.3% | 4.1% |
| duplicated | 12.5% | 3.4% | 58.8% | 93.2% |
| fragmented | 2.1% | 1.5% | 2.8% | 0.8% |
| missing | 3.0% | 2.3% | 5.1% | 1.9% |
| total | 2,326 | 2,326 | 2,326 | 2,326 |

559 *From the Plant DNA C-values database: https://cvalues.science.kew.org/; **After gap filling;
560 ***Denoeud et al, 2014; ****Scalabrin et al, 2020.
561
562 **ACKNOWLEDGMENTS**

13

581

## AUTHOR CONTRIBUTIONS

583 Conceived the study: AdK, DC, PD. Provided genetic resources: AA, AN, CK, EC, GHS, HR,
584 LB, LFP, LP, MS, MTB, OGF, PaM, PM, PH, US. Carried out DNA sequencing: AC, CF, DM,
585 GL, JeS, LB, MK, ND, PD, SM. Sequencing of the Linnaean accession: EG. Carried out
586 genome assembly: SS, CW, JS, SP, LM. Genetic mapping: PR, MR, JS. Genome annotation:
587 AR, SS, LM, JS, SR, VP, ZQW, DD, SIS, MM, RA, SMCL, ML, MP, CT-D, GG. Annotation
588 of non-coding RNA: ARP, JE, PS. Transposable element annotation and analysis: SOA, AG,
589 RG. Telomere identification: VAA, WCM. Analyzed genome evolution: ZY, ZC, DSM, RG,
590 JM, DS, LC-P, TL, TK, VAA, SOA, AG, JS. Gene family analysis: ZQW, VP, DD, GG, SF,
591 VAA, SiR, JS. RNA-seq data analysis: AR, SP, SiR, JS. Provided RNA-seq data: RH. Analyzed
592 population data: JS. Analysed GBS data: YB, RG. Arranged online data access: LM, SR. Wrote
593 the first draft: JS, completed with input from GG, DS, VAA, LFP, RG, SR, AdK, PD, VP, LM,
594 DC, DD, SP, AA, as well as PM, YB, TR, YVdP, and all co-authors.

595

## COMPETING INTERESTS

597 The authors declare no competing financial interests.

598

## REFERENCES

600 1. Van de Peer, Y., Mizrachi, E. & Marchal, K. The evolutionary significance of
601 polyploidy. *Nature Reviews Genetics* **18**, 411-424 (2017).
602 2. Van de Peer, Y., Ashman, T.-L., Soltis, P.S. & Soltis, D.E. Polyploidy: an evolutionary
603 and ecological force in stressful times. *The Plant Cell* **33**, 11-26 (2021).
604 3. Leebens-Mack, J.H. *et al.* One thousand plant transcriptomes and the phylogenomics
605 of green plants. *Nature* **574**, 679-685 (2019).
606 4. Sun, H. *et al.* Chromosome-scale and haplotype-resolved genome assembly of a
607 tetraploid potato cultivar. *Nature Genetics* **54**, 342-348 (2022).
608 5. Athiyannan, N. *et al.* Long-read genome sequencing of bread wheat facilitates disease
609 resistance gene cloning. *Nature Genetics* **54**, 227-231 (2022).
610 6. Wu, S. *et al.* Genome sequences of two diploid wild relatives of cultivated
611 sweetpotato reveal targets for genetic improvement. *Nature Communications* **9**, 4580
612 (2018).
613 7. Wang, T. *et al.* A complete gap-free diploid genome in Saccharum complex and the
614 genomic footprints of evolution in the highly polyploid Saccharum genus. *Nature
615 Plants* **9**, 554-571 (2023).
616 8. Edger, P.P. *et al.* Origin and evolution of the octoploid strawberry genome. *Nature
617 Genetics* **51**, 541-547 (2019).
618 9. Li, F. *et al.* Genome sequence of cultivated Upland cotton (Gossypium hirsutum TM-

619           1) provides insights into genome evolution. *Nature Biotechnology* **33**, 524-530
620           (2015).

621 10.    Chalhoub, B. *et al.* Early allopolyploid evolution in the post-Neolithic *Brassica napus*
622           oilseed genome. *Science* **345**, 950-953 (2014).

623 11.    Sattler, M.C., Carvalho, C.R. & Clarindo, W.R. The polyploidy and its key role in
624           plant breeding. *Planta* **243**, 281-296 (2016).

625 12.    McClintock, B. The Significance of Responses of the Genome to Challenge. *Science*
626           **226**, 792-801 (1984).

627 13.    Sha, Y. *et al.* Genome shock in a synthetic allotetraploid wheat invokes subgenome-
628           partitioned gene regulation, meiotic instability, and karyotype variation. *Journal of*
629           *Experimental Botany*, erad247 (2023).

630 14.    Thomas, B.C., Pedersen, B. & Freeling, M. Following tetraploidy in an Arabidopsis
631           ancestor, genes were removed preferentially from one homeolog leaving clusters
632           enriched in dose-sensitive genes. *Genome Research* **16**, 934-946 (2006).

633 15.    Schnable, J.C., Springer, N.M. & Freeling, M. Differentiation of the maize
634           subgenomes by genome dominance and both ancient and ongoing gene loss.
635           *Proceedings of the National Academy of Sciences* **108**, 4069 (2011).

636 16.    Gaeta, R.T., Pires, J.C., Iniguez-Luy, F., Leon, E. & Osborn, T.C. Genomic Changes
637           in Resynthesized Brassica napus and Their Effect on Gene Expression and
638           Phenotype. *The Plant Cell* **19**, 3403-3417 (2007).

639 17.    Burns, R. *et al.* Gradual evolution of allopolyploidy in Arabidopsis suecica. *Nature*
640           *Ecology & Evolution* **5**, 1367-1381 (2021).

641 18.    Conant, G.C., Birchler, J.A. & Pires, J.C. Dosage, duplication, and diploidization:
642           clarifying the interplay of multiple models for duplicate gene evolution over time.
643           *Current Opinion in Plant Biology* **19**, 91-98 (2014).

644 19.    Carvalho, A. *et al.* Melhoramento do cafeeiro: IV - Café Mundo Novo. *Bragantia* **12**,
645           97-130 (1952).

646 20.    Scalabrin, S. *et al.* A single polyploidization event at the origin of the tetraploid
647           genome of *Coffea arabica* is responsible for the extremely low genetic variation in
648           wild and cultivated germplasm. *Sci Rep* **10**, 4642 (2020).

649 21.    Cenci, A., Combes, M.-C. & Lashermes, P. Genome evolution in diploid and
650           tetraploid *Coffea* species as revealed by comparative analysis of orthologous genome
651           segments. *Plant Molecular Biology* **78**, 135-145 (2012).

652 22.    Bawin, Y. *et al.* Phylogenomic analysis clarifies the evolutionary origin of *Coffea*
653           *arabica*. *Journal of Systematics and Evolution* **59**, 953-963 (2020).

654 23.    Yu, Q. *et al.* Micro-collinearity and genome evolution in the vicinity of an ethylene
655           receptor gene of cultivated diploid and allotetraploid coffee species (*Coffea*). *The*
656           *Plant Journal* **67**, 305-317 (2011).

657 24.    Merot-L'anthoene, V. *et al.* Development and evaluation of a genome-wide Coffee
658           8.5K SNP array and its application for high-density genetic mapping and for
659           investigating the origin of *Coffea arabica* L. *Plant Biotechnology Journal* **17**, 1418-
660           1430 (2019).

661 25.    Wellman, F.L. *Coffee: botany, cultivation and utilization*, (L. Hill, London, 1961).

662 26.    Lécolier, A., Besse, P., Charrier, A., Tchakaloff, T.-N. & Noirot, M. Unraveling the
663           origin of *Coffea arabica* 'Bourbon pointu' from La Réunion: a historical and
664           scientific perspective. *Euphytica* **168**, 1-10 (2009).

665 27.    Clarindo, W.R., Carvalho, C.R., Caixeta, E.T. & Koehler, A.D. Following the track of
666           "Híbrido de Timor" origin by cytogenetic and flow cytometry approaches. *Genetic*
667           *Resources and Crop Evolution* **60**, 2253-2259 (2013).

668 28.    Bertrand, B., Guyot, B., Anthony, F. & Lashermes, P. Impact of the *Coffea canephora*
669           gene introgression on beverage quality of *C. arabica*. *Theoretical and Applied*
670           *Genetics* **107**, 387-394 (2003).

671 29.    Marie, L. *et al.* G × E interactions on yield and quality in *Coffea arabica*: new F1
672           hybrids outperform American cultivars. *Euphytica* **216**, 78 (2020).

673 30.    Bertrand, B., Villegas Hincapié, A.M., Marie, L. & Breitler, J.-C. Breeding for the

674        main agricultural farming of *Arabica* coffee. *Frontiers in Sustainable Food Systems*
675        **5**(2021).
676  31.  Breitler, J.-C. *et al.* CRISPR/Cas9-mediated efficient targeted mutagenesis has the
677        potential to accelerate the domestication of *Coffea canephora. Plant Cell, Tissue and*
678        *Organ Culture (PCTOC)* **134**, 383-394 (2018).
679  32.  Berthaud, J. Etude cytogénétique d'un haploïde de *Coffea arabica* L. *Café, Cacao,*
680        *Thé (Francia) v. 20 (2) p. 91-96 (1976).*
681  33.  Denoeud, F. *et al.* The coffee genome provides insight into the convergent evolution
682        of caffeine biosynthesis. *Science* **345**, 1181-1184 (2014).
683  34.  Pellicer, J. & Leitch, I.J. The Plant DNA C-values database (release 7.1): an updated
684        online repository of plant genome size data for comparative studies. *New Phytologist*
685        **226**, 301-305 (2020).
686  35.  Manni, M., Berkeley, M.R., Seppey, M., Simão, F.A. & Zdobnov, E.M. BUSCO
687        Update: Novel and Streamlined Workflows along with Broader and Deeper
688        Phylogenetic Coverage for Scoring of Eukaryotic, Prokaryotic, and Viral Genomes.
689        *Molecular Biology and Evolution* **38**, 4647-4654 (2021).
690  36.  Petit, M. *et al.* Mobilization of retrotransposons in synthetic allotetraploid tobacco.
691        *New Phytologist* **186**, 135-147 (2010).
692  37.  Sarilar, V. *et al.* Allopolyploidy has a moderate impact on restructuring at three
693        contrasting transposable element insertion sites in resynthesized Brassica napus
694        allotetraploids. *New Phytologist* **198**, 593-604 (2013).
695  38.  Bird, K.A., VanBuren, R., Puzey, J.R. & Edger, P.P. The causes and consequences of
696        subgenome dominance in hybrids and recent polyploids. *New Phytologist* **220**, 87-93
697        (2018).
698  39.  Göbel, U. *et al.* Robustness of transposable element regulation but no genomic shock
699        observed in interspecific *Arabidopsis* hybrids. *Genome Biology and Evolution* **10**,
700        1403-1415 (2018).
701  40.  Birchler, J.A. & Veitia, R.A. The gene balance hypothesis: implications for gene
702        regulation, quantitative traits and evolution. *The New phytologist* **186**, 54-62 (2010).
703  41.  Zeiss, D.R., Piater, L.A. & Dubery, I.A. Hydroxycinnamate Amides: Intriguing
704        Conjugates of Plant Protective Metabolites. *Trends in Plant Science* **26**, 184-195
705        (2021).
706  42.  Bird, K.A. *et al.* Replaying the evolutionary tape to investigate subgenome
707        dominance in allopolyploid Brassica napus. *New Phytologist* **230**, 354-371 (2021).
708  43.  Combes, M.-C., Joët, T., Stavrinides, A.K. & Lashermes, P. New cup out of old
709        coffee: contribution of parental gene expression legacy to phenotypic novelty in
710        coffee beans of the allopolyploid Coffea arabica L. *Annals of Botany* **131**, 157-170
711        (2023).
712  44.  Yoo, M.J., Szadkowski, E. & Wendel, J.F. Homoeolog expression bias and expression
713        level dominance in allopolyploid cotton. *Heredity* **110**, 171-180 (2013).
714  45.  Meyer, F.G., Fernie, L.M., Narasimhaswami, R.L., Monaco, L.C. & Greathead, D.J.
715        FAO Coffee Mission to Ethiopia, 1964-1965. (1968).
716  46.  Halle, F. Echantillonnage du matériel *Coffea arabica* récolté en Ethiopie. *Bulletin -*
717        *IFCC*, 13-18 (1978).
718  47.  Chalhoub, B. *et al.* Early allopolyploid evolution in the post-Neolithic Brassica napus
719        oilseed genome. *Science* **345**, 950-953 (2014).
720  48.  Krug, C.A.M., A.J.T. Cytological observations in *Coffea* – IV. *J Genet* **39**, 189–203
721        (1940).
722  49.  Cros, J. *et al.* Phylogenetic Analysis of Chloroplast DNA Variation in *Coffea* L.
723        *Molecular Phylogenetics and Evolution* **9**, 109-117 (1998).
724  50.  Lashermes, P. *et al.* Molecular characterisation and origin of the Coffea arabica L.
725        genome. *Molecular and General Genetics MGG* **261**, 259-266 (1999).
726  51.  Wu, Y. *et al.* Genomic mosaicism due to homoeologous exchange generates extensive
727        phenotypic diversity in nascent allopolyploids. *National Science Review* **8**, nwaa277
728        (2021).

| 729 | 52. | Terhorst, J., Kamm, J.A. & Song, Y.S. Robust and scalable inference of population |
| 730 | | history from hundreds of unphased whole genomes. *Nature Genetics* **49**, 303-309 |
| 731 | | (2017). |
| 732 | 53. | Moat, J., Gole, T.W. & Davis, A.P. Least concern to endangered: Applying climate |
| 733 | | change projections profoundly influences the extinction risk assessment for wild |
| 734 | | *Arabica* coffee. *Global Change Biology* **25**, 390-403 (2019). |
| 735 | 54. | Kuper, R. & Kröpelin, S. Climate-controlled holocene occupation in the Sahara: |
| 736 | | motor of Africa's evolution. *Science* **313**, 803-807 (2006). |
| 737 | 55. | Excoffier, L. *et al.* fastsimcoal2: demographic inference under complex evolutionary |
| 738 | | scenarios. *Bioinformatics* **37**, 4882-4885 (2021). |
| 739 | 56. | Lambeck, K. *et al.* Sea level and shoreline reconstructions for the Red Sea: isostatic |
| 740 | | and tectonic considerations and implications for hominin migration out of Africa. |
| 741 | | *Quaternary Science Reviews* **30**, 3542-3574 (2011). |
| 742 | 57. | Montagnon, C., Mahyoub, A., Solano, W. & Sheibani, F. Unveiling a unique genetic |
| 743 | | diversity of cultivated *Coffea arabica* L. in its main domestication center: Yemen. |
| 744 | | *Genetic Resources and Crop Evolution* **68**, 2411-2422 (2021). |
| 745 | 58. | Nordborg, M. & Donnelly, P. The coalescent process with selfing. *Genetics* **146**, 1185 |
| 746 | | (1997). |
| 747 | 59. | Hu, G. *et al.* Two divergent haplotypes from a highly heterozygous lychee genome |
| 748 | | suggest independent domestication events for early and late-maturing cultivars. |
| 749 | | *Nature Genetics* **54**, 73-83 (2022). |
| 750 | 60. | Manichaikul, A. *et al.* Robust relationship inference in genome-wide association |
| 751 | | studies. *Bioinformatics* **26**, 2867-2873 (2010). |
| 752 | 61. | Lan, T. *et al.* Insights into bear evolution from a Pleistocene polar bear genome. |
| 753 | | *Proceedings of the National Academy of Sciences* **119**, e2200016119 (2022). |
| 754 | 62. | Molloy, E.K., Durvasula, A. & Sankararaman, S. Advancing admixture graph |
| 755 | | estimation via maximum likelihood network orientation. *Bioinformatics* **37**, i142-i150 |
| 756 | | (2021). |
| 757 | 63. | Pfeifer, B. & Kapan, D.D. Estimates of introgression as a function of pairwise |
| 758 | | distances. *BMC Bioinformatics* **20**, 207 (2019). |
| 759 | 64. | Gaut, B.S., Seymour, D.K., Liu, Q. & Zhou, Y. Demography and its effects on |
| 760 | | genomic variation in crop domestication. *Nature Plants* **4**, 512-520 (2018). |
| 761 | 65. | dos Santos, T.B., Baba, V.Y., Vieira, L.G.E., Pereira, L.F.P. & Domingues, D.S. The |
| 762 | | urea transporter DUR3 is differentially regulated by abiotic and biotic stresses in |
| 763 | | coffee plants. *Physiology and Molecular Biology of Plants* **27**, 203-212 (2021). |
| 764 | 66. | Wang, W. *et al.* Structural basis of salicylic acid perception by Arabidopsis NPR |
| 765 | | proteins. *Nature* **586**, 311-316 (2020). |
| 766 | 67. | Mukhtar, M.S. *et al.* Independently evolved virulence effectors converge onto hubs in |
| 767 | | a plant immune system network. *Science* **333**, 596-601 (2011). |
| 768 | 68. | Jousimo, J. *et al.* Ecological and evolutionary effects of fragmentation on infectious |
| 769 | | disease dynamics. *Science* **344**, 1289-1293 (2014). |
| 770 | 69. | Cooley, M.B., Pathirana, S., Wu, H.J., Kachroo, P. & Klessig, D.F. Members of the |
| 771 | | *Arabidopsis* HRT/RPP8 family of resistance genes confer resistance to both viral and |
| 772 | | oomycete pathogens. *The Plant cell* **12**, 663-676 (2000). |
| 773 | 70. | Mohr, T.J. *et al.* The *Arabidopsis* downy mildew resistance gene *RPP8* is induced by |
| 774 | | pathogens and salicylic acid and is regulated by W-box *cis* elements. *Molecular* |
| 775 | | *Plant-Microbe Interactions®* **23**, 1303-1315 (2010). |
| 776 | 71. | MacQueen, A. *et al.* Population genetics of the highly polymorphic *RPP8* gene |
| 777 | | family. *Genes* **10**(2019). |
| 778 | 72. | Cheng, Y.T. *et al.* Stability of plant immune-receptor resistance proteins is controlled |
| 779 | | by SKP1-CULLIN1-F-box (SCF)-mediated protein degradation. *Proceedings of the* |
| 780 | | *National Academy of Sciences*, 201105685 (2011). |
| 781 | 73. | Hedtmann, C. *et al.* The Plant Immunity Regulating F-Box Protein CPR1 Supports |
| 782 | | Plastid Function in Absence of Pathogens. *Frontiers in Plant Science* **8**(2017). |
| 783 | 74. | Feuillet, C., Schachermayr, G. & Keller, B. Molecular cloning of a new receptor-like |

784      kinase gene encoded at the Lr10 disease resistance locus of wheat. *The Plant Journal*
785      **11**, 45-52 (1997).
786   75.   Zhou, H. *et al.* Molecular analysis of three new receptor-like kinase genes from
787      hexaploid wheat and evidence for their participation in the wheat hypersensitive
788      response to stripe rust fungus infection. *The Plant Journal* **52**, 420-434 (2007).
789   76.   Xia, T. *et al.* Efficient expression and function of a receptor-like kinase in wheat
790      powdery mildew defence require an intron-located MYB binding site. *Plant*
791      *Biotechnology Journal* **19**, 897-909 (2021).
792   77.   Florez, J.C. *et al.* High throughput transcriptome analysis of coffee reveals
793      prehaustorial resistance in response to *Hemileia vastatrix* infection. *Plant Molecular*
794      *Biology* **95**, 607-623 (2017).
795   78.   Gaut, B.S., Díez, C.M. & Morrell, P.L. Genomics and the Contrasting Dynamics of
796      Annual and Perennial Domestication. *Trends in Genetics* **31**, 709-719 (2015).
797   79.   Chen, Z.J. Molecular mechanisms of polyploidy and hybrid vigor. *Trends in Plant*
798      *Science* **15**, 57-71 (2010).

799