

1 Genetic differentiation and intrinsic genomic features
2 explain variation in recombination hotspots among
3 cocoa tree populations

4 Enrique J. Schwarzkopf¹, Juan C. Motamayor², and Omar E. Cornejo^{1,*}

5 ¹School of Biological Sciences, Washington State University, 100 Dairy Road,
6 Pullman, WA 99164, USA

7 ²Universal Genetic Solutions, LLC

8 *Corresponding author: ocornejo@gmail.com

9 April 20, 2019

10 Abstract

11 Our study investigates the possible drivers of recombination hotspots in *Theobroma cacao*
12 using ten genetically differentiated populations. This constitutes the first time that recom-
13 bination rates from more than two populations of the same species have been compared,
14 providing a novel view of recombination at the population-divergence time-scale. For each
15 population, a fine-scale recombination map was generated using under the coalescent with a
16 standard method based on linkage disequilibrium (LD). They revealed higher recombination
17 rates in a domesticated population and a population that has undergone a recent bottleneck.
18 We address whether the pattern of recombination rate variation along the chromosome is
19 sensitive to the uncertainty in the per-site estimates. We find that uncertainty, as assessed
20 from the Markov chain Monte Carlo iterations is orders of magnitude smaller than the scale of
21 variation of the recombination rates genome-wide. We inferred hotspots of recombination for
22 each population and find that the genomic locations of these hotspots correlate with genetic
23 differentiation between populations (F_{ST}). We developed novel randomization approaches
24 to generate appropriate null models for understanding the association between hotspots of
25 recombination and both DNA sequence motifs and genomic features. Hotspot regions con-
26 tained fewer known retroelement sequences than expected, and were overrepresented near
27 transcription start and termination sites. Our findings indicate that recombination hotspots
28 are evolving in a way that is consistent with genetic differentiation, but are also preferentially
29 driven to regions of the genome that are up or downstream from coding regions.

30 Introduction

31 Genetic variation is fundamental for evolutionary forces like selection and genetic drift to
32 act. Selection and drift also contribute to a loss of variation, which means that they must act
33 in conjunction with forces that maintain variation along the genome in order for populations
34 to continue evolving over prolonged periods of time. Recombination's rearranging of genetic
35 material onto different backgrounds generates a larger set of haplotype combinations on
36 which selection can act, reducing the magnitude of Hill-Robertson interference (Felsenstein,
37 1974). Different regimes of recombination can strongly influence how efficient selection is
38 at purging deleterious mutations and increasing the frequency of beneficial mutations in the
39 population (Felsenstein, 1974).

40 One way to elucidate the distribution of recombination events along the genome is by
41 using fine-scale recombination maps (Myers et al., 2005; Auton et al., 2012; Brunshwig et al.,
42 2012; Paape et al., 2012; Choi et al., 2013; Hellsten et al., 2013; Singhal et al., 2015; Stevison
43 et al., 2016). These maps are constructed with methods that leverage current patterns
44 of linkage disequilibrium (LD) using the coalescent, in order to estimate historical rates of
45 recombination between sites along the genome (Auton and McVean, 2007). Studies in a wide
46 range of species have shown that recombination rates are not uniform along the genome and
47 general patterns of variation have been described (Begun and Aquadro, 1992; Akhunov et al.,
48 2003; Wu et al., 2003; Anderson et al., 2004; McVean et al., 2004; Mézard, 2006; Kim et al.,
49 2007; Gore et al., 2009; Schnable et al., 2009; Branca et al., 2011; Paape et al., 2012). One of
50 these patterns is the reduced recombination rate in centromeric regions of the chromosomes
51 and the progressive increase of recombination rates as the physical distance to telomeres
52 decreases (Begun and Aquadro, 1992; Akhunov et al., 2003; Wu et al., 2003; Anderson et al.,
53 2004; Gore et al., 2009; Schnable et al., 2009). This pattern has also been shown to arise in
54 simulation studies (e.g. Mackiewicz et al., 2010). Another interesting pattern that has been

55 observed is that of regions with unusually high rates of recombination spread throughout
56 chromosomes: recombination hotspots (McVean et al., 2004; Brunschwrig et al., 2012; Paape
57 et al., 2012; Hellsten et al., 2013; Stevison et al., 2016; Shanfelter et al., 2018). In this study,
58 we define hotspots locally, requiring that their recombination rate be unusually high when
59 compared to neighboring regions. The importance of recombination hotspots lies in their
60 ability to shuffle genetic variation at higher rates than the rest of the genome, profoundly
61 impacting the dynamics of selection for or against specific mutations (Felsenstein, 1974).

62 A variety of genomic features have been identified as being associated with regions of
63 high recombination. Recombination hotspots have been linked to transcriptional start sites
64 (TSSs) and transcriptional termination sites (TTSs) in *Arabidopsis thaliana*, *Taeniopygia*
65 *guttata*, *Poephila acuticauda*, and humans (Myers et al., 2005; Choi et al., 2013; Singhal et al.,
66 2015). In *Mimulus guttatus* hotspots were found to be associated with CpG islands (short
67 segments of cytosine and guanine rich DNA, associated with promoter regions) (Hellsten
68 et al., 2013). CpG islands were also associated with increased recombination rates in humans
69 and chimpanzees (Auton et al., 2012). These patterns point to recombination occurring
70 frequently near, but not within coding regions. The formation of chiasmata is important
71 for the proper disjunction of chromosomes during meiosis (Martinez-perez et al., 2008),
72 but repeated double-strand breaks can lead to an increased mutation rate (Rodgers and
73 McVey, 2015). In coding regions in particular, this excess mutation rate can have a high
74 evolutionary cost, due to the likelihood of novel deleterious mutations being higher than that
75 of beneficial ones (Haldane, 1937; Crow and Kimura, 1970; Wloch et al., 2001; Sanjuán et al.,
76 2004; Eyre-Walker and Keightley, 2007). Recombination hotspots have also been found to be
77 correlated with particular DNA sequence motifs. In some mammals, including *Mus musculus*
78 (Brunschwrig et al., 2012) and apes (Auton et al., 2012; Stevison et al., 2016) binding sites for
79 PRDM9, a histone trimethylase with a DNA zinc-finger binding domain, have been found to
80 correlate with recombination hotspots. In *Arabidopsis thaliana*, proteins that limit overall

81 recombination rate have been identified, leading to a genome-wide increase in recombination
82 rate in knockout mutants (Fernandes et al., 2018). However, these *Arabidopsis* proteins have
83 not been shown to direct recombination to particular regions, and are therefore not expected
84 to affect the location of recombination hotspots.

85 Comparisons of recombination hotspots between pairs of populations have yielded vary-
86 ing results. Hinch et al. (2011) found that, at finer scales, the genetic maps of European and
87 African human populations were significantly different. They also found that, when look-
88 ing at hotspots in the major histocompatibility complex, the African populations showed a
89 hotspot that was not present in Europeans, but all European hotspots were found in African
90 populations (Hinch et al., 2011). Recent work on recombination in apes (Stevison et al.,
91 2016) found little correlation of recombination rates in orthologous hotspot regions when
92 looking between species, but a strong correlation when comparing between two populations
93 of the same species. Other studies have also found very little sharing of hotspots between
94 humans and chimpanzees Ptak et al. (2005); Winckler et al. (2005). Additionally, the dy-
95 namic of changing hotspot locations observed in humans and other apes has been observed
96 in simulations Mackiewicz et al. (2013). This suggests that recombination hotspots are po-
97 tentially changing in ways that match demographic patterns, differentiating at a similar rate
98 as genomic sequences.

99 The identification of ten genetically differentiated populations of the cocoa tree, *Theo-*
100 *broma cacao* (Motamayor et al., 2008; Cornejo et al., 2018) can be leveraged to study
101 population-level drivers of recombination hotspots. These ten populations originate from
102 different regions of South and Central America, and include one fully domesticated pop-
103 ulation (Criollo), used in the production of fine chocolate, and nine wilder, more resilient
104 populations which generate higher cocoa yield than the Criollo variety (Fig. 1) (Motamayor
105 et al., 2008; Henderson et al., 2007; Cornejo et al., 2018). These ten populations have been
106 shown to have strong signatures of differentiation between them (F_{ST} values ranging from

107 0.16 to 0.65) and they separate into clear clusters of ancestry (Cornejo et al., 2018). Com-
108 paring the locations of hotspots between these ten populations of *T. cacao* can contribute
109 to the understanding of hotspot turnover at the population-divergence time-scale. These
110 comparisons also contribute to our understanding of how demographics impact the turnover
111 of recombination hotspot locations.

112 Fine-scale, LD-based recombination maps have been constructed for a number of plant
113 models (Paape et al., 2012; Choi et al., 2013; Hellsten et al., 2013), identifying a variety of
114 features correlated to recombination rate. Unlike these model plants with short generation
115 times, *T. cacao* is a perennial woody plant with a five-year generation time (Henderson
116 et al., 2007). The size and long generation time of *T. cacao* makes direct measurements
117 of recombination impractical. However, historical recombination can be estimated for *T.*
118 *cacao* using coalescent based methods (Auton and McVean, 2007). Theoretical studies have
119 shown that population structure can generate artificially inflated measures of LD (Ohta,
120 1982; Li and Nei, 1974), which would be detrimental to our estimates of recombination.
121 For this reason recombination maps were constructed independently for each population.
122 In contrast to previous studies, which have focused primarily on recombination rates, this
123 study attempts to describe the relationship between recombination hotspots and a variety
124 of factors.

125 We used an LD-based method to estimate recombination rates, which we then analyzed
126 with a maximum likelihood statistical framework to infer the location of recombination
127 hotspots. The location of hotspots were compared across populations and a novel resam-
128 pling scheme tailored to the genomic architecture of *T. cacao* was used to generate null
129 assumptions for the distribution of hotspots along the genome. These null distributions
130 were used to identify differential representation of known DNA sequence motifs in ubiqui-
131 tous recombination hotspots, and of overlap between recombination hotspots and genomic
132 traits for each population. The re-sampling schemes used to identify these associations are

133 novel in the context of this work and were designed to take into account the size and distri-
134 bution of elements in the genome. In this work we aimed to answer the following questions:
135 (i) How are recombination rates distributed within 10 highly differentiated populations of *T.*
136 *cacao*, and how do they compare to each other? (ii) How are hotspots distributed along the
137 genome of each of the ten populations of *T. cacao*, and can these distributions be explained
138 by patterns of population genetic differentiation? (iii) Are there identifiable DNA sequence
139 motifs that are associated with the location of recombination hotspots along the *T. cacao*
140 genome? (iv) Are there genomic features (e.g. TSSs, TTSs, exons, introns) consistently
141 associated with recombination hotspot locations across *T. cacao* populations? Our findings
142 suggest that recombination hotspot locations generally follow patterns of diversification be-
143 tween populations, while also having a strong tendency to occur close to TSSs and TTSs.
144 Moreover, we find a strong negative association between the occurrence of recombination
145 hotspots and the presence of retroelements.

146

147 Results

148 *Comparing recombination rates between populations*

149

150 Populations show a mean recombination rate r/kb between 2.1×10^{-5} and 5.25×10^{-3} (Ta-
151 ble 1), with a variety of distributions (Fig. 2). We observe a higher mean than median r/kb
152 for all populations, indicating that extreme high values are present for all populations. The
153 extreme recombination rate values affect the mean, driving it to values consistently higher
154 than the median. The pattern of recombination rates along the genome varied between pop-
155 ulations, as can be seen in the comparison of the Nanay and Purus third chromosome (Fig.
156 3). Purus appears to have a higher average recombination rate than Nanay for chromosome

157 three. More specifically, particular regions of the chromosome present peaks in one popula-
158 tion that are absent in the other. A similar patten can also be observed for the density of
159 recombination hotspots, e.g. Purus presenting a high density of hotspots in certain regions
160 that is not observed in Nanay. The median 95% probability interval for recombination rate
161 across the genome for each population was found to be several orders of magnitude larger
162 than the uncertainty per site, estimated as the median 95% Credibility Interval of the trace
163 for each position in the genome for that population (Table 2).

164 Overall, the mean recombination rate for most of the populations was higher than esti-
165 mated mutation rates for multicellular eukaryotes of 10^{-6} changes per kb per generation
166 (Lynch, 2010; Exposito-Alonso et al., 2018) (Table 1). Two populations, Guianna and Criollo,
167 were notable exceptions, having higher average recombination rates than the other popula-
168 tions by one and two orders of magnitude respectively. Guianna and Criollo also have been
169 estimated to have a lower effective population size (N_e) (Cornejo et al., 2018) by one and
170 two orders of magnitude respectively. However, there was no significant linear trend be-
171 tween mean N_e and r/kb ($p = 0.1119$), indicating that, for a high enough N_e , the ability to
172 detect recombination events is not dictated by the effective population size. When Criollo
173 and Guianna were excluded, the relationship was also not present ($p = 0.3886$). When all
174 populations were included, the inbreeding coefficient (F , from Cornejo et al., 2018) showed
175 no significant linear association with mean r/kb ($p = 0.3361$). We also found no linear trend
176 between sample size and mean r/kb ($p = 0.2333$). The average recombination rate per popu-
177 lation was transformed from r/kb to cM/Mb (Table 1) using the Kosambi mapping function
178 Kosambi (1943). The average cM/Mb was 4.6×10^{-04} .

179 In order to compare the average recombination rates (r/kb) of the different populations,
180 a Kruskal-Wallis test was performed for every pair of populations. The only pair of popu-
181 lations that did not show a significant difference in mean recombination rate was the pair
182 of Nacional and Nanay ($p = 0.3$). All other pairwise comparisons were highly significant

183 ($p < 2 \times 10^{-16}$).

184

185 *Comparing recombination hotspot locations between populations*

186

187 The majority (55.5%) of hotspots identified were not shared between populations. The 25
188 most numerous sets of hotspots are represented in Fig. 5. The nine largest of these are sets
189 of hotspots unique to single populations. The hotspots unique to the remaining population
190 (Criollo) formed the eleventh largest set. Effective population size (N_e) is not a good linear
191 predictor of the amount of detected hotspots ($p = 0.1489$), nor is sample size ($p = 0.351$).

192 The recombination rate in hotspot regions for nine of the populations was on average be-
193 tween 22 and 237% higher than the average recombination rate of the genome. The exception
194 was Guianna, which only showed an approximately 1% increase in average recombination
195 rate in hotspot regions when compared to that of the non-hotspot regions. A 1% higher
196 average recombination rate in hotspots may be due to an increased ability to detect hotspots
197 in regions of low recombination for this population. Additionally, Guianna presents unusu-
198 ally large hotspots (average 8.9 kb, Table 6), which points to an especially low resolution in
199 hotspot detection for this population.

200 Despite the majority of hotspots not being shared between populations, we conducted
201 pairwise Fisher's exact tests to verify whether there was significantly more hotspot overlap
202 than expected (if hotspots were randomly distributed along the genome) between popu-
203 lations. For most pairs of populations we found significantly more hotspot overlap than
204 expected (Table 3). There were three comparisons that did not show significantly more
205 overlap than expected: Amelonado-Nacional, Amelonado-Purus, and Criollo-Nacional. A
206 Mantel test comparing distances between populations based on shared hotspots and F_{ST}
207 values between populations resulted in a significant correlation between them ($r = 0.66$,
208 $p = 0.002$). The correlation between eigenvectors from a correlation matrix and those of the

209 genetic covariance matrix were also explored. When all populations were included, we found
210 that the first eigenvector from the genetic covariance matrix was not significantly correlated
211 with the first eigenvector from the hotspot correlation matrix ($p = 0.7055$), but the second
212 genetic eigenvector was ($p = 0.009007$, $r = 0.7711638$). However, the first eigenvector of the
213 genetic covariance matrix captured the difference between the Criollo population (the only
214 domesticated variety) and the rest of the populations. The second eigenvector explains most
215 of the natural differentiation across populations (Cornejo et al., 2018). For that reason, we
216 decided to exclude Criollo and repeat the analysis. We found that the first eigenvector from
217 the correlation matrix constructed from shared hotspot information was not significantly
218 correlated with either of the first two eigenvectors of the genetic covariance matrix when
219 Criollo was excluded (eigenvector 1: $p = 0.1314$, eigenvector2: $p = 0.3376$).

220 To study the effects of demographic history more closely, shared hotspots were converted
221 to dimensions of a multiple correspondence analysis and modeled along a previously con-
222 structed drift tree (Cornejo et al., 2018). Modeling the dimension as a Brownian motion was a
223 better fit (AIC=79.4) than modeling it as an Ornstein-Uhlenbeck (OU) process (AIC=81.4),
224 which is consistent with the small number of hotspots shared between populations. The
225 model assuming Brownian motion is consistent with pure drift driving differentiation of a
226 trait along a genealogy, while an OU process is consistent with a higher trait maintenance
227 (stabilizing selection).

228

229 *Identifying DNA sequence motifs associated with the locations of recombination hotspots*

230

231 RepeatMasker was used to analyze the set of recombination hotspots that were present
232 in at least eight *T. cacao* populations (17 total hotspots), as well as the consensus set
233 of recombination hotspots, and the reference genome. In order to determine whether a
234 particular set of DNA sequence repeats was overrepresented in the regions of ubiquitous

235 recombination hotspots, the percentage of DNA sequence that was identified as potentially
236 being from retroelements or DNA transposon was compared to an empirical distribution.
237 The percentage of observations from the distribution which were greater than the observed
238 are reported in Table 4. While retroelements were found to be underrepresented in the ubiq-
239 uitous hotspots, DNA transposons were marginally overrepresented.

240

241 *Identifying genomic features associated with the location of recombination hotspots*

242

243 An overrepresentation of recombination hotspots was found in all ten of the populations
244 at transcriptional start sites (TSSs) and transcriptional termination sites (TTSs)(Table 5).
245 The level of overrepresentation of hotspots in particular regions was compared to a null ex-
246 pectation based on simulations of hotspots of the same size as the ones detected, distributed
247 randomly along the chromosomes. For all populations, all 1000 simulations showed a lower
248 proportion of overlap with TSSs and TTSs than the observed. In the case of exons and
249 introns, seven populations (Contamana, Criollo, Iquitos, Maranon, Nacional, Nanay, Purus)
250 had an observed value that was lower than all, or almost all (Purus for exons), simulations.
251 Three of the remaining four populations (Amelonado, Curaray, and Nanay) had no clear
252 trend in either direction (Table 5). The final population (Guianna) showed an overrepresen-
253 tation of hotspots in both exons and introns.

254

255 Discussion

256 Understanding how recombination rates vary between genetically differentiated populations
257 of the same species is an important step toward disentangling the role of recombination in
258 genetic differentiation. This set of *T. cacao* populations presents a unique opportunity to

259 infer recombination in wild, long- and recently established populations, as well as a domesti-
260 cated population (Criollo) (Cornejo et al., 2018; Bartley, 2005). This system has allowed us
261 to explore differences in recombination hotspot locations between populations of the same
262 species. Our results point to a conservation of hotspots between populations that generally
263 mirrors the patterns of genetic differentiation between populations. Also, we find that TSSs
264 and TTSs are strongly associated with recombination hotspots in all populations, which is
265 consistent with previous findings in plants (Paape et al., 2012; Choi et al., 2013; Hellsten
266 et al., 2013). This factor seems to play an important role in determining the location of
267 novel hotspots. Finally, hotspots that are shared by at least eight populations appear to be
268 associated with DNA transposons, pointing to a potential mechanism for the maintenance
269 of recombination hotspots at the population-divergence time-scale.

270

271 *Comparing recombination rates between populations*

272

273 We found that the eight long-established, wild *T. cacao* populations show an average
274 recombination rate (r/kb) greater than multicellular eukaryotic mutation rates (Table 1),
275 while the other two populations (Criollo and Guianna) show unusually high average recom-
276 bination rates in comparison. Despite a small sample for some populations, we found no
277 linear trend between sample size and recombination rate. Additionally, the rates calculated
278 for the two wild, small-sample populations (Curaray and Nacional) were consistent with
279 those of other wild populations. This makes us confident in our estimates, particularly for
280 the domesticated Criollo population. For all populations, the mean recombination rate was
281 found to be greater than the median. This is consistent with high rate outlier values; an
282 expected result in the presence of recombination hotspots. Using the effective population
283 size for *Medicago truncatula* from Siol et al., 2007 and the estimate of ρ from Paape et al.,
284 2012, we calculated r/kb ($= 4 \times 10^{-3}$) and found that it was comparable with the rate found

285 for the Criollo population (Table 1). We also calculated the median recombination rate
286 in cM/Mb for each chromosome using the Kosambi mapping function (Kosambi, 1943) over
287 non-overlapping, 100 SNP windows. The average cM/Mb for all populations was 4.6×10^{-04} ,
288 which is lower than has been measured for any Malvale (Kundu et al., 2015), but not as low
289 as the lowest measured for conifers (Chen et al., 2010; Stapley et al., 2017). This places
290 *T. cacao* on the high end of known recombination rates for its order but comfortably in
291 the range of other long-lived, woody plants. Average recombination rates in cM/Mb varied
292 between populations from Amelonado (4.04×10^{-06}) to Criollo (3.91×10^{-03}). Previous work
293 has shown that Criollo is the only population showing a strong signature of domestication, as
294 revealed by much higher drift parameter than that observed for other populations (Cornejo
295 et al., 2018). Domestication has been observed to increase recombination rates, particularly
296 in plants (Ross-Ibarra, 2004), and is a possible explanation for the higher recombination
297 rate observed for the Criollo population. The high recombination rate observed in Guianna
298 can be explained in a similar way; while Guianna does not show a strong signature of do-
299 mestication, it is the most recently established population (Bartley, 2005), and it has also
300 undergone a recent bottleneck (Cornejo et al., 2018). We hypothesize from this result that
301 the Guianna population is undergoing the initial stages of domestication and its increased
302 recombination is an early indicator of this. It is possible that the high recombination rates
303 estimated for Criollo and Guianna can be explained by biases in estimation caused by errors
304 associated to small samples or low genetic variation; yet, the recombination rates for Amel-
305 onado (another population with low variation) or Purus (a population with small sample
306 size) did not present this problem. Analyses exploring mutations of putative recombination
307 suppression genes (Fernandes et al., 2018) could help disentangle the nature of this extreme
308 variation in recombination rate in the Criollo and Guianna populations.

309 Despite recombination rates for eight of the ten populations being of the same order
310 of magnitude, pairwise comparisons of average rates indicated that most populations have

311 a significantly different rate of recombination from the others. The only exception were
312 Nacional and Nanay whose average rates were not significantly different from each other.
313 These two populations, however, are not more closely related to each other than they are
314 to other populations, based on genetic differentiation (Cornejo et al., 2018). We interpret
315 this result as suggestive that their similarity is not due to genetic similarity, but some other
316 factors, e.g. epigenetics.

317 The likelihood of detecting hotspots of recombination in the genome will likely be affected
318 by the amount uncertainty in the estimates of recombination across sites or regions. Yet,
319 we have been unable to identify any study where the magnitude of the uncertainty in the
320 estimates of recombination are assessed to address this issue. We have performed careful
321 comparisons and assessed the magnitude of the uncertainty in the estimation of recombi-
322 nation rates to show that this uncertainty is several orders of magnitude smaller than the
323 variation in recombination rates across the genome (Table 2).

324

325 *Comparing recombination hotspot locations between populations*

326

327 Similarly to recombination rates, the location of recombination hotspots can be very
328 informative to questions of divergence between populations. Understanding the pattern and
329 rate of change of recombination hotspots at the population level can elucidate their role
330 in shaping genome architecture, impacting how effectively selection operates (Felsenstein,
331 1974). We found that a large proportion (55.5%) of hotspots detected are unique to a single
332 population. While we do not detect all the hotspots in these populations and not all the
333 hotspots detected are necessarily true positives, this proportion of unique hotspots can be
334 seen as an indicator that the turnover rate for hotspots is faster than the time it took the
335 10 populations to differentiate. The detection rate for LDhot is approximately 55% under
336 constant population conditions, and greater when a recent bottleneck has occurred (Auton

337 et al., 2014; Dapper and Payseur, 2017). Only two of the populations in this study (Criollo
338 and Guiana) have a known recent bottleneck (Cornejo et al., 2018). However, Criollo was
339 the only one of these two with an unusually low hotspot count (Table 6). Criollo's low
340 number of detected hotspots can be a product of its increased genome-wide recombination
341 rate, making the signal of hotspots less pronounced. The observed variability of hotspot
342 location between populations points to demographic history not being the main driver of
343 recombination hotspot location. However, the hotspots tend to appear in similar regions,
344 as demonstrated by the Fisher's exact tests (Table 3). This dichotomy can be explained
345 by considering that the proportion of the genome occupied by recombination hotspots is
346 very low, so even a small proportion of hotspots from two different populations being in the
347 same region is enough for the Fisher's exact test to recognize them as significantly similar.
348 This small but significant similarity can occur by recombination being limited in its possible
349 positioning along the genome, but not to the point of forcing hotspots to occur consistently
350 in the same locations, and thus maintaining some level of stochasticity. It is important to
351 note that our hotspots are unusually large (Table 6). This is likely a product of our low
352 sample size leading to low resolution when resolving hotspot regions.

353 Given the significant proportion of overlapping hotspots between populations, it was still
354 important to explore whether the similarities can be explained by shared genetic history. If
355 demographic history explains the evolution of hotspot location, we would expect that more
356 closely related populations would have a higher percent of overlapped hotspots. A significant
357 relationship was found between population differentiation (F_{ST}) and the distance between
358 populations based on shared hotspots (Mantel test, $r = 0.66$, $p = 0.002$). The comparison
359 between the hotspot correlation matrix and the genetic covariance matrix supports what
360 was found when comparing the hotspot correlation matrix to the F_{ST} matrix. One caveat
361 is that the first genetic eigenvector, which separated Criollo from the other populations,
362 was not correlated with the first hotspot correlation eigenvector, indicating that Criollo's

363 domestication generated a genetic pattern that deviates from the pattern of shared hotspots.
364 This indicates that, to some extent, the genetic differentiation and the location of hotspots
365 are mirroring each other, which could be due to recombination hotspots being a product of the
366 shared history between the populations. However, since recombination rates were estimated
367 using a coalescent-based method, we expect historical relationships to be represented in
368 our findings. We transformed the information of hotspot overlap to model hotspots as
369 quantitative traits changing along a population tree (Cornejo et al., 2018). Our results, show
370 that a Brownian motion model (AIC=79.4) better fits the data than a model with stabilizing
371 selection Ornstein-Uhlenbeck model (AIC=81.4) and suggest that, in principle, drift alone
372 could explain the evolution of the location of recombination hotspots. However, the absolute
373 number of hotspots that are shared among populations indicates that demographic history
374 alone is insufficient to explain the evolution of recombination hotspots in this species.

375 One conclusion that follows from these results is that, while shared recombination hotspots
376 can to some extent be explained by patterns of genetic differentiation, some of the sharing
377 can simply be due to a tendency for hotspots to arise near TSSs and TTSs. It has been ob-
378 served in other organisms that hotspots of recombination are frequently associated to specific
379 genomic features (including TSSs and TTSs) (Auton et al., 2013; Choi et al., 2013; Hellsten
380 et al., 2013; Myers et al., 2005; Singhal et al., 2015) or DNA sequence motifs (Auton et al.,
381 2012; Brunshwig et al., 2012; Stevison et al., 2016). These factors can affect the landscape
382 of recombination, contributing to the patterns of shared hotspot locations between popula-
383 tions that we are observing in *T. cacao*. Previous studies looking at apes and finches have
384 explored recombination hotspots in multiple species and as many as two populations of the
385 same species (Singhal et al., 2015; Stevison et al., 2016; Shanfelter et al., 2018), but this
386 study is the first to compare hotspots in more than two populations of the same species at
387 once. The increased number of populations allows us to analyze the relationship between
388 population genetic processes and recombination. Our results suggest that the pattern of

389 gains and losses of recombination hotspots is very dynamic and the landscape of recombina-
390 tion changes rapidly during the process of diversification within a species. This dynamism
391 can have a tremendous impact on the adaptive dynamics of a species, and it should be taken
392 into account, considering that theoretical studies tend to assume that recombination rates
393 are constant during the evolution of populations (Hudson and Kaplan, 1988; Donnelly and
394 Kurtz, 1999).

395

396 *Identifying DNA sequence motifs associated with the locations of recombination hotspots*

397

398 The analysis of 17 hotspots shared between at least eight populations of *T. cacao* found an
399 underrepresentation of retroelements and a marginal overrepresentation of DNA transposons
400 when compared to the entire genome (Table 4). These results are not entirely surprising as
401 it has already been suggested that transposable elements (TEs) tend to be enriched in ar-
402 eas of low recombination in *Drosophila* as a consequence of selection against TEs (Rizzon
403 et al., 2002). However, the marginal over-representation of DNA transposons in the most
404 conserved recombination hotspot is unexpected, given that all previous observations have
405 shown a reduced representation of mobile elements in areas with high recombination rate
406 (Rizzon et al., 2002). It is possible that DNA transposons are at least partly responsible for
407 the maintenance of recombination hotspots as populations diverge, from which we expect
408 that site-directed recombination is more frequent in these locations of the genome. However,
409 the low percentage of these sequences observed in the set of all hotspots (Table 4) indicates
410 that these sequences only have a small effect on the maintenance of hotspots. It has been
411 observed in humans that short DNA motifs enriched for repeat sequences determine the loca-
412 tion of 40 per cent of hotspots enriched for recurrent non-allelic homologous recombination
413 (McVean, 2010). One potential explanation for why natural selection does not eliminate
414 hotspots in these regions is the possibility that these regions do not produce a large enough

415 mutational load for natural selection to remove them from the population (McVean, 2010).

416

417 *Identifying genomic features associated with the location of recombination hotspots*

418

419 For all ten populations, an overrepresentation of hotspots was found in the areas im-
420 mediately preceding and following transcribed regions of the chromosome. This matches
421 the findings of previous studies in *Arabidopsis thaliana* (Choi et al., 2013), *Taenipygia gut-*
422 *tata* and *Poephila acuticauda* (Singhal et al., 2015), and humans (Myers et al., 2005). The
423 most likely explanation is that recombination events within genes are selected against. The
424 rationale being that a recombinant chromosome that undergoes a double-strand break in
425 the middle of a coding region will have a higher risk of being inviable, and therefore not
426 represented in the current set of chromosomes for its population. Recombination occurring
427 in transcription start and stop sites, on the other hand, does a much better job at break-
428 ing up haplotypes or shuffling alleles in different genomic backgrounds, while preserving the
429 functionality of coding regions. This rationale is supported by previous findings of increased
430 recombination rates in these regions (Choi et al., 2013). It is also supported by results from
431 PRDM9 knock-out *Mus musculus*, which has shown a reversion to hotspots located near
432 TSSs (Brick Kevin et al., 2012). The enrichment of *T. cacao* hotspots in TSSs and TTSs
433 is thus a reasonable result given that zinc-finger binding motifs and potential modifiers like
434 PRDM9 have not been identified in this species.

435

436 *Implications for the evolutionary history of T. cacao*

437

438 Overall, our results show a large consistent pattern where recombination rates in the ten
439 populations of *T. cacao* are of a similar magnitude as mutation rates, but show a high di-
440 versity in location and number of hotspots of recombination that cannot be explained solely

441 by the process of diversification of the populations. In fact, the results are indicative of the
442 turnover rate of hotspots being faster than the process of divergence among populations.
443 A potential hypothesis that could explain the rapid turnover of hotspots of recombination
444 and the relative differences in recombination among populations is that epigenetic changes
445 are involved in controlling the turnover of recombination in plants. This hypothesis is not
446 unreasonable given the recent observation of epigenetic control of recombination in plants
447 (Yelina et al., 2015). Further theoretical and simulation work should be done in order to bet-
448 ter understand the implications of the rapidly changing recombination hotspots in adaptive
449 dynamics. We also show that there is an overall underrepresentation of hotspots in exons
450 and introns for most populations, which is consistent with purifying selection acting against
451 changes that could result in disruptions of gene function. On the other hand, we observed an
452 overrepresentation of hotspots in TTSs and TSSs for all ten populations. This could impact
453 the maintenance and spread of beneficial traits in the population by shuffling allelic variants
454 of genes without causing disruption of their function. We hypothesize that the enrichment of
455 hotspots of recombination in TTSs and TSSs can have an important impact in the spread of
456 beneficial mutations across different genomic backgrounds; increasing the rate of adaptation
457 to selective pressures (e.g. selection for improved pathogen response).

458

459 **Materials and Methods**

460 *Comparing recombination rates between populations*

461

462 Sequence data were downloaded from the Cacao Genome Database and NCBI (Accession
463 PRJNA486011), including the reference sequence for each chromosome and the full genome
464 annotation (*Theobroma cacao* cv. Matina 1-6 v1.1)(Motamayor et al., 2013). Processing

465 was done using the pipeline from (Cornejo et al., 2018) available at the github repository
466 *oeco28/Cacao_Genomics*. Full genome data was used from a total of 73 individuals across
467 10 populations (Cornejo et al., 2018): Criollo (N = 4, #SNPs = 309,818), Curaray (N =
468 5, #SNPs = 1,106,871), Contamana (N = 9, #SNPs = 2,097,618), Amelonado (N = 11,
469 #SNPs = 373,789), Maranon (N = 14, #SNPs = 1,783,226), Guianna (N = 9, #SNPs =
470 770,729), Iquitos (N = 7, #SNPs = 1,575,711), Purus (N = 6, #SNPs = 1,184,181), Nanay
471 (N = 10, #SNPs = 830,885), and Nacional (N = 4, #SNPs = 718,099). We filtered the single
472 nucleotide polymorphism data and excluded rare variants (minor allele frequency ≤ 0.05)
473 per population. Separate variant files per population per chromosome were then phased
474 using default conditions with SHAPEIT2 (Delaneau et al., 2011) under default parameters.
475 Haplotype files were converted back to phased variant calling format (vcf) for its downstream
476 analysis. We have also phased the data with Beagle (Browning and Browning, 2007), using
477 a burnin of 10000 iterations, and estimations done over 10000 iterations. No appreciable
478 differences were observed between the two methods and Beagle phasing was maintained for
479 the analyses. The reason for performing the phasing separately for each population is that
480 linkage disequilibrium patterns are expected to be affected by population structure. The ten
481 populations have been shown to be unique clusters with very little admixture between them
482 Cornejo et al. (2018), and the individuals used in this study were those whose ancestry was
483 clearly from a single population. VCFTools (Danecek Petr et al., 2011) was used to remove
484 all singletons and doubletons. Only bi-allelic single nucleotide polymorphisms (SNPs) were
485 retained and were exported in LDhat format.

486 In order to estimate recombination rates we used the *interval* routine of LDhat (Auton
487 and McVean, 2007), a program that implements coalescent resampling methods to estimate
488 historical recombination rates from SNP data. To reduce computation time, each chromo-
489 some was split into windows, each containing 2000 SNPs. To counteract the overestimation
490 of recombination rate produced at the ends of the windows, an overlap of 500 SNPs was left

491 between consecutive windows. The final window for each chromosome did not always match
492 the general scheme, so the final 2000 SNPs were taken (making the overlap with the second
493 to last window variable, but never less than 500 SNPs) (Fig. 6). Once these windows were
494 generated, LDhat was run over each window using 100 million iterations, sampling every
495 10000 iterations (10000 total points sampled), with a block penalty of 5. Lookup tables
496 with a grid of 100 points, a population mutation rate parameter (θ) of 0.1 and a number
497 of sequences (n) of 50 were used for all populations. We used the same θ for all popula-
498 tions since estimates from Cornejo et al. (2018) ranged from $\pi = 0.27\%$ to $\pi = 0.37\%$, all
499 comfortably within an order of magnitude of each other. The first 50 million iterations were
500 discarded as burn-in. Once recombination rates were calculated, 250 positions were cut off
501 from both windows involved in each overlap, so that the estimates for the first half of the
502 overlap was taken from the end of the preceding window and the estimates for the second
503 half of the overlap were taken from the beginning of the following window. The final overlap
504 in each chromosome was split in order to remove 250 SNPs from the second to last window,
505 regardless of the remaining size of the last window. The remaining rate estimates were then
506 merged in order to obtain recombination rates for the entire chromosome. This was done for
507 each chromosome of each population.

508 The estimation of recombination rates with LDhat is approximated using a sampling
509 scheme with a Markov Chain Monte Carlo (MCMC) algorithm as implemented in the *interval*
510 routine. The inference of recombination rates is the result of the integration of estimated
511 parameter values across iterations with the routine *stats*. In the majority of recent studies
512 where LDhat or LDhelmet are used (Myers et al., 2005; Auton et al., 2012; Brunshwig et al.,
513 2012; Paape et al., 2012; Auton et al., 2013; Choi et al., 2013; Singhal et al., 2015; Stevison
514 et al., 2016), whether there is convergence of the Markov chains has not been explicitly
515 investigated. One study that we are aware of has used simulations to asses whether their
516 small sample size affected their ability to obtain reliable estimates of recombination using

517 LDhelmet (Booker et al., 2017), but did not assess the uncertainty of the estimates from
518 the MCMC process itself. We argue that evaluation of convergence is important to assess
519 the confidence in the estimated reported values, especially if there is interest in analyzing
520 the differences in recombination rate along the genome. Visual inspection of pilot runs of
521 the analysis demonstrated that convergence was not achieved after running 40M iterations,
522 which is why the length of the chains was increased to 100M iterations. Additionally we
523 explored the uncertainty in the estimates of recombination site-wise by integrating over the
524 trace of the estimates for recombination rate to infer the 95% Credibility Interval. We then
525 estimated the 95% interval of recombination estimates range across all sites in the genome
526 to have an overall measure of uncertainty that we compared to the median 95% Credibility
527 Interval for the trace of each position.

528 In order to compare recombination rates, the effective population size (N_e) calculated for
529 each population (Cornejo et al., 2018) was used to convert rates in $N_e r/kb$ to r/kb . Differ-
530 ences in the mean genome-wide recombination rate between populations were then tested
531 using the Kruskal-Wallis test (`kruskal.test` function from the `stats` package in R) (R Core
532 Team, 2018). There were 45 comparisons, making the Bonferroni correction cutoff value:
533 $\alpha = 0.0011$. To transform per population recombination rates from r/kb to cM/Mb , we
534 divided each chromosome into windows of 100 SNPs and used the Kosambi mapping func-
535 tion (Kosambi, 1943). The median for the windows of a chromosome was then calculated,
536 and the average of each population's chromosomes was taken as that population's average
537 recombination rate in cM/Mb .

538

539 *Comparing recombination hotspot locations between populations*

540

541 Recombination hotspots were estimated with LDhot (Auton and McVean, 2007), a likelihood-
542 based program that tests whether a single distribution model or a two distribution model

543 better explains the observed recombination rates in 1 kb sliding windows (default), for each
544 chromosome. Each chromosome was run in its entirety, with the number of simulations
545 (nsims) set to 1000. The resulting potential hotspots were refined by an alpha of 0.001, and
546 overlapping hotspots were merged. This method therefore detects hotspots by comparing
547 rates in 1 kb windows to the rates in the surrounding regions.

548 To determine the set of consensus hotspots, the hotspots from all populations were
549 merged. Two hotspots from different populations were considered to be shared if they both
550 overlapped with the same hotspot in the consensus set. To summarize all shared hotspots, a
551 Boolean matrix was constructed, in which a population having a hotspot that overlaps with
552 a hotspot in the consensus list leads to an indication of presence of the consensus hotspot
553 in that population. This matrix was used to determine hotspots shared by two or more
554 populations.

555 A Fisher's exact test was run for each pair of populations in order to determine whether
556 hotspots for the pair of populations overlap significantly more than expected. The BED
557 files containing the location of the recombination hotspots for each pair of populations were
558 compared using `Bedtools:fisher` (Quinlan and Hall, 2010). The number of comparisons was
559 45, making the the Bonferroni correction cutoff value: $\alpha = 0.0011$.

560 In order to compare the relationships between populations based on shared hotspots we
561 calculated Jaccard distances (`distance` function, `phylentropy` package, R) (Drost, 2018)
562 and compared them to a published F_{ST} matrix (Cornejo et al., 2018) using a Mantel test
563 (`mantel.rtest` function, `ade4` package, R) (Chessel et al., 2004; Dray and Dufour, 2007;
564 Dray et al., 2007; Bougeard and Dray, 2018). The F_{ST} estimates from Cornejo et al. (2018)
565 were generated using Weir and Cockerham's estimator Weir and Cockerham (1984).

566 The Boolean matrix for shared hotspots was also used to explore the relationship be-
567 tween hotspot similarities and genetic covariances from a previous study (Cornejo et al.,
568 2018). Singletons were removed from the hotspot matrix, which was converted to a corre-

569 lation matrix using the `mixed.cor` function from the `psych` package in R (Revelle, 2018).
570 The `mixed.cor` function was used due to its ability to calculate Pearson correlations from
571 dichotomous data. We then used the `eigen` function in R (R Core Team, 2018) to generate
572 eigenvectors for the hotspot correlation matrix and the genetic covariance matrix. Pearson
573 correlations between the first and second eigenvector of the genetic covariance matrix and the
574 hotspot correlation matrix were then calculated (`cor.test` function, `stats` package, R)(R
575 Core Team, 2018). This analysis was done once with all populations included, and once with
576 the Criollo population excluded before correlations were calculated.

577 In order to model the presence or absence of hotspots along a drift tree, a multiple
578 correspondance analysis was used on the Boolean matrix of shared hotspots using the `MCA`
579 function from the `FactoMineR` package in R Lê et al. (2008). Nine dimensions were retained
580 and used as traits along a previously generated drift tree (Cornejo et al., 2018). Using the
581 `Rphylopars` package in R (Goolsby et al., 2016), the dimensions were modeled as Brownian
582 motion and as an Ornstein-Uhlenbeck process. The fit of the two models were compared
583 using the AIC values for the best fitting models of each type.

584

585 *Identifying DNA sequence motifs associated with the locations of recombination hotspots*

586

587 Motifs associated with hotspots were found using `RepeatMasker` (Smith et al., 2016).
588 The entire genome, the set of consensus hotspots, and a set of ubiquitous hotspots (hotspots
589 shared by at least eight of the populations) were examined with `RepeatMasker`, using normal
590 speed and "theobroma cacao" in the species option. In order to determine whether ubiqui-
591 tous hotspots were enriched for particular DNA sequences, a set of the same number and size
592 of sequences was randomly selected from the genome using `Bedtools:shuffle` (Quinlan and
593 Hall, 2010) and examined with `RepeatMasker`. This simulation was repeated one thousand
594 times and a null distribution against which observed values were compared was constructed

595 from the results.

596

597 *Identifying genomic features associated with the location of recombination hotspots*

598

599 Testing whether recombination hotspots were overrepresented near particular genomic
600 features was done by using a resampling scheme to establish null expectations and then
601 comparing the observed value to the empirical distribution. For each feature, locations were
602 retrieved and the number of observed hotspots that overlap with this feature were counted.
603 To determine whether this amount of overlapping hotspots was unusually high or low, a set of
604 hotspots that matched the number of hotspots and the size of each hotspot was simulated.
605 These simulated hotspots were placed randomly along the chromosome, using a uniform
606 distribution. The simulation was run 1000 times and the number of simulated hotspots that
607 overlap with the true genomic features was measured for each simulation. The simulations
608 generate an expected distribution of overlap with the genomic feature, and the true value
609 was then compared to the distribution. When simulated hotspots overlapped, the location
610 of one of them was sampled again. Features tested were: Transcriptional start sites (TSSs),
611 transcriptional termination sites (TTSs), exons, and introns. TSSs and TTSs are considered
612 to be the 500bp upstream and downstream of coding regions respectively.

613 The reason for the proposed novel resampling scheme is that, if the size and distribution
614 of genomic features and hotspots were not taken into account, it would set unrealistic expect-
615 tations for the overlap between features under a null model of no association. In this sense,
616 the null model would be inappropriate and potentially inflate the false positive rate.

617

618 *Data and code availability*

619

620 Rate and summary files from LDhat runs as well as hotspots for each population will

621 be placed in a Dryad repository. Scripts for LDhat and LDhot runs as well as the re-
622 sampling schemes used and additional analysis is available in the following github repository
623 *ejschwarzkopf/recombination – map*.

624

625 **1 Acknowledgments**

626 The authors would like to thank the Noe Higinbotham endowment and the WSU College
627 of Arts and Science for travel funds to EJS to present earlier versions of this work. We
628 would like to thank the Kamiak High Performance Computing Cluster at WSU for the in-
629 frastructure support to run the analyses, and the Cornejo, Kelley, and Busch labs at WSU
630 for feedback and edits on the manuscript.

631

632 **Tables**

Population	Mean $4N_e r/kb$	Mean N_e	Mean r/kb	Median r/kb	Lower Bound (Mean r/kb)	Upper Bound (Mean r/kb)	Mean cM/Mb
Amelonado	1.58	15744	2.51e-05	2.40e-09	2.48e-05	2.54e-05	4.04e-06
Contamana	8.53	61102	3.49e-05	4.92e-06	3.48e-05	3.50e-05	7.74e-05
Criollo	14.60	695	5.25e-03	4.27e-03	5.23e-03	5.27e-03	3.91e-03
Curaray	10.36	58213	4.45e-05	1.78e-05	4.44e-05	4.46e-05	1.18e-04
Guianna	8.66	4651	4.65e-04	7.74e-06	4.63e-04	4.67e-04	2.74e-04
Iquitos	4.23	49984	2.11e-05	5.88e-09	2.10e-05	2.12e-05	1.84e-05
Maranon	4.09	34037	3.01e-05	1.64e-08	2.99e-05	3.02e-05	1.68e-05
Nacional	4.66	26060	4.47e-05	9.76e-08	4.44e-05	4.49e-05	4.10e-05
Nanay	6.82	42429	4.02e-05	1.51e-07	4.00e-05	4.04e-05	1.33e-05
Purus	5.95	17357	8.57e-05	7.74e-06	8.54e-05	8.60e-05	1.23e-04

Table 1: Recombination rates in $4N_e r/kb$, r/kb , and cM/Mb for all ten *T. cacao* populations. The N_e that was used for the transformation is also reported for each population, as are the lower and upper bounds of a 95% confidence interval for r/kb .

633

634

635

636

637

Population	Position L95	Position U95	Genome L95	Genome U95	Position Range Quotient	Genome Range Quotient
Amelonado	6.35e-10	7.67e-08	2.33e-10	3.13e-04	120.75	1.34e+06
Contamana	1.63e-06	1.34e-05	1.40e-09	2.64e-04	8.22	1.88e+05
Criollo	8.75e-04	4.31e-03	5.35e-07	1.66e-02	4.92	3.11e+04
Curaray	5.15e-06	2.63e-05	2.72e-09	2.02e-04	5.11	7.40e+04
Guianna	9.76e-06	1.26e-04	1.81e-09	2.96e-03	12.90	1.63e+06
Iquitos	3.50e-10	6.52e-08	1.58e-10	2.45e-04	186.29	1.55e+06
Maranon	1.98e-09	7.40e-07	2.31e-10	3.52e-04	373.35	1.52e+06
Nacional	4.80e-10	6.50e-08	2.76e-10	3.66e-04	135.60	1.33e+06
Nanay	1.65e-09	3.06e-07	2.06e-10	3.52e-04	185.32	1.71e+06
Purus	3.98e-08	5.24e-06	2.00e-09	6.35e-04	131.87	3.18e+05

Table 2: The median of the upper and lower bounds of the 95% Credibility Interval for the trace of estimates of r from all positions in the genome are presented for each population (i.e. Position L95 and Position U95). The upper and lower bounds of the 95% probability interval for the median estimate of r for each population is also presented (i.e. Genome L95 and Genome U95). The quotients of the upper and lower bounds for each of the two intervals point to a much larger genome-wide variation in r than per-position variation in the trace for the estimate of r .

Population	Ame	Con	Cri	Cur	Gui	Iqu	Mar	Nac	Nan
Amelonado	-	-	-	-	-	-	-	-	-
Contamana	<2e-07	-	-	-	-	-	-	-	-
Criollo	<9e-05	<5e-13	-	-	-	-	-	-	-
Curaray	<3e-05	<3e-37	<5e-08	-	-	-	-	-	-
Guianna	<3e-06	<1e-37	<7e-07	<4e-20	-	-	-	-	-
Iquitos	<4e-08	<6e-87	<2e-11	<3e-16	<2e-29	-	-	-	-
Maranon	<6e-13	<7e-77	<2e-11	<2e-20	<5e-33	<4e-64	-	-	-
Nacional	0.0015	<2e-43	0.0212	<7e-14	<3e-06	<6e-14	<3e-13	-	-
Nanay	0.0004	<2e-44	<9e-11	<4e-16	<2e-21	<3e-39	<2e-38	<9e-06	-
Purus	0.1782	<4e-117	<2e-05	<2e-29	<1e-33	<2e-39	<8e-43	<6e-27	<2e-21

Table 3: Fisher’s exact test p-values for pairwise comparisons of recombination hotspot locations between populations of *T. cacao*

Measures	Observed % ubiquitous HS	Observed % all HS	Observed % whole genome	Mean % Sim	% Sim >ubiquitous HS
Retroelements	2.34	9.45	11.12	11.11	99.9
DNA transposons	1.94	1.64	1.10	1.10	5.4
Total	4.28	11.09	12.21	12.22	99.7

Table 4: Percentage of DNA sequences identified as either retroelements or DNA transposons, and total interspersed repeats. Observed values for the entire *T. cacao* genome, for all recombination hotspots (HS), and ubiquitous hotspots (hotspots in the same location in at least eight different populations). Also presented are mean percentage of these sequences for 1000 simulations of hotspots equivalent in size and count as the ubiquitous set and the percentile at which the observed value for the ubiquitous set is found in the distribution of the simulated set (Sim).

	TSSs (500bp)	TTSs (500bp)	Exon	Intron
Amelonado	1	1	0.602	0.527
Contamana	1	1	0.000	0.000
Criollo	1	1	0.000	0.000
Curaray	1	1	0.346	0.058
Guiana	1	1	1.000	1.000
Iquitos	1	1	0.000	0.000
Maranon	1	1	0.000	0.000
Nacional	1	1	0.000	0.000
Nanay	1	1	0.027	0.237
Purus	1	1	0.004	0.000

Table 5: Proportion of simulated chromosomes that presented a lower amount of hotspots intersecting with TSSs, TTSs, exons, and introns than the observed chromosomes. TSSs and TTSs are considered to be the 500bp upstream and downstream of transcribed regions, respectively.

Population	Mean Hotspot Size (kb)	Hotspot Count
Amelonado	6.9	1324
Contamana	6.1	5184
Criollo	6.1	887
Curaray	5.8	2303
Guianna	8.6	3655
Iquitos	7.0	3258
Maranon	6.8	3296
Nacional	6.9	2202
Nanay	7.6	3818
Purus	6.3	3972
Average	6.9	2989.9

Table 6: Average hotspot size (in kb) and count for hotspots detected in each population and average for all populations.

638 **Figures**

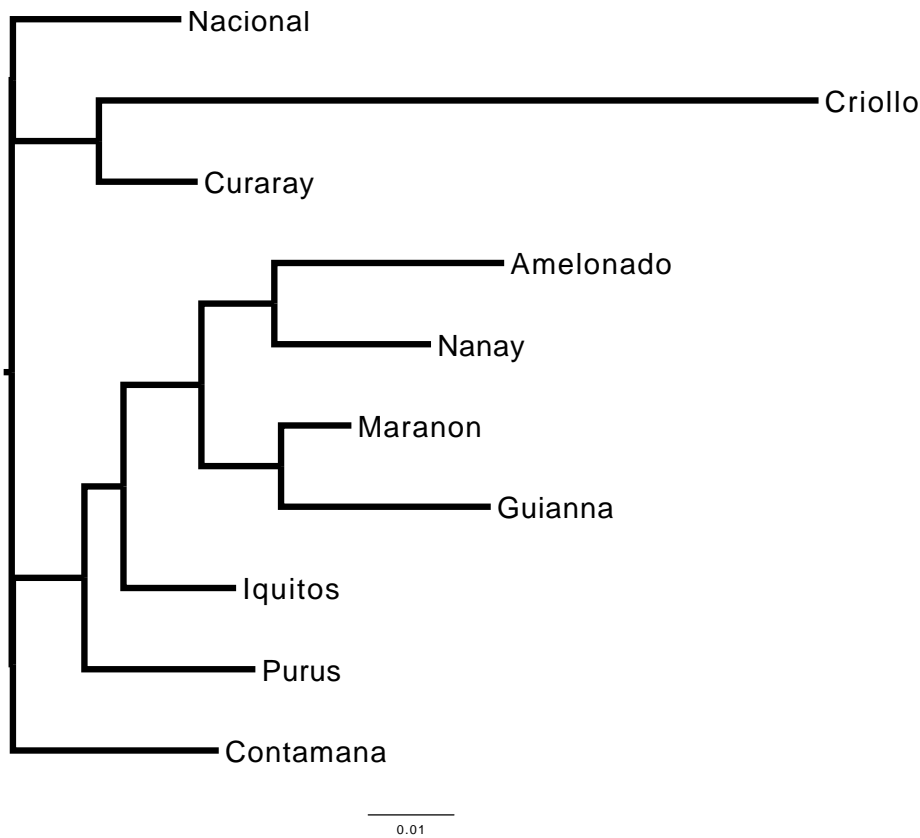


Figure 1: Drift tree constructed using TreeMix (Pickrell and Pritchard, 2012) for the 10 *T. cacao* populations. Distances between populations are based on the drift parameter. Modified from Cornejo et al. (2018)

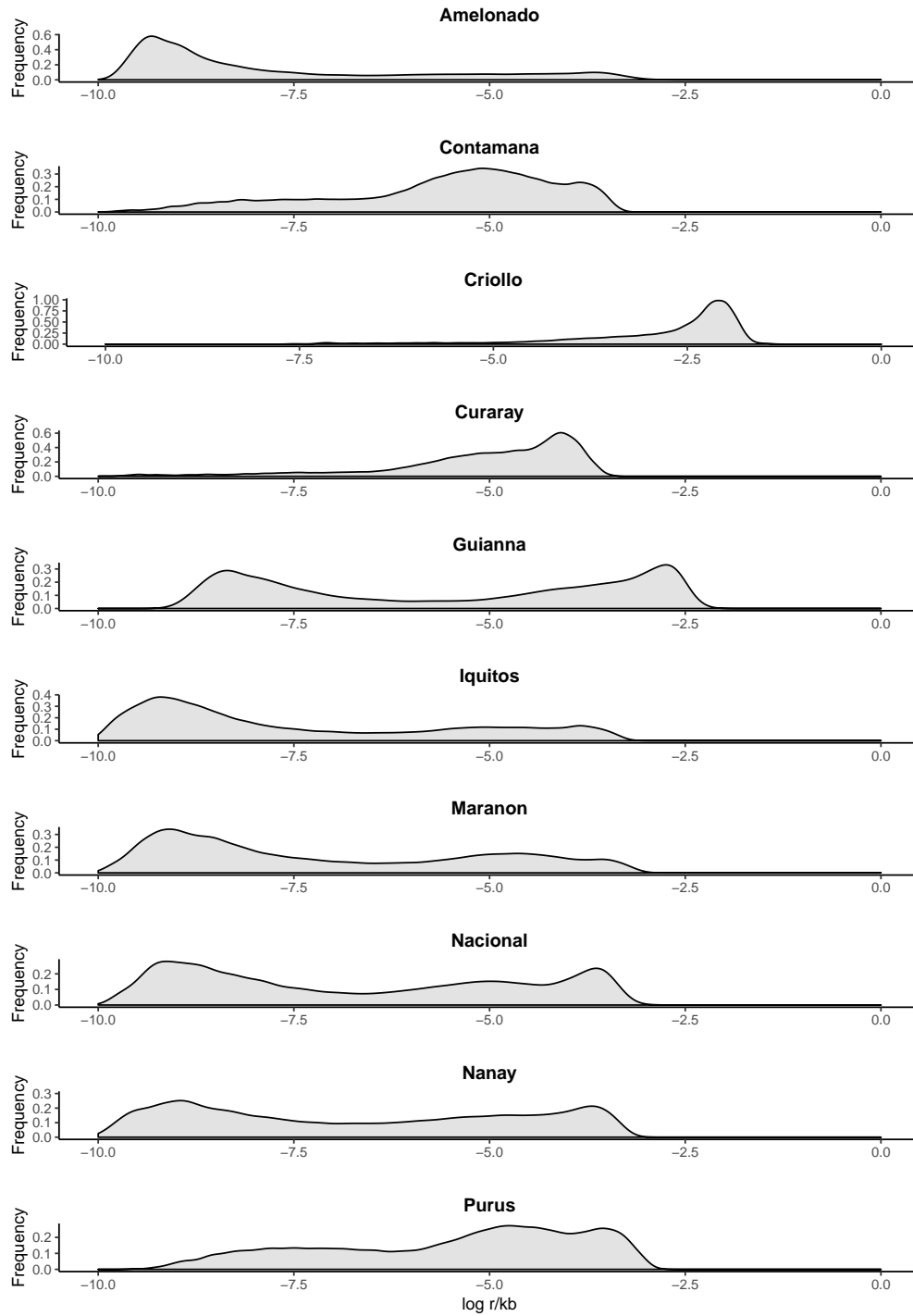


Figure 2: Distribution of \log_{10} recombination rates ($\log(r/kb)$) along the genomes of the ten *T. cacao* populations.

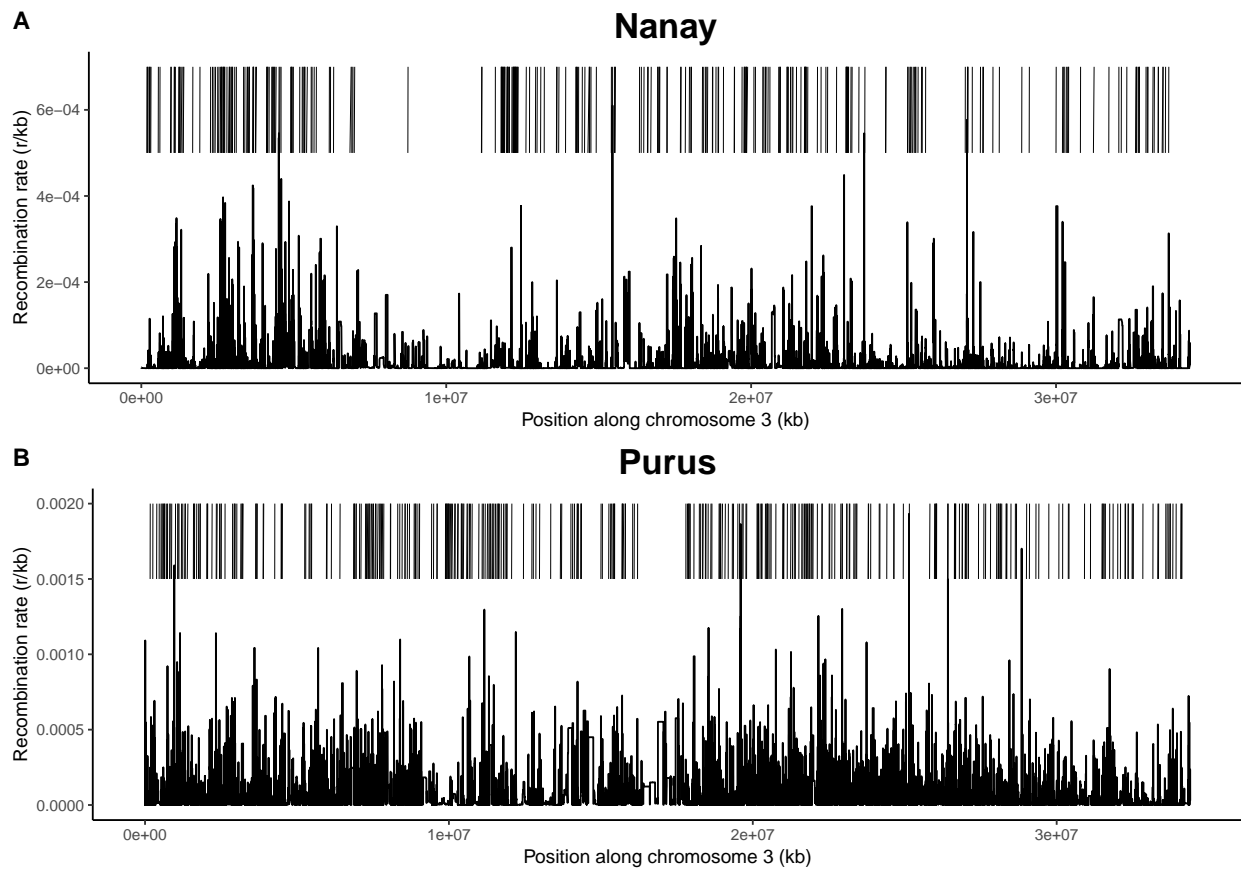


Figure 3: The third chromosomes of the Nanay (A) and Purus (B) populations were selected to exemplify the differences between populations in recombination rates (r/kb) and recombination hotspot locations (vertical bars above rates).

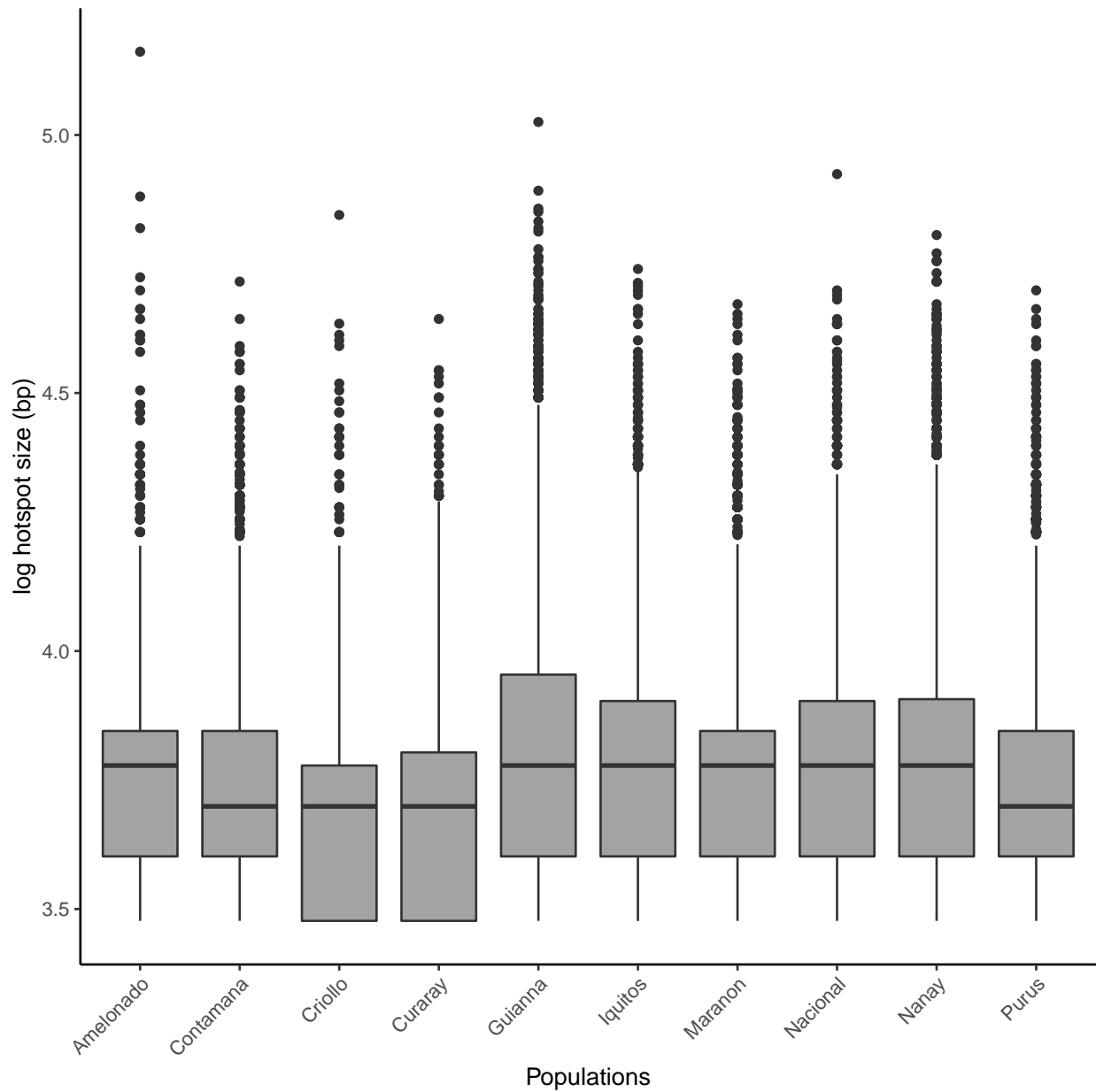


Figure 4: Boxplots of recombination hotspot sizes ($\log_{10}(\text{bp})$) by population. The horizontal line in the box represents the median value, while the points represent potential outliers.

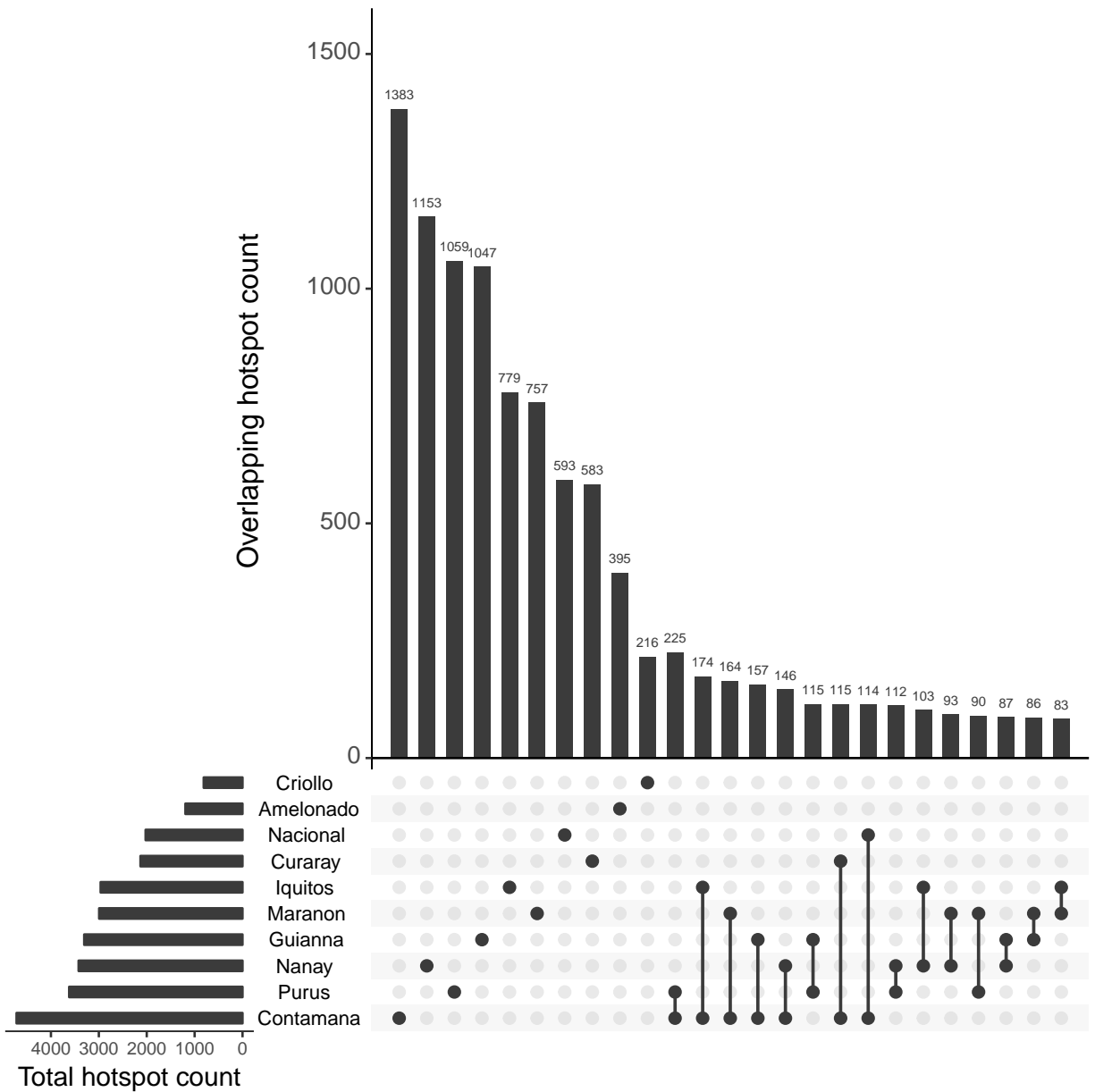


Figure 5: Upset plot showing number of hotspots in different subsets. Horizontal bars represent total hotspots detected in a population, each dot on the matrix indicate that the vertical bar above it is the count of hotspots unique to that population, connected dots indicate that the vertical bar above them represents hotspots shared between the populations represented by the connected dots. The 25 largest subsets are shown.

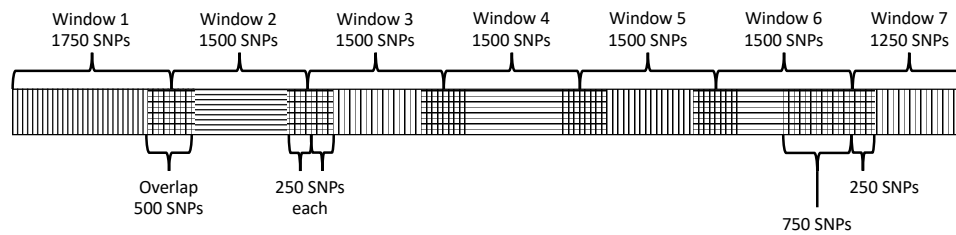


Figure 6: Example of the window layout for a 10,750 SNP chromosome. The 2,000 SNP long windows are represented by alternating horizontal and vertical lines and the overlaps between them are represented by square crosshatches. Braces above the chromosome indicate the regions from which recombination rates are extracted to generate the chromosome-wide recombination rates.

References

- 639
- 640 Akhunov, E. D., and Coauthors, 2003: The organization and rate of evolution of wheat
641 genomes are correlated with recombination rates along chromosome arms. *Genome re-*
642 *search*, **13** (5).
- 643 Anderson, L. K., N. Salameh, H. W. Bass, L. C. Harper, W. Z. Cande, G. Weber, and S. M.
644 Stack, 2004: Integrating genetic linkage maps with pachytene chromosome structure in
645 maize. *Genetics*, **166** (4), 1923–1933, doi:10.1534/genetics.166.4.1923, URL [http://www.](http://www.genetics.org/content/166/4/1923)
646 [genetics.org/content/166/4/1923](http://www.genetics.org/content/166/4/1923), <http://www.genetics.org/content/166/4/1923.full.pdf>.
- 647 Auton, A., and G. McVean, 2007: Recombination rate estimation in the presence of hotspots.
648 *Genome Research*, **17** (8), 1219–1227, URL [http://www.ncbi.nlm.nih.gov/pmc/articles/](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC1933511/)
649 [PMC1933511/](http://www.ncbi.nlm.nih.gov/pmc/articles/PMC1933511/).
- 650 Auton, A., S. Myers, and G. McVean, 2014: Identifying recombination hotspots using pop-
651 ulation genetic data. *arXiv preprint arXiv:1403.4264*.
- 652 Auton, A., and Coauthors, 2012: A Fine-Scale Chimpanzee Genetic Map from Popula-
653 tion Sequencing. *Science*, **336** (6078), 193–198, doi:10.1126/science.1216872, URL [http:](http://science.sciencemag.org/content/336/6078/193)
654 [//science.sciencemag.org/content/336/6078/193](http://science.sciencemag.org/content/336/6078/193), [http://science.sciencemag.org/content/](http://science.sciencemag.org/content/336/6078/193.full.pdf)
655 [336/6078/193.full.pdf](http://science.sciencemag.org/content/336/6078/193.full.pdf).
- 656 Auton, A., and Coauthors, 2013: Genetic recombination is targeted towards gene promoter
657 regions in dogs. *PLOS Genetics*, **9** (12), 1–9, doi:10.1371/journal.pgen.1003984, URL
658 <https://doi.org/10.1371/journal.pgen.1003984>.
- 659 Bartley, B., 2005: *The Genetic Diversity of Cacao and Its Utilization*. CAB books, CABI,
660 URL https://books.google.com/books?id=_I40iGVJD64C.

- 661 Begun, D. J., and C. F. Aquadro, 1992: Levels of naturally occurring dna polymorphism
662 correlate with recombination rates in *d. melanogaster*. *Nature*, **356** (6369).
- 663 Booker, T. R., R. W. Ness, and P. D. Keightley, 2017: The recombination landscape in
664 wild house mice inferred using population genomic data. *Genetics*, **207** (1), 297–309,
665 doi:10.1534/genetics.117.300063, URL <http://www.genetics.org/content/207/1/297>, <http://www.genetics.org/content/207/1/297.full.pdf>.
- 667 Bougeard, S., and S. Dray, 2018: Supervised multiblock analysis in R with the ade4 package.
668 *Journal of Statistical Software*, **86** (1), 1–17, doi:10.18637/jss.v086.i01.
- 669 Branca, A., and Coauthors, 2011: Whole-genome nucleotide diversity, recombination, and
670 linkage disequilibrium in the model legume *medicago truncatula*. *Proceedings of the Na-*
671 *tional Academy of Sciences*, **108** (42).
- 672 Brick Kevin, Smagulova Fatima, Khil Pavel, Camerini-Otero R. Daniel, and Petukhova
673 Galina V., 2012: Genetic recombination is directed away from functional ge-
674 nomic elements in mice. *Nature*, **485** (7400), 642–645, doi:[http://dx.doi.org/](http://dx.doi.org/10.1038/nature11089)
675 [10.1038/nature11089](http://dx.doi.org/10.1038/nature11089), URL [http://www.nature.com/nature/journal/v485/n7400/abs/](http://www.nature.com/nature/journal/v485/n7400/abs/nature11089.html#supplementary-information)
676 [nature11089.html#supplementary-information](http://www.nature.com/nature/journal/v485/n7400/abs/nature11089.html#supplementary-information), [10.1038/nature11089](http://dx.doi.org/10.1038/nature11089).
- 677 Browning, S. R., and B. L. Browning, 2007: Rapid and accurate haplotype phasing and
678 missing-data inference for whole-genome association studies by use of localized haplo-
679 type clustering. *The American Journal of Human Genetics*, **81** (5), 1084 – 1097, doi:
680 <https://doi.org/10.1086/521987>, URL [http://www.sciencedirect.com/science/article/pii/](http://www.sciencedirect.com/science/article/pii/S0002929707638828)
681 [S0002929707638828](http://www.sciencedirect.com/science/article/pii/S0002929707638828).
- 682 Brunshwig, H., L. Levi, E. Ben-David, R. W. Williams, B. Yakir, and S. Shifman, 2012:
683 Fine-scale Map of Recombination Rates and Hotspots in the Mouse Genome. *Genet-*
684 *ics*, doi:10.1534/genetics.112.141036, URL <http://www.genetics.org/content/early/2012/>

685 05/04/genetics.112.141036, [http://www.genetics.org/content/early/2012/05/04/genetics.](http://www.genetics.org/content/early/2012/05/04/genetics.112.141036.full.pdf)
686 112.141036.full.pdf.

687 Chen, M.-M., F. Feng, X. Sui, M.-H. Li, D. Zhao, and S. Han, 2010: Construction of
688 a framework map for *Pinus koraiensis* Sieb. et Zucc. using SRAP, SSR and ISSR markers.
689 *Trees*, **24** (4), 685–693, doi:10.1007/s00468-010-0438-5, URL [https://doi.org/10.1007/](https://doi.org/10.1007/s00468-010-0438-5)
690 s00468-010-0438-5.

691 Chessel, D., A.-B. Dufour, and J. Thioulouse, 2004: The ade4 package – I: One-table meth-
692 ods. *R News*, **4** (1), 5–10, URL <https://cran.r-project.org/doc/Rnews/>.

693 Choi, K., and Coauthors, 2013: *Arabidopsis* meiotic crossover hot spots overlap with
694 h2a.z nucleosomes at gene promoters. *Nature Genetics*, **45** (11), 1327–1336, doi:
695 10.1038/ng.2766.

696 Cornejo, O. E., and Coauthors, 2018: Population genomic analyses of the chocolate tree,
697 *Theobroma cacao* L., provide insights into its domestication process. *Communications*
698 *Biology*, doi:10.1038/s42003-018-0168-6.

699 Crow, J. F., and M. Kimura, 1970: (Harper international editions.). An introduction to
700 population genetics theory. Harper & Row, New York.

701 Danecek Petr, and Coauthors, 2011: The variant call format and VCFtools. *Bioinformatics*,
702 **27** (15), 2156–2158, URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3137218/>.

703 Dapper, A. L., and B. A. Payseur, 2017: Effects of Demographic History on the Detection of
704 Recombination Hotspots from Linkage Disequilibrium. *Molecular Biology and Evolution*,
705 **35** (2), 335–353, doi:10.1093/molbev/msx272, URL [https://doi.org/10.1093/molbev/](https://doi.org/10.1093/molbev/msx272)
706 msx272, <http://oup.prod.sis.lan/mbe/article-pdf/35/2/335/24367711/msx272.pdf>.

- 707 Delaneau, O., J. Marchini, and J.-F. Zagury, 2011: A linear complexity phasing method for
708 thousands of genomes. *Nature Methods*, **9** (2).
- 709 Donnelly, P., and T. G. Kurtz, 1999: Genealogical processes for fleming-viot models with
710 selection and recombination. *Ann. Appl. Probab.*, **9** (4), 1091–1148, doi:10.1214/aoap/
711 1029962866, URL <https://doi.org/10.1214/aoap/1029962866>.
- 712 Dray, S., and A.-B. Dufour, 2007: The ade4 package: Implementing the duality diagram for
713 ecologists. *Journal of Statistical Software*, **22** (4), 1–20, doi:10.18637/jss.v022.i04.
- 714 Dray, S., A.-B. Dufour, and D. Chessel, 2007: The ade4 package – II: Two-table and K-table
715 methods. *R News*, **7** (2), 47–52, URL <https://cran.r-project.org/doc/Rnews/>.
- 716 Drost, H.-G., 2018: *philentropy: Similarity and Distance Quantification Between Probabil-*
717 *ity Functions*. URL <https://CRAN.R-project.org/package=philentropy>, r package version
718 0.2.0.
- 719 Exposito-Alonso, M., and Coauthors, 2018: The rate and potential relevance of new mu-
720 tations in a colonizing plant lineage. *PLOS Genetics*, **14** (2), 1–21, doi:10.1371/journal.
721 pgen.1007155, URL <https://doi.org/10.1371/journal.pgen.1007155>.
- 722 Eyre-Walker, A., and P. D. Keightley, 2007: The distribution of fitness effects of new muta-
723 tions. *Nature Reviews Genetics*, **8** (8).
- 724 Felsenstein, J., 1974: The evolutionary advantage of recombination. *Genetics*, **78** (2),
725 737–756, URL <http://www.genetics.org/content/78/2/737>, [http://www.genetics.org/
726 content/78/2/737.full.pdf](http://www.genetics.org/content/78/2/737.full.pdf).
- 727 Fernandes, J. B., M. Séguéla-Arnaud, C. Larchevêque, A. H. Lloyd, and R. Mercier, 2018:
728 Unleashing meiotic crossovers in hybrid plants. *Proceedings of the National Academy of Sci-*

729 *ences*, **115 (10)**, 2431–2436, doi:10.1073/pnas.1713078114, URL [http://www.pnas.org/](http://www.pnas.org/content/115/10/2431)
730 [content/115/10/2431](http://www.pnas.org/content/115/10/2431.full.pdf), <http://www.pnas.org/content/115/10/2431.full.pdf>.

731 Goolsby, E. W., J. Bruggeman, and C. Ane, 2016: *Rphylopars: Phylogenetic Comparative*
732 *Tools for Missing Data and Within-Species Variation*. URL [https://CRAN.R-project.org/](https://CRAN.R-project.org/package=Rphylopars)
733 [package=Rphylopars](https://CRAN.R-project.org/package=Rphylopars), r package version 0.2.9.

734 Gore, M. A., and Coauthors, 2009: A first-generation haplotype map of maize. *Science*,
735 **326 (5956)**, 1115–1117, doi:10.1126/science.1177837, URL [http://science.sciencemag.](http://science.sciencemag.org/content/326/5956/1115)
736 [org/content/326/5956/1115](http://science.sciencemag.org/content/326/5956/1115.full.pdf), [http://science.sciencemag.org/content/326/5956/1115.full.](http://science.sciencemag.org/content/326/5956/1115.full.pdf)
737 [pdf](http://science.sciencemag.org/content/326/5956/1115.full.pdf).

738 Haldane, J. B. S., 1937: The effect of variation on fitness. *The American Naturalist*, **71 (735)**,
739 337–349, URL <http://www.jstor.org/stable/2457289>.

740 Hellsten, U., and Coauthors, 2013: Fine-scale variation in meiotic recombination in
741 *Mimulus* inferred from population shotgun sequencing. *PNAS*, **110 (48)**, 19478–19482,
742 doi:10.1073/pnas.1319032110.

743 Henderson, J. S., R. A. Joyce, G. R. Hall, W. J. Hurst, and P. E. McGovern, 2007: Chemical
744 and archaeological evidence for the earliest cacao beverages. *Proceedings of the National*
745 *Academy of Sciences*, **104 (48)**, 18937–18940, doi:10.1073/pnas.0708815104, URL [http:](http://www.pnas.org/content/104/48/18937.abstract)
746 [//www.pnas.org/content/104/48/18937.abstract](http://www.pnas.org/content/104/48/18937.abstract), [http://www.pnas.org/content/104/48/](http://www.pnas.org/content/104/48/18937.full.pdf)
747 [18937.full.pdf](http://www.pnas.org/content/104/48/18937.full.pdf).

748 Hinch, A. G., and Coauthors, 2011: The landscape of recombination in african americans.
749 *Nature*, **476 (7359)**.

750 Hudson, R. R., and N. L. Kaplan, 1988: The coalescent process in models with selection and
751 recombination. *Genetics*, **120 (3)**, 831–840, URL [http://www.genetics.org/content/120/](http://www.genetics.org/content/120/3/831)
752 [3/831](http://www.genetics.org/content/120/3/831), <http://www.genetics.org/content/120/3/831.full.pdf>.

- 753 Kim, S., and Coauthors, 2007: Recombination and linkage disequilibrium in arabidopsis
754 thaliana. *Nature Genetics*, **39** (9).
- 755 Kosambi, D. D., 1943: The estimation of map distances from recombination val-
756 ues. *Annals of Eugenics*, **12** (1), 172–175, doi:10.1111/j.1469-1809.1943.tb02321.x,
757 URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1469-1809.1943.tb02321.x>, <https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1469-1809.1943.tb02321.x>.
- 759 Kundu, A., A. Chakraborty, N. A. Mandal, D. Das, P. G. Karmakar, N. K. Singh, and
760 D. Sarkar, 2015: A restriction-site-associated dna (rad) linkage map, comparative genomics
761 and identification of qtl for histological fibre content coincident with those for retted
762 bast fibre yield and its major components in jute (*corchorus olitorius* l., malvaceae s. l.).
763 *Molecular Breeding*, **35** (1), 19, doi:10.1007/s11032-015-0249-x, URL [https://doi.org/10.](https://doi.org/10.1007/s11032-015-0249-x)
764 [1007/s11032-015-0249-x](https://doi.org/10.1007/s11032-015-0249-x).
- 765 Lê, S., J. Josse, and F. Husson, 2008: FactoMineR: A package for multivariate analysis.
766 *Journal of Statistical Software*, **25** (1), 1–18, doi:10.18637/jss.v025.i01.
- 767 Li, W. H., and M. Nei, 1974: Stable linkage disequilibrium without epistasis in subdivided
768 populations. *Theoretical population biology.*, **6** (2), 173–183.
- 769 Lynch, M., 2010: Evolution of the mutation rate. *Trends in genetics : TIG*, **26** (8), 345–352,
770 URL <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC2910838/>.
- 771 Mackiewicz, D., A. J. Lustig, P. M. C. de Oliveira, S. Moss de Oliveira, and S. Cebrat, 2013:
772 Distribution of recombination hotspots in the human genome – a comparison of computer
773 simulations with real data. *PLoS ONE*, **8** (6).
- 774 Mackiewicz, D., M. Zawierta, W. Waga, and S. Cebrat, 2010: Genome analyses and
775 modelling the relationships between coding density, recombination rate and chromo-
776 some length. *Journal of Theoretical Biology*, **267** (2), 186 – 192, doi:[https://doi.](https://doi.org/10.1016/j.jtbi.2010.05.011)

777 [org/10.1016/j.jtbi.2010.08.022](https://doi.org/10.1016/j.jtbi.2010.08.022), URL [http://www.sciencedirect.com/science/article/pii/](http://www.sciencedirect.com/science/article/pii/S0022519310004315)
778 [S0022519310004315](http://www.sciencedirect.com/science/article/pii/S0022519310004315).

779 Martinez-perez, E., M. Schvarzstein, C. Barroso, J. M. Lightfoot, A. F. Dernburg, and A. M.
780 Villeneuve, 2008: Crossovers trigger a remodeling of meiotic chromosome axis composition
781 that is linked to two-step loss of sister chromatid cohesion. *Genes & development*, **22** **20**,
782 2886–901.

783 McVean, G., 2010: What drives recombination hotspots to repeat dna in humans? *Philo-*
784 *sophical Transactions of the Royal Society B*, **365** (**1544**), 1213–1218.

785 McVean, G. A. T., S. R. Myers, S. Hunt, P. Deloukas, D. R. Bentley, and P. Donnelly, 2004:
786 The fine-scale structure of recombination rate variation in the human genome. *Science*,
787 **304** (**5670**), 581–584, doi:10.1126/science.1092500, URL [http://science.sciencemag.org/](http://science.sciencemag.org/content/304/5670/581)
788 [content/304/5670/581](http://science.sciencemag.org/content/304/5670/581), <http://science.sciencemag.org/content/304/5670/581.full.pdf>.

789 Mézard, C., 2006: Meiotic recombination hotspots in plants. *Biochemical Society Transac-*
790 *tions*, **34** (**4**), 531–534, doi:10.1042/BST0340531, URL [http://www.biochemsoctrans.org/](http://www.biochemsoctrans.org/content/34/4/531)
791 [content/34/4/531](http://www.biochemsoctrans.org/content/34/4/531), <http://www.biochemsoctrans.org/content/34/4/531.full.pdf>.

792 Motamayor, J. C., P. Lachenaud, da Silva e Mota Jay Wallace, R. Looor, D. N. Kuhn, S. J.
793 Brown, and R. J. Schnell, 2008: Geographic and Genetic Population Differentiation of
794 the Amazonian Chocolate Tree (*Theobroma cacao* L). *PLOS ONE*, **3** (**10**), e3311, doi:
795 <https://doi.org/10.1371/journal.pone.0003311>.

796 Motamayor, J. C., and Coauthors, 2013: The genome sequence of the most widely
797 cultivated cacao type and its use to identify candidate genes regulating pod color.
798 *Genome Biology*, **14** (**6**), r53, doi:10.1186/gb-2013-14-6-r53, URL [https://doi.org/10.](https://doi.org/10.1186/gb-2013-14-6-r53)
799 [1186/gb-2013-14-6-r53](https://doi.org/10.1186/gb-2013-14-6-r53).

- 800 Myers, S., L. Bottolo, C. Freeman, G. McVean, and P. Donnelly, 2005: A Fine-Scale Map
801 of Recombination Rates and Hotspots Across the Human Genome. *Science*, **310** (5746),
802 321–324, doi:10.1126/science.1117196, URL [http://science.sciencemag.org/content/310/](http://science.sciencemag.org/content/310/5746/321)
803 [5746/321](http://science.sciencemag.org/content/310/5746/321.full.pdf), <http://science.sciencemag.org/content/310/5746/321.full.pdf>.
- 804 Ohta, T., 1982: Linkage disequilibrium due to random genetic drift in finite subdivi-
805 ded populations. *Proceedings of the National Academy of Sciences*, **79** (6), 1940–
806 1944, doi:10.1073/pnas.79.6.1940, URL <http://www.pnas.org/content/79/6/1940>, <http://www.pnas.org/content/79/6/1940.full.pdf>.
- 808 Paape, T., P. Zhou, A. Branca, R. Briskine, N. Young, and P. Tiffin, 2012: Fine-Scale
809 Population Recombination Rates, Hotspots, and Correlates of Recombination in the
810 *Medicago truncatula* Genome. *Genome Biology and Evolution*, **4** (5), 726–737, URL
811 <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3381680/>.
- 812 Pickrell, J. K., and J. K. Pritchard, 2012: Inference of population splits and mixtures from
813 genome-wide allele frequency data. *PLOS Genetics*, **8** (11), 1–17, doi:10.1371/journal.
814 [pgen.1002967](https://doi.org/10.1371/journal.pgen.1002967), URL <https://doi.org/10.1371/journal.pgen.1002967>.
- 815 Ptak, S. E., and Coauthors, 2005: Fine-scale recombination patterns differ between chim-
816 panzees and humans. *Nature Genetics*, **37** (4).
- 817 Quinlan, A. R., and I. M. Hall, 2010: Bedtools: a flexible suite of utilities for comparing
818 genomic features. *Bioinformatics*, **26** (6), 841–842, doi:10.1093/bioinformatics/btq033,
819 URL [+http://dx.doi.org/10.1093/bioinformatics/btq033](http://dx.doi.org/10.1093/bioinformatics/btq033), [/oup/backfile/content_public/](http://oup/backfile/content_public/journal/bioinformatics/26/6/10.1093_bioinformatics_btq033/3/btq033.pdf)
820 [journal/bioinformatics/26/6/10.1093_bioinformatics_btq033/3/btq033.pdf](http://oup/backfile/content_public/journal/bioinformatics/26/6/10.1093_bioinformatics_btq033/3/btq033.pdf).
- 821 R Core Team, 2018: *R: A Language and Environment for Statistical Computing*. Vienna,
822 Austria, R Foundation for Statistical Computing, URL <https://www.R-project.org/>.

823 Revelle, W., 2018: *psych: Procedures for Psychological, Psychometric, and Personality Re-*
824 *search*. Evanston, Illinois, Northwestern University, URL [https://CRAN.R-project.org/](https://CRAN.R-project.org/package=psych)
825 [package=psych](https://CRAN.R-project.org/package=psych), r package version 1.8.10.

826 Rizzon, C., G. Marais, M. Gouy, and C. Biémont, 2002: Recombination rate and the distri-
827 bution of transposable elements in the drosophila melanogaster genome. *Genome Research*,
828 **12 (3)**, 400–407, doi:10.1101/gr.210802, URL <http://genome.cshlp.org/content/12/3/400>.
829 abstract, <http://genome.cshlp.org/content/12/3/400.full.pdf+html>.

830 Rodgers, K., and M. McVey, 2015: Error-prone repair of dna double-strand breaks. *Journal*
831 *of Cellular Physiology*, **231**.

832 Ross-Ibarra, J., 2004: The Evolution of Recombination under Domestication: A Test of
833 Two Hypotheses. *The American Naturalist*, **163 (1)**, 105–112, doi:10.1086/380606, URL
834 <https://doi.org/10.1086/380606>, PMID: 14767840, <https://doi.org/10.1086/380606>.

835 Sanjuán, R., A. Moya, and S. F. Elena, 2004: The distribution of fitness effects caused
836 by single-nucleotide substitutions in an rna virus. *Proceedings of the National Academy*
837 *of Sciences*, **101 (22)**, 8396–8401, doi:10.1073/pnas.0400146101, URL [http://www.pnas.](http://www.pnas.org/content/101/22/8396)
838 [org/content/101/22/8396](http://www.pnas.org/content/101/22/8396), <http://www.pnas.org/content/101/22/8396.full.pdf>.

839 Schnable, P. S., and Coauthors, 2009: The b73 maize genome: complexity, diversity, and
840 dynamics. *Science (New York, N.Y.)*, **326 (5956)**.

841 Shanfelter, A., S. L. Archambeault, and M. A. White, 2018: Fine-scale recombination
842 landscapes between a freshwater and marine population of threespine stickleback fish.
843 *bioRxiv*, doi:10.1101/430249, URL [https://www.biorxiv.org/content/early/2018/09/29/](https://www.biorxiv.org/content/early/2018/09/29/430249)
844 [430249](https://www.biorxiv.org/content/early/2018/09/29/430249), <https://www.biorxiv.org/content/early/2018/09/29/430249.full.pdf>.

845 Singhal, S., and Coauthors, 2015: Stable recombination hotspots in birds. *Science*,

846 **350 (6263)**, 928–932, doi:10.1126/science.aad0843, URL [http://science.sciencemag.org/](http://science.sciencemag.org/content/350/6263/928)
847 [content/350/6263/928](http://science.sciencemag.org/content/350/6263/928.full.pdf), <http://science.sciencemag.org/content/350/6263/928.full.pdf>.

848 Siol, M., I. Bonnin, I. Oliveri, J. M. Prospero, and J. Ronfort, 2007: Effective population
849 size associated with self-fertilization: lessons from temporal changes in allele frequencies
850 in the selfing annual medicago truncatula. *Journal of Evolutionary Biology*, **20 (6)**, 2349–
851 2360, doi:10.1111/j.1420-9101.2007.01409.x, URL [https://onlinelibrary.wiley.com/doi/](https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1420-9101.2007.01409.x)
852 [abs/10.1111/j.1420-9101.2007.01409.x](https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1420-9101.2007.01409.x), [https://onlinelibrary.wiley.com/doi/pdf/10.1111/](https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1420-9101.2007.01409.x)
853 [j.1420-9101.2007.01409.x](https://onlinelibrary.wiley.com/doi/pdf/10.1111/j.1420-9101.2007.01409.x).

854 Smith, A., R. Hubley, and P. Green, 2016: Repeatmasker open-4.0.(2013-2015).

855 Stapley, J., P. G. D. Feulner, S. E. Johnston, A. W. Santure, and C. M. Smadja, 2017: Vari-
856 ation in recombination frequency and distribution across eukaryotes: patterns and pro-
857 cesses. *Philosophical Transactions of the Royal Society B: Biological Sciences*, **372 (1736)**,
858 doi:10.1098/rstb.2016.0455, URL [http://rstb.royalsocietypublishing.org/content/372/](http://rstb.royalsocietypublishing.org/content/372/1736/20160455)
859 [1736/20160455](http://rstb.royalsocietypublishing.org/content/372/1736/20160455.full.pdf), [http://rstb.royalsocietypublishing.org/content/372/1736/20160455.full.](http://rstb.royalsocietypublishing.org/content/372/1736/20160455.full.pdf)
860 [pdf](http://rstb.royalsocietypublishing.org/content/372/1736/20160455.full.pdf).

861 Stevison, L. S., and Coauthors, 2016: The time scale of recombination rate evolution in great
862 apes. *Molecular Biology and Evolution*, **33 (4)**, 928–945, doi:10.1093/molbev/msv331.

863 Weir, B. S., and C. C. Cockerham, 1984: Estimating f-statistics for the analysis of population
864 structure. *Evolution*, **38 (6)**, 1358–1370, URL <http://www.jstor.org/stable/2408641>.

865 Winckler, W., and Coauthors, 2005: Comparison of fine-scale recombination rates in humans
866 and chimpanzees. *Science (New York, N.Y.)*, **308 (5718)**, 107–111, URL [http://search.](http://search.proquest.com/docview/67570665/)
867 [proquest.com/docview/67570665/](http://search.proquest.com/docview/67570665/).

868 Wloch, D. M., K. Szafraniec, R. H. Borts, and R. Korona, 2001: Direct estimate of the
869 mutation rate and the distribution of fitness effects in the yeast *saccharomyces cerevisiae*.

870 *Genetics*, **159** (2), 441–452, URL <http://www.genetics.org/content/159/2/441>, [http://](http://www.genetics.org/content/159/2/441.full.pdf)
871 www.genetics.org/content/159/2/441.full.pdf.

872 Wu, J., and Coauthors, 2003: Physical maps and recombination frequency of six rice chro-
873 mosomes. *Plant Journal*, **36** (5), 720–730.

874 Yelina, N., P. Diaz, C. Lambing, and I. R. Henderson, 2015: Epigenetic control of
875 meiotic recombination in plants. *Science China Life Sciences*, **58** (3), 223–231, doi:
876 [10.1007/s11427-015-4811-x](https://doi.org/10.1007/s11427-015-4811-x), URL <https://doi.org/10.1007/s11427-015-4811-x>.