



HAL
open science

Testing for population decline using maximal linkage disequilibrium blocks

Elise Kerdoncuff, Amaury Lambert, Guillaume Achaz

► **To cite this version:**

Elise Kerdoncuff, Amaury Lambert, Guillaume Achaz. Testing for population decline using maximal linkage disequilibrium blocks. *Theoretical Population Biology*, 2020, 134, pp.171 - 181. 10.1016/j.tpb.2020.03.004 . hal-03492554

HAL Id: hal-03492554

<https://hal.science/hal-03492554>

Submitted on 22 Aug 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

Testing for population decline using maximal linkage disequilibrium blocks

Elise Kerdoncuff^{1,2,*}, Amaury Lambert^{2,3} and Guillaume Achaz^{1,2}

January 23, 2021

1: Atelier de Bioinformatique, UMR 7205 ISYEB, Sorbonne Université, CNRS, EPHE, Muséum National d'Histoire Naturelle, Paris, France

2: SMILE (Stochastic Models for the Inference of Life Evolution), UMR 7241 CIRB, Collège de France, CNRS, INSERM, PSL Research University, Paris, France

3: Laboratoire de Probabilités, Statistique et Modélisation (LPSM), UMR 8001, CNRS, Sorbonne Université, Paris, France

* corresponding author: elise.kerdoncuff@college-de-france.fr

Keywords: coalescent theory, recombination, demography, conservation biology.

Abstract

Only 6% of known species have a conservation status. Methods that assess conservation statuses are often based on individual counts and are thus too laborious to be generalized to all species. Population genomics methods that infer past variations in population size are easy to use but limited to the relatively distant past. Here we propose a population genomics approach that tests for recent population decline and may be used to assess species conservation statuses. More specifically, we study Maximal Recombination Free (MRF) blocks, that are segments of a sequence alignment inherited from a common ancestor without recombination. MRF blocks are relatively longer in small than in large populations. We use the distribution of MRF block lengths rescaled by their mean to test for recent population decline. However, because MRF blocks are difficult to detect, we also consider Maximal Linkage Disequilibrium (MLD) blocks, which are runs of single nucleotide polymorphisms compatible with a single tree. We develop a new method capable of inferring a very recent decline (e.g. with a detection power of 50% for populations which size was halved to N , $0.05 \times N$ generations ago) from rescaled MLD block lengths. Our framework could serve as a basis for quantitative tools to assess conservation status in a wide range of species.

1 Introduction

The severe and rapid changes imposed by human activities upon living organisms are suspected to be a major factor leading to short-term mass extinctions (Barnosky et al., 2011). The most comprehensive list of endangered species is the Red List of the International Union for Conservation of Nature (IUCN) (Rodrigues et al., 2006). Criteria used in the list to assess the species conservation status are based on geographical range, population trends, threats to habitat and ecology. Despite being very robust and reliable, these criteria are hard to establish for many species. To quantify the ongoing crisis for a wider range of organisms, there is a crucial need to develop quantitative measures of extinction risk to efficiently monitor species in real time and at a global scale. Previous attempts were developed to estimate quantitatively extinction rates, including by two of the present authors, based on occurrence data (Régnier et al., 2015; Ceballos et al., 2017; Sánchez-Bayo and Wyckhuys, 2019) or genetic data (from museum specimen Díez-del Molino et al. (2018); van der Valk et al. (2019)). The genetic methods measure the genetic diversity at different time to estimate the population size at these times and conclude on a general trend. The limitation of these methods is the difficulty to obtain time series data.

A handful of genomes sampled in a population at a single time point can help infer the past demography of this population (Gutenkunst et al., 2009; Li and Durbin, 2011; Excoffier et al., 2013; Harris and Nielsen, 2013; Sheehan et al., 2013; Schiffels and Durbin, 2014; Lapierre et al., 2017; Ringbauer et al., 2017; Terhorst et al., 2017; Beichman et al., 2018). In standard population genetic inferences, the periods when variations of population size can be estimated are of the order of N_e generations back in time. N_e denotes the so-called effective population size (Wright, 1931). Recent methods such as MSMC (Schiffels and Durbin, 2014) can provide inferences on more recent past but hardly scale up to large datasets of complete genomes because of their computational load. With the development of next generation sequencing, complete genomes from multiple individuals of the same species are now released routinely (Gibbs et al., 2015; Alonso-Blanco et al., 2016). Actual methods can not be applied to test for recent decline of populations, the models and methods we present in this manuscript specifically target very recent past when considering small populations and are meant to be applied to datasets of arbitrary size.

Methods using whole genome sequences to infer demography use different measures of genomic polymorphism. One of these measures is the so-called Site Frequency Spectrum,

62 or SFS (Fu, 1995). The SFS, that is the genome wide distribution of the frequencies of
63 polymorphic alleles in a sample of the population, is strongly distorted by the demographic
64 history of the species (Adams and Hudson, 2004; Marth et al., 2004). SFS-based methods
65 (e.g. Gutenkunst et al. (2009)) can handle arbitrarily large numbers of loci and genomes
66 but disregard correlations between sites caused by genetic linkage. Using genetic linkage
67 information may help overcoming the SFS-based methods limitations (e.g. difficulty to
68 discriminate between different scenarios Lapierre et al. (2017) and to infer recent demog-
69 raphy).

70 Recombination is the process by which two DNA sequences are intermixed to create
71 a new sequence that combines segments of different ancestries. When two homologous
72 regions of the genome are inherited from the same ancestor without having undergone
73 recombination, they are said IBD: *Identical By Descent*. The probability distribution and
74 the length of IBD regions passed through generations have been studied (Stam, 1980;
75 Chapman and Thompson, 2003; Stefanov, 2000).

76 Recombination patterns are also characterized by Linkage Disequilibrium (LD). LD
77 arises when individuals of a finite population share chunks of DNA inherited from a com-
78 mon ancestor (IBD blocks). Specifically, two variants located at two distinct sites are
79 in linkage disequilibrium (LD) when their joint frequency differs from what is expected
80 under independence. More specifically, LD is defined as the covariance $f_{A_1B_1} - f_{A_1}f_{B_1}$,
81 where f_{A_1} is the frequency of allele 1 at locus A (Lewontin and ichi Kojima, 1960). When
82 $f_{A_1} \times f_{B_1} = f_{A_1B_1}$, the two variants are said in complete linkage equilibrium. On average,
83 LD decreases exponentially with genetic distance due to recombination. The pattern of
84 LD is distorted by demography (Hill and Robertson, 1968) and thus can be used to infer
85 the past demography of a population (Hollenbeck et al., 2016; Patin et al., 2014).

86 Importantly, despite the fact that breakpoints between IBD blocks are usually not
87 observable when comparing two homologous regions, “long enough” IBD blocks can be
88 retrieved by applying one of several recent methods to a pair of sequences (Purcell et al.,
89 2007; Gusev et al., 2009; Browning and Browning, 2010). These methods are based on
90 detecting long identical shared segments (Gusev et al., 2009) or shared regions that harbor
91 multiple rare variants (Purcell et al., 2007; Browning and Browning, 2010). If two individ-
92 uals share the same rare variant, they may also share the surrounding chromosomal region,
93 particularly because rarer variants are more likely to be relatively recent. Most methods
94 take sequencing errors into account, allowing IBD blocks to not be totally identical. The

95 accuracy of IBD block detection depends on the algorithm used (Browning and Browning,
96 2013).

97 Some demographic inference methods are based on the distribution of lengths of inferred
98 pairwise IBD blocks in a population. Palamara et al. (2012) have calculated the distribution
99 of expected lengths of pairwise IBD blocks for a given parameterized demographic model.
100 Browning and Browning (2015) have calculated the expected time to the most recent
101 common ancestor (TMRCA) of an IBD block as a function of its length. Then they use
102 the empirical density of IBD block lengths to estimate the distribution of TMRCA and
103 thus the variations of effective population size through time.

104 Other methods use the length of identical shared segments of chromosome within a
105 diploid individual (Hayes et al., 2003). Two identical shared segments may be inherited
106 from a common ancestor without recombination event (and then be IBD) or may not be
107 IBD as there are invisible recombination events that may have occurred within it. The
108 probability that the two haplotypes of an individual share identical alleles for a given
109 number of adjacent positions can be predicted (Hayes et al., 2003; MacLeod et al., 2009).
110 Tools have been developed to apply these methods to infer demographic inference from
111 genomic data (MacLeod et al., 2013; Harris and Nielsen, 2013).

112 Yet another approach to infer demographic history from IBD blocks is to reconstruct
113 the genome-wide distribution of the TMRCA between two haploid genomes. In PSMC,
114 Li and Durbin (2011) devised a Hidden Markov Model that infers the TMRCA from the
115 positions of heterozygous sites along a pair of sequences and then estimate a step-wise
116 demographic pattern. MSMC, the extension of PSMC (Schiffels and Durbin, 2014), ana-
117 lyzes the heterozygosity pattern from multiple individuals and uses first coalescence events
118 between any two haploid genomes of the sample. These methods are computationally
119 intensive (as of today, MSMC cannot infer the demographic history of more than 8 indi-
120 viduals) and pool the diversity on windows of 100 bp, that are assumed to form a single
121 locus with two states, heterozygous or homozygous.

122 Importantly, the previous methods infer stepwise changes of the “effective population
123 size” ($N_e(t)$) that are estimated from the density of coalescence events. This motivated
124 Mazet et al. (2015, 2016); Chikhi et al. (2018); Rodríguez et al. (2018) to propose to replace
125 $N_e(t)$ by the more explicit *Inverse Instantaneous Coalescence Rate*. IICR only matches the
126 instantaneous population size when the population is panmictic. It is nonetheless always
127 possible to find a population model with constant size but spatial structure that corre-

128 sponds to any IICR of a size-changing population for the TMRCA of 2 sequences (Chikhi
129 et al., 2018). For larger samples, the joint distribution of coalescence events $[T_2, T_3, \dots]$
130 can be used, in theory, to disentangle structure from demography (Grusea et al., 2019).

131 Existing methods for demographic inference using recombination information **often** use
132 the whole genome of few individuals (less than 10) or use a smaller part of the genome.
133 These methods only consider the joint history of two individuals (e.g. the pairwise IBD
134 length distribution or the time of the first coalescence event between any two haploid
135 genomes) which algorithmic complexity increases drastically with the number of individuals
136 (e.g. detection of pairwise IBD blocks is quadratic) and generates a computational load
137 limiting **in most cases** the application of the methods to a larger number of individuals. On
138 the other hand, with few individuals, demographic inferences are unable to detect recent
139 changes of population size.

140 Following the idea of Turet and Hospital (2017), we decided to study the IBD concept
141 extended to a multilocus segment and a larger number of individuals ($n > 2$). Some studies
142 have been conducted on the amount of genetic material shared IBD with $n > 2$, consider-
143 ing closely related individuals (Donnelly, 1983; Ball and Stefanov, 2005). We extend the
144 concept at a population level while relaxing the need for *identical* sequence (without mu-
145 tation), which is why we decided to define a new term. We call ‘MRF blocks’ homologous
146 segments that are entirely inherited from the same ancestor without recombination; these
147 segments may or may not harbour different alleles, because of mutations. An MRF block
148 is a segment of an alignment of haploid genomes that share the same coalescent tree. A
149 recombination event along the sequence cuts the genome alignment into two MRF blocks,
150 one on each side of the recombination point. By definition there is no recombination
151 within MRF blocks so that all variants located within an MRF block are necessarily in
152 complete linkage disequilibrium. The reciprocal is not true, as variants in complete LD do
153 not necessarily belong to the same MRF block. MRF blocks carry the information of any
154 recombination event that happened among the sampled individuals. As for IBD blocks,
155 MRF blocks are usually not observable.

156 **Outline** We have developed a new test to detect very recent population declines of en-
157 dangered species. We first consider the full length distribution of MRF blocks in a sample

158 of haploid genomes ($n \geq 2$). Second, as MRF blocks are not directly observable from
159 sequence alignments, we devised a simple and efficient algorithm to chop an alignment of
160 $n \geq 4$ haploid genomes in Maximal Linkage Disequilibrium (MLD) blocks, that are seg-
161 ments which variants are in complete LD. From the length distributions of MRF blocks
162 or MLD blocks, we devised a summary statistic to test whether a population has been
163 declining in the very recent past. Our method is not limited by the number of genomes in
164 the sample.

165 **2 Model and Methods**

166 In the absence of recombination, ancestral relationships between genomes can be repre-
167 sented in the form of a genealogical tree. Individual haploid genomes at present time are
168 the leaves of the tree, the MRCA of these individuals is the root. The fusion of two lin-
169 eages into one (a common ancestor) is named coalescence event (Kingman, 1982), hence
170 the name of “coalescent tree”. The sum of all branch lengths that separates two genomes
171 up to their common ancestor is the time of divergence between them, usually expressed in
172 generations. In the Wright-Fisher model with constant population size N , branch lengths
173 measured in number of generations scale like N . In particular, if we define T as the total
174 length of the coalescent tree, the expectation of T is proportional to N . Large popula-
175 tions generate coalescent trees with deep nodes, whereas small populations have shallow
176 coalescent trees.

177 In the presence of recombination, two loci of an alignment have the same coalescent
178 tree only if no recombination event happened since their MRCA. We name MRF block,
179 a maximal interval along the alignment of sites sharing the same coalescent tree. MRF
180 blocks are consequently separated by recombination points, corresponding to recombination
181 events. It is standard to assume that conditional on the total length T of the coalescent
182 tree of a site, the length L of its MRF block is exponentially distributed with rate ρT ,
183 where ρ is the recombination rate (expressed in a arbitrary unit proportional to Morgan).
184 Then for a fixed ρ , T and L are negatively correlated: recombination is more likely to
185 occur in deep trees, which thus are carried by shorter blocks. As mentioned above, as T
186 is proportional to N , MRF blocks are also shorter in larger populations. More accurately,
187 because the law of T/N does not depend on N , neither does the law of NL (for $n = 2$
188 it alludes to results of Carmi et al. (2014)). In other words, if population 1 has size N_1

189 and population 2 has size N_2 , the distribution of MRF block lengths in population 2 can
190 be deduced from that in population 1 by a scaling factor N_1/N_2 , both populations having
191 identical demography otherwise. For example if $N_2 = 2N_1$, the MRF blocks in population
192 2 are twice smaller than those of population 1.

193 Note that for a given N the lengths (L_1, L_2, \dots) of successive adjacent blocks have the
194 same distribution, but they are not independent, because the coalescent trees of adjacent
195 MRF blocks are not. The dependencies between these trees is encoded in the so-called
196 Ancestral Recombination Graph (ARG) (Griffiths and Marjoram, 1997). Because these
197 dependencies have a complex structure (Wiuf and Hein, 1999), a popular way of approxi-
198 mating them is the Sequentially Markovian Coalescent (SMC) (McVean and Cardin, 2005;
199 Marjoram and Wall, 2006). This approximation neglects coalescences between lineages
200 with no overlapping ancestral material and assumes Markovian dependencies of coalescent
201 trees along the sequence: the genealogy of an MRF block only depends on the genealogy
202 of the adjacent ones.

203 Although genealogies of different MRF blocks are not independent, they are asymptot-
204 ically independent as the distance between them increases.

205 Throughout this article, we use msprime to generate MRF blocks directly from the
206 ARG (Kelleher et al., 2016) but very similar results were obtained with a local SMC
207 implementation. We assume constant recombination and mutation rate along the genome.
208 We simulated the alignment of $n = 10$ haploid genomes at present time.

209 **Demographic scenario.** We consider a single change of population size (Fig 1a). Here
210 N_t represents the population size at time t , $t = 0$ is the present time and positive values
211 represent the past. We denote by κ the ratio of the two sizes: $\kappa = N_\infty/N_0$, and by τ the
212 time at which the population size changes in coalescent units of N_0 generations. If $\kappa = 1$,
213 $N_\infty = N_0$: there is no change. If $\kappa = 10$, $N_\infty = 10N_0$: the population size has been divided
214 by 10, τN_0 generations in the past.

215 3 MRF blocks

216 3.1 Distribution of block lengths

217 **Impact of population decline on tree length.** For declining populations ($\kappa > 1$), the
218 coalescent trees have two distinct time scales: a first one for the shallow part of the tree
219 ($t < \tau$), expressed in N_0 generations, and a second one for the deep part of the tree ($t > \tau$)
220 that is expressed in κN_0 generations. When the declining population tree is compared to a
221 standard coalescent tree (constant population size), it has shorter external branches if the
222 reference time scale is expressed in κN_0 generations *or* longer internal ones if the reference
223 time scale is expressed in N_0 generations. When it is compared to a reference tree with
224 population size chosen so as to have the same $T_{MRC A}$, its external branches are too short
225 *and* its internal branches are too long. Similarly, the distribution of the total length T of
226 the tree is overdispersed when compared to the length of the standard coalescent tree with
227 the same mean.

228 **Impact of population decline on lengths of MRF blocks.** For a declining population,
229 the distribution of the length L of MRF blocks will depend not only on ρ and N_0 but also
230 on κ and τ . As the tree relative branch lengths are distorted and the distribution of T
231 is overdispersed, so is the distribution of L . In a declining population, the distribution
232 of L can be seen as a mixture of the two distributions that correspond to the two pop-
233 ulation sizes, N_0 and κN_0 . The strength of the decline (κ) tunes the difference between
234 the distributions; the date of decline (τ) tunes in what proportion the two distributions
235 are mixed. When $\tau \rightarrow 0$ (practically, $\tau < 10^{-4}$ times N_0 generations for a sample size
236 $n \in [10, 100]$), the distribution of L is indistinguishable from that of block lengths in a
237 population with constant size equal to κN_0 . At the opposite, for $\tau \rightarrow \infty$ (practically,
238 $\tau > 10$ times N_0 generations for a sample size $n \in [10, 100]$), the distribution of L is indis-
239 tinguishable from that of block lengths in a population with constant size equal to N_0 . As
240 a result for $\tau \in [10^{-4}, 10]$, the distribution of L has an excess of MRF blocks smaller than
241 the N_0 reference and an excess of MRF blocks longer than the N_∞ reference (Fig 1b). The
242 small blocks correspond to the trees which total length T is mostly driven by the distant
243 N_∞ time scale and the long ones to the trees which total length T is mostly driven by the
244 recent N_0 time scale.

245 As mentioned in the previous section, in a population with constant size N , the dis-
246 tribution of L , briefly denoted L_N , scales like $1/N$, in the sense that the distribution of
247 $\tilde{L} := NL_N$ does not depend on N . In particular, the distribution of $L' := L/E[L]$ does not

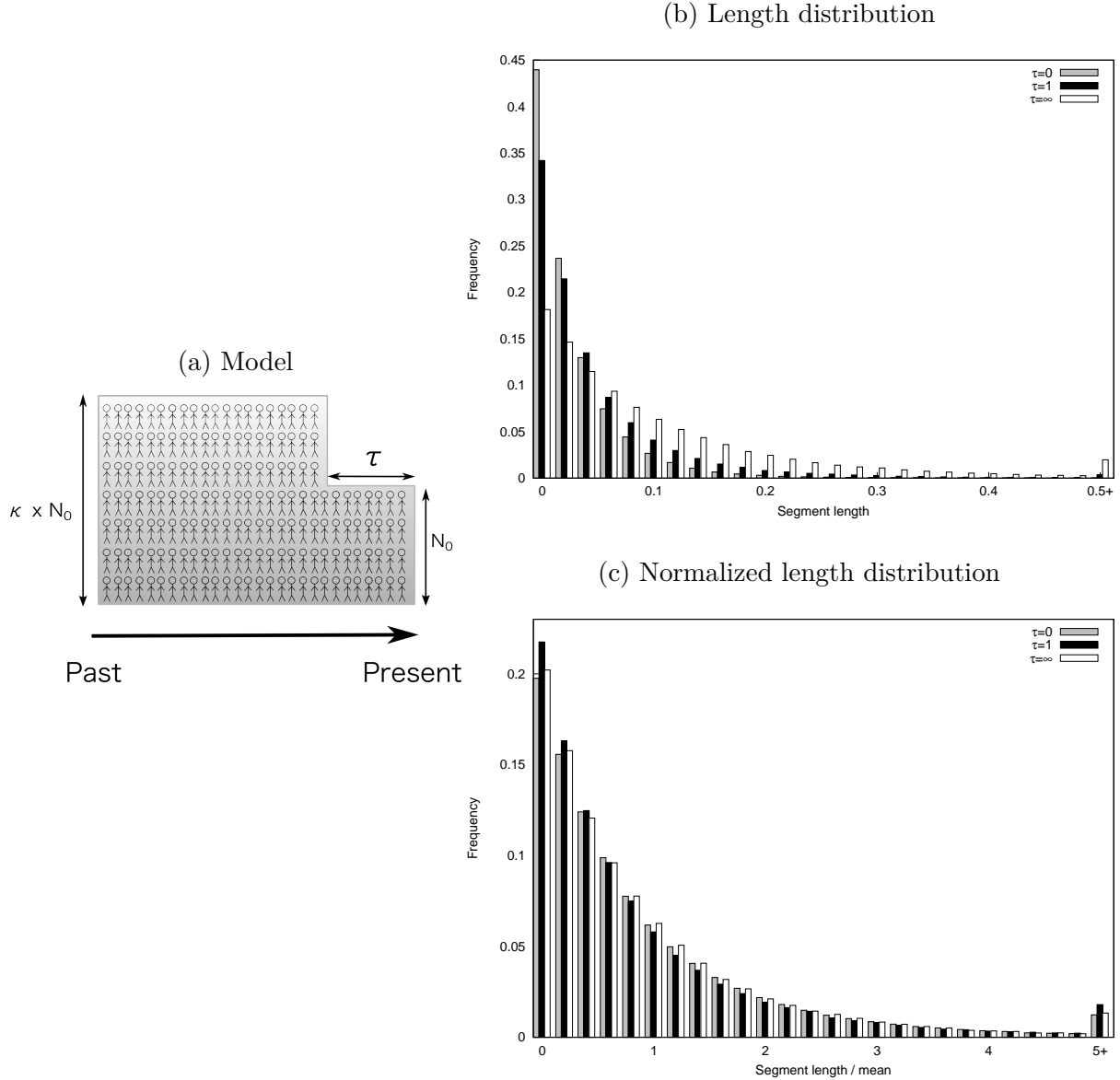


Figure 1: Impact of the demography on the distribution of MRF block lengths. (a) The demography considered here is a sudden size change tuned by 3 parameters : N_0 , the actual population size, τ the date of decline (in backward time) and κ the strength of decline. Time is expressed in N_0 generations. (b) Distribution of L for $\rho = 1$, $\kappa = 3$ with $\tau = \{0, 1, \infty\}$. When $\tau = 0$ (grey) or $\tau = \infty$ (white), population size is constant. (c) Distribution of $L' = L/E[L]$ under the same values of $\rho = 1$, κ and τ . In case of a decline, the distribution is overdispersed, with an excess of both short and long normalized MRF blocks.

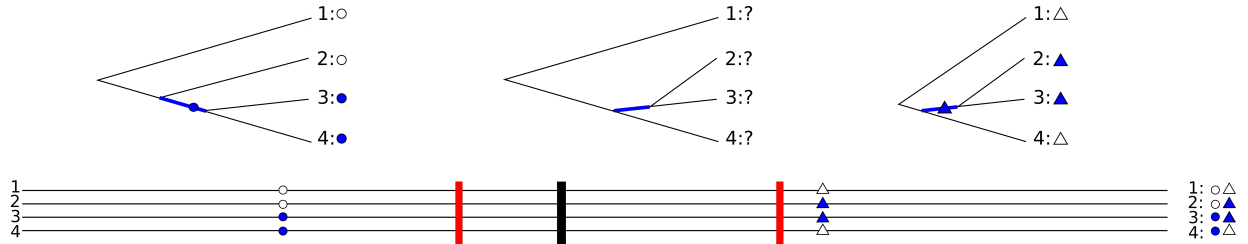


Figure 2: **Detection of recombination in three MRF blocks.** The four lines represent four haploid genomes, circles and triangles are mutation events, red lines are the true recombination events delimiting MRF blocks and above each MRF block is represented its true tree. Mutation events occur on certain lineages as represented on the trees. The first recombination event generates an incompatibility between the blue branches of the first two MRF blocks, but as no mutation occurs on the second MRF block, this recombination event cannot be detected. The second recombination event does not change the topology of the tree and thus this second event cannot be detected either. However, the first and the third MRF blocks carry mutations that are not compatible; thus a minimum of one recombination event can be inferred between the two mutations, as indicated by a vertical thick black line arbitrarily placed in the middle.

248 depend on N in a population with constant size and follows the law of $\tilde{L}/E[\tilde{L}]$. However,
 249 the distribution of L' is distorted when there is a size change. For a declining population,
 250 the distribution of L' is overdispersed, it has an excess of small blocks (*i.e.* less than 0.2)
 251 and an excess of long blocks (*i.e.* more than 5), as can be seen on Fig 1c.

252 Note that we always have $E[L'] = 1$, but here $E[L]$ has

$$\frac{1}{N_\infty}E[\tilde{L}] = E[L_{N_\infty}] < E[L] < E[L_{N_0}] = \frac{1}{N_0}E[\tilde{L}].$$

253 As the block distribution of L_N is a mixture of the one of L_{N_0} and L_{N_∞} , $E[L]$ is bounded
 254 by $E[L_{N_\infty}]$ and $E[L_{N_0}]$ that depend on the population size.

255 4 MLD blocks

256 4.1 Definition

257 All recombination events are not directly visible in a genome alignment. First, adjacent
 258 MRF blocks may have coalescent trees sharing the same topology and the same branch
 259 lengths, so that mutations occurring on either tree show exactly the same pattern on
 260 either block. Second, adjacent MRF blocks may have coalescent trees sharing the same

261 topology but not the same branch lengths, so that mutations occurring on either tree
262 display the same bipartitions (compare the second and third tree in Figure 2). Third, even
263 if two adjacent MRF blocks have trees with different topology, it is possible that branches
264 distinguishing these topologies do not carry mutations (see the second block in Figure 2).

265 Importantly, recombination events that happen between the two oldest lineages do not
266 impact the topology of the tree, so are never detectable because they do not impact the
267 possible bipartitions.

268 A possibility used in the literature to detect breakpoints between MRF blocks is to
269 detect the changes in the density of polymorphic sites along the sequence due to the
270 change of coalescent tree (like in PSMC, Li and Durbin (2011)).

271 Here we used instead the incompatibilities between bipartitions displayed by polymor-
272 phic sites to place the minimal number of recombination events on the alignment. Two
273 bipartitions are said incompatible when they are not compatible with a common tree.

274 In what follows, we will assume that the mutation rate μ is constant through time and
275 along the genome.

276 4.1.1 The four-gamete test

277 From now on, we assume that each site can be hit at most once by a mutation, so that
278 a polymorphic site is always bi-allelic, an assumption known as the “infinitely-many sites
279 model”. The four-gamete test (Hudson and Kaplan, 1985) serves to detect incompatibilities
280 between bipartitions displayed by two polymorphic sites. For any two biallelic sites (A/a
281 and B/b) there are at most four gamete haplotypes in the population (A-B A-b a-B and a-
282 b). Under the infinitely-many sites model, the four possible haplotypes cannot be observed
283 in a sample if the two sites share the same genealogy. Then if the four possible haplotypes
284 are observed in the sample, a recombination event must have occurred between them –
285 but not necessarily the other way round. This property can be used to compute a lower
286 bound for the number of recombination events in a genome alignment (Hudson and Kaplan,
287 1985) or even to estimate the recombination rate (Hey and Wakeley, 1997). We used it to
288 compute and place the minimal number of breakpoints in a genome alignment.

289 Two polymorphic sites are said *incompatible* if the four possible haplotypes are present
290 in the sample. When a sequence of adjacent polymorphic sites contains no pairwise incom-
291 patibility, we speak of a sequence of compatible sites. Note that a sequence of compatible
292 sites are in complete linkage disequilibrium. We thus define an MLD block, for *Maximal*

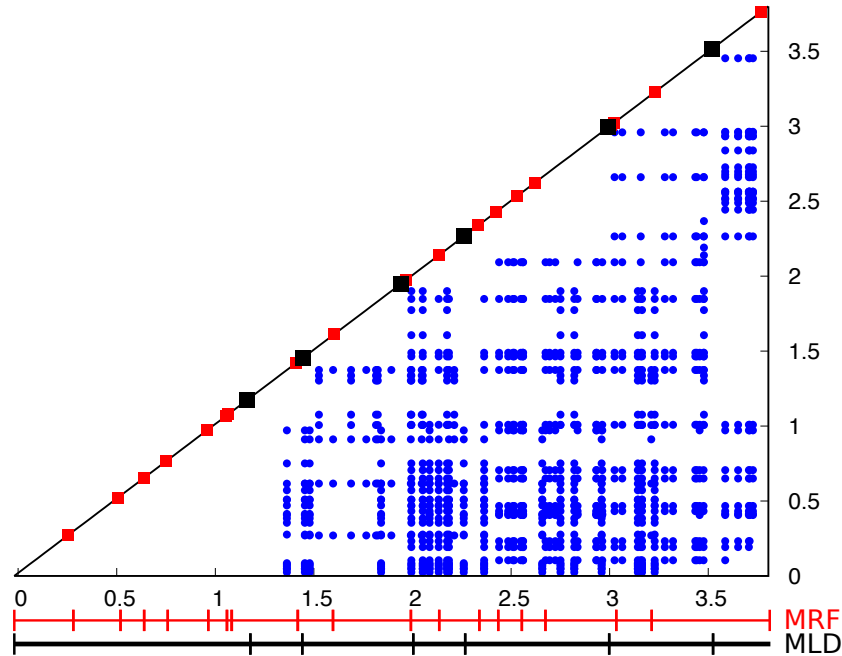


Figure 3: The incompatibility matrix and the chopping algorithm. X and Y-axis are positions on the genome alignment. Blue dots represent a pair (x,y) of incompatible sites. The red squares are the true positions of recombination events (MRF breakpoints) and the black squares are MLD breakpoints inferred by the chopping algorithm.

293 *Linkage Disequilibrium* block, as any maximal sequence of compatible sites.

294 We now explain how to extend this notion originally designed for haploid genomes (or
 295 phased diploid genomes) to an unphased diploid genome, that is, a diploid sequence lacking
 296 the linkage information. For an unphased diploid genome, the two original haplotypes can
 297 be determined if the diploid genome is homozygous at at least one the two sites:

- 298 • When the genome is homozygous at both loci ($A/A-B/B$), both haplotypes must be
 299 $A-B$.
- 300 • When the genome is homozygous at one locus and heterozygous at the second one
 301 ($A/A-B/b$), the haplotypes must be $A-B$ and $A-b$.

302 The four-gamete test can then be extended to a sample of unphased diploid genomes
 303 by saying that two sites are incompatible in this sample if they are incompatible in the
 304 subsample of haplotypes that have been inferred thanks to the previous remark. When
 305 the haplotype is ambiguous, the sites are considered compatible and do not bring more
 306 information about a recombination event.

307 4.1.2 The chopping algorithm

308 We used the four-gamete test to detect incompatibilities in the genome alignment and to
309 chop it into MLD blocks (Fig 3). To avoid computing the full matrix of pairwise incompati-
310 bilities between all polymorphic sites of the genome, we only compute the incompatibilities
311 for sequences of P adjacent polymorphic sites (by default $P = 150$). Each pair of incom-
312 compatible sites (i, j) defines an interval that contains at least one MLD breakpoint. To place
313 the MLD breakpoint, we seek the shortest interval that is sufficient to explain the incom-
314 patibilities.

315 **Algorithm:** We retrieve all intervals and sort them in increasing order of site positions
316 along the genome (first by i the first site position and when equal, by j the second site
317 position). As we scan two times the list of intervals, the algorithm complexity is linear
318 with the number of polymorphic sites:

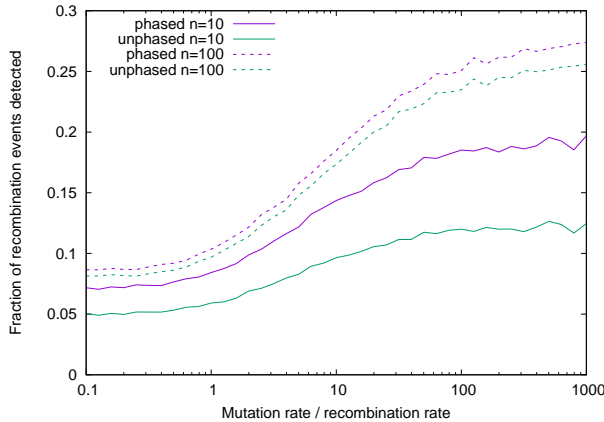
- 319 1. **Discarding and shortening.** For this step, we scan the list in reverse order, from
320 the last (N) to the first interval. (The algorithm can be done in the forward order,
321 the distribution will be slightly different but it will not affect the study.) Each
322 interval containing another entire interval is discarded: for two intervals (i_N, j_N)
323 and (i_{N-1}, j_{N-1}) , if $i_N \leq i_{N-1} \leq j_{N-1} \leq j_N$, then (i_N, j_N) is discarded. When two
324 intervals overlap, they are replaced by their intersection (the two original ones are
325 discarded): for the two intervals (i_N, j_N) and (i_{N-1}, j_{N-1}) , if $i_{N-1} \leq i_N \leq j_{N-1} \leq j_N$,
326 both are replaced by a new interval (i_N, j_{N-1}) , that is then compared to (i_{N-2}, j_{N-2}) ...
- 327 2. **Positioning.** From the final list of disjoint intervals, we place an MLD breakpoint
328 at the middle of each interval.

329 MLD breakpoints partition the genome alignment into MLD blocks.

330 4.2 Length distribution

331 The distribution of the length L_c of a typical MLD block does not only depend on the
332 distribution of L (MRF block length) but also on the fraction p of recombination events
333 that are detected. This fraction increases with the ratio μ/ρ , as illustrated in Figure 4a.
334 When many mutations occur in two different MRF blocks ($\mu \gg \rho$), the probability that
335 they occur on incompatible branches of their respective coalescent trees increases and so

(a) Factors impacting the power of detection



(b) MLD block lengths

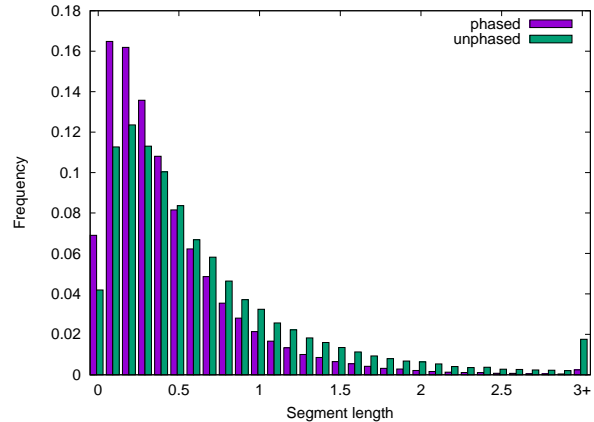


Figure 4: Detection of recombination events and its impact on MLD block length (L_c) distribution under a constant population. (a) Fraction of recombination events that are detected as a function of μ/ρ for different sample sizes ($n = 100$, dashed lines and $n = 10$ plain lines) and for phased (purple) or unphased (green) diploid genomes. (b) Distribution of MLD block lengths for phased (purple) and unphased (green) diploid genomes in a population of constant size ($\mu = 10$, $\rho = 1$, $n = 10$).

336 does the detection efficiency, up to a point of saturation due to cases when these MRF
 337 blocks share the same tree topology. The number of sampled individuals also impacts the
 338 efficiency of detection (Fig 4a): the larger the sample size, the higher the probability to
 339 observe incompatible mutations. The four-gamete test for unphased diploid genomes has
 340 obviously less power to detect recombination than for phased genomes (Fig 4a).

341 The lower the power to detect recombination, the longer the MLD blocks. In particular,
 342 phased genomes have smaller MLD blocks than unphased ones (Fig 4b). Furthermore
 343 increasing the sample size results in more detectable recombination points and thus smaller
 344 MLD blocks. In Figure 4b, the average block length, in our arbitrary unit for $n = 10$ phased
 345 haploid genomes is $\bar{L}_c = 0.497$ ($\mu = 10$, $\rho = 1$). Considering smaller sample size will result
 346 in larger MLD blocks (e.g. $\bar{L}_c = 1.32$ for $n = 5$). This implies that the total number of
 347 blocks can be limiting for small sample size, and that these long blocks will be harder to
 348 detect in scaffolds of partial genomes. In (very) large samples, MLD blocks are shorter:
 349 $\bar{L}_c = 0.132$ for $n = 600$ (ten times smaller than for $n = 5$) and $\bar{L}_c = 0.103$ for $n = 6,000$.
 350 On a side note, the theoretical pitfall of having too small “undetectable” blocks can always

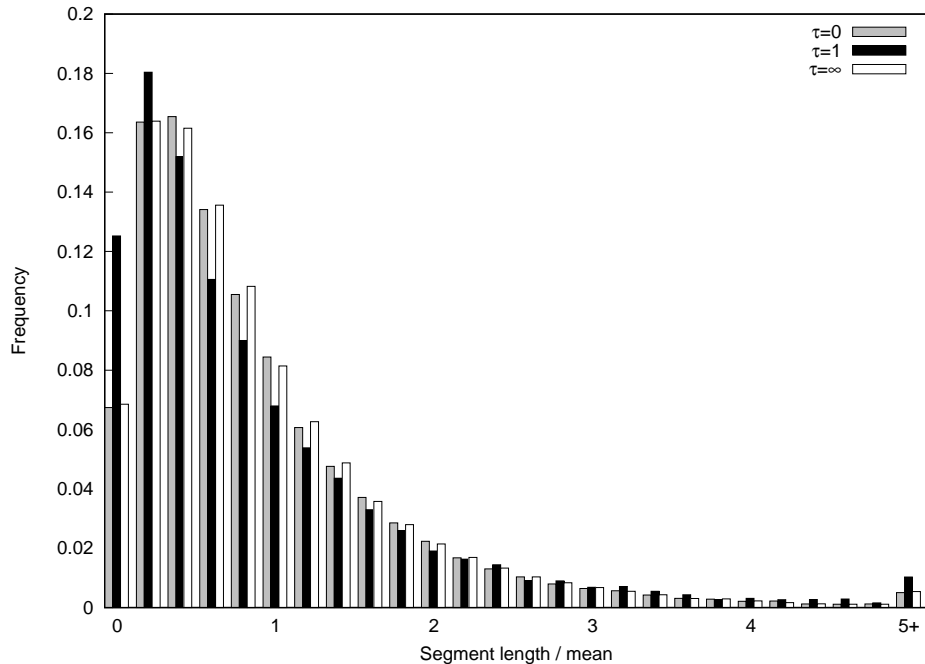


Figure 5: Distribution of L'_c for a population of constant size (white, $N = N_0 = 1$ and grey, $N = \kappa N_0 = 3$) and for a declining population (black for $\tau = 1$) with $\rho = 1$, $\mu = 10$ and $n = 10$.

351 be overcome by subsampling.

352 Here, we consider the block lengths normalized by the average length $L'_c = L_c / \bar{L}_c$.
 353 Similarly to the MRF blocks, the distribution of L'_c does not depend on the value of N but
 354 does depend on the demographic scenario (Fig 5). However, it still depends on our ability
 355 to detect recombination and so on the ratio μ/ρ and n the number of sampled individuals.
 356 To compare distributions, it is then important that they have the same ratio μ/ρ and the
 357 same n .

358 Similar to what we have observed for MRF blocks, a declining population exhibits both
 359 an excess of small blocks ($L'_c < 0.2$) and large blocks ($L'_c > 5$) (Fig 5). The shape of the
 360 distribution of L'_c (Fig 5) differs from the one for L' (Fig 1c): MLD blocks are longer than
 361 MRF blocks. Indeed, they contain a variable number of MRF blocks and below a certain
 362 size, MRF blocks are not detectable as recombination points at the edges of an MRF block
 363 can be detected only when mutations have occurred inside the block. MLD blocks are
 364 always longer than MRF blocks.

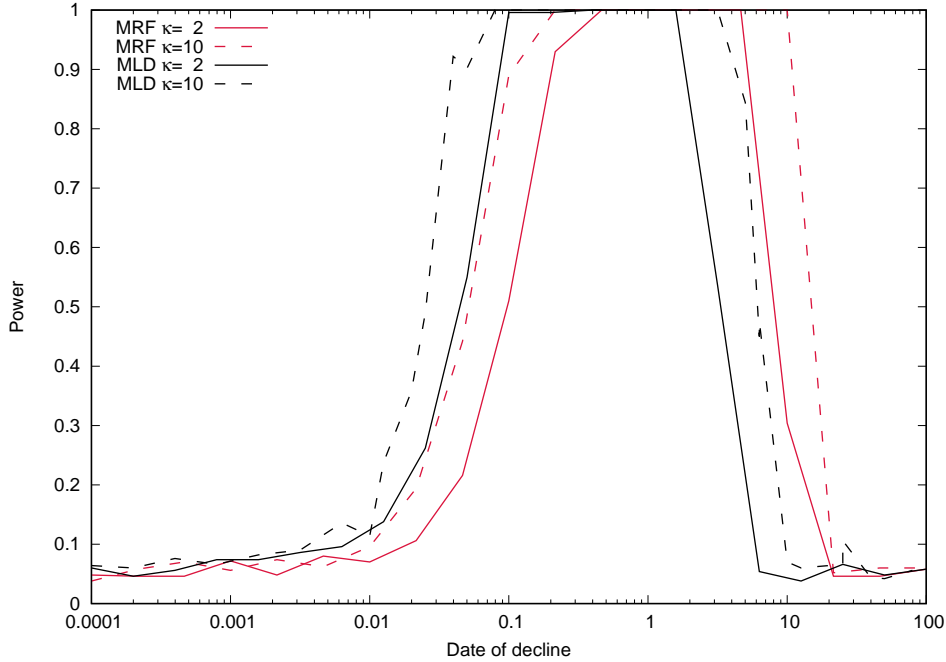


Figure 6: Power to detect population decline. The test based on MRF blocks (f) is pictured in red, whereas the one based on MLD blocks (f_c) is represented in black. We assess the power of the two tests for $\kappa = 2$ (plain line) and $\kappa = 10$ (dashed line) with $\tau \in [0.0001, 100]$.

365 5 Statistical tests for population decline

366 5.1 Test

367 To test for population decline, we use the excess of small and large blocks that we observe
 368 when comparing samples from a declining *vs* a constant population size. More specifically,
 369 we compute the fraction of blocks which normalized length is either smaller than 0.2 or
 370 larger than 5, both in the case of MRF blocks ($f = f_{L' < 0.2} + f_{L' > 5}$) and of MLD blocks
 371 ($f_c = f_{L'_c < 0.2} + f_{L'_c > 5}$). To set an empirical threshold value under H_0 , we simulate 10,000
 372 genomes of 10^5 MRF blocks under a constant population size for 10 haploid genomes
 373 and compute both $f^{5\%}$ and $f_c^{5\%}$ as upper limits for one-tailed tests: $f^{5\%} = 0.214236$
 374 and $f_c^{5\%} = 0.075824$. As the threshold is empirical, simulations need to be redone for
 375 a change in sampled size or in null model. Time needed for simulations depends on the
 376 algorithm/software used and the specific features of the model.

377 **5.2 Power**

378 To assess the power of this test, we simulated 1,000 replicates under population decline
379 (H_1 with various τ and κ) and report the fraction of runs where $f^{H_1} > f^{5\%}$ for MRF blocks
380 or $f_c^{H_1} > f_c^{5\%}$ for MLD blocks. When the power is 1, the decline was significant in all runs.
381 When the power is 5%, the decline is not detectable, the test can not differentiate H_0 and
382 H_1 .

383 Without surprise, results show that the power of the test to detect population decline
384 depends on both the decline strength (κ) and the date of decline (τ) (Fig 6). For both
385 tests based either on MRF or on MLD blocks, the power outreaches the 5% risk only for
386 a range of τ . The type I error of the test is 5% as expected. For both tests, the range
387 of detection is wider when the decline is stronger (compare dashed to solid lines in Fig
388 6). The surprise is that the test based on MLD blocks (f_c) detects more recent declines
389 than the test based on MRF blocks (f). Therefore, we recommend using the f_c test when
390 searching for very recent decline even if MRF blocks are known (which is generally not the
391 case).

392 **6 Application to data: the case of the western lowland** 393 **Gorillas**

394 **6.1 Handling the low quality of real genomes**

395 Genomic data sets often include sequencing errors and regions that are not genotyped.
396 Consequently, the f_c test cannot be run as is on these data sets. We present some mod-
397 ifications to our test to handle the poor quality of data. We show in this section that
398 adjustments can be made to get the information from the L'_c distribution.

399 **Simulations with lower quality**

400 Difficulties in applying the f_c test to real data sets can stem from the low quality of
401 DNA sequences. We replicated in the simulated genomes the two main issues, namely
402 the interruptions of DNA tracts and the absence of genotyping for some SNPs in some
403 individuals.

404 DNA tract interruptions truncate MLD blocks and make their detection difficult. The

405 number, size and location of these interruptions will have an effect on the detection of
406 MLD blocks and thus will alter the L'_c distribution. To handle the effect of interruptions,
407 we placed the interruptions at the same positions in our simulated chromosome as in the
408 real chromosome.

409 As for the partial genotyping issue, we artificially lowered the genotyping quality in the
410 simulated chromosomes. We used the empirical distribution of missing individuals (e.g.
411 chr1 of *Gorilla gorilla*, Fig 7) to pick random positions in the simulated chromosome and
412 erase the genotypes of some individuals.

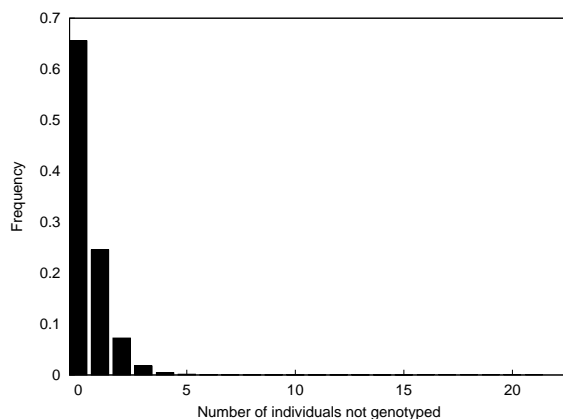


Figure 7: Distribution of the number of individuals not genotyped per SNP on chromosome 1 of the Gorillas dataset (Prado-Martinez et al., 2013).

413 **Mutation rate and recombination rate**

414 To cope properly with the issue of genotyping, we simulated chromosomes with the same
415 number of mutations and the same MLD length mean as in the studied data set. We use the
416 Watterson estimator (Watterson, 1975) for the mutation rate and fixed the recombination
417 rate so that simulated and real chromosomes had the same average length of MLD block.

418 **6.2 Application to Chr1 of *Gorilla gorilla gorilla***

419 We applied this methodology on chromosome 1 of twenty-three unrelated western lowland
420 Gorillas (*Gorilla gorilla gorilla*) from the Great Ape Genome Project (Prado-Martinez
421 et al., 2013). The chromosomes have 247,249,719 base pairs. The 23.1% of sites that are
422 considered “low coverage”(Prado-Martinez et al., 2013) divide the chromosome alignment

423 into 6,277,293 uninterrupted stretches. The 5,388,083 interruptions due to a single site
 424 were not considered as *interruptions*. To speed up simulations, we considered stretches
 425 longer than 499 sites, as smaller stretches often carry no entire MLD block. We chopped
 426 chromosome 1 using a window of 150 polymorphic sites, into 7,082 MLD blocks with an
 427 average length of 307.897 bp.

428 **Distribution of L'_c**

429 The distribution of L'_c for our sample of gorilla sequences has an excess of small and
 430 long MLD blocks compared to the L'_c distribution of a constant population with the same
 431 characteristics (same number of mutations and same average length of MLD blocks) (Fig
 432 8). The excess of small blocks is even larger than what we see in simulated declines
 433 (see above). The truncation of long MLD blocks due to the inclusion of low quality of
 434 genotyping can potentially inflate this excess.

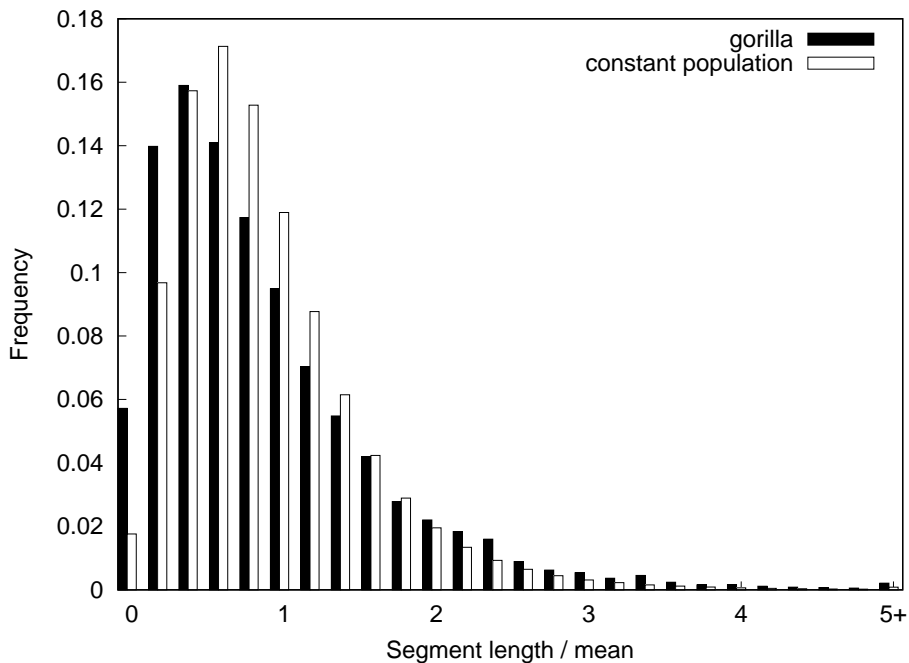


Figure 8: Distribution of L'_c for a population of constant size (white, mutation rate=0.000375 , recombination rate= 0.012) and for the chromosome 1 of the gorillas (black)

435 With the low quality of genotyping and the chosen mutation and recombination rates,
 436 the threshold value $f_c^{5\%}$ is 0.041627. As we measure $f_c^{gor} = 0.0631178$ for the gorillas,

437 the test significantly rejects H_0 . However, it is possible that other designs of similar tests
438 (tweaking the lower or upper bounds) may be more relevant to analyze demography from
439 low quality chromosome alignments.

440 However, and this may be even more important, misspecifications of the model can also
441 make the test significant. Among all, we have chosen to explore the impact of recovery
442 after the decline and of spatial structure.

443 **7 Misspecification of H_1**

444 To appreciate how the f_c test, that was specifically designed to detect population decline, is
445 sensitive to other violations of H_0 , we explored their sensitivity to a scenario of bottleneck
446 (decline followed by recovery) and to a scenario with structure but no demography.

447 **7.1 Bottleneck**

448 In the bottleneck scenario, we model a population that experienced a sudden strong decline
449 ($\kappa = 10$) at time $\tau = 1$ in the past and recovered to its original size after a duration of
450 $x \in [0, 1]$. If $x = 0$, there is no population decline. If $x = 1$, the population has not
451 recovered and the bottleneck scenario is identical to our original H_1 . When the bottleneck
452 lasts long enough ($x > 0.02$), it is detected by the f_c test (Fig 9b). On the contrary, when
453 the bottleneck is too short ($x < 0.02$), the distribution of L'_c is similar to the one under H_0
454 (Fig 9a). This shows that even if the population has recovered, the signal of decline will
455 be observable in the excess of short and long MLD blocks.

456 **7.2 Island-mainland structure**

457 Structured populations generate signals of population size change, even when the popula-
458 tion is stationary (Mazet et al., 2015). For example if the size of the sample is $n = 2$, for
459 any population model with spatial structure, there exists a model without structure but
460 specifically designed variations of population size which has the same distributions of coa-
461 lescence times (Mazet et al., 2016). We consider here a larger sample size ($n = 10$, as in the
462 other scenarios). We assume that genomes are sampled from an island with population size
463 N and the island receives migrants with individual rate m from the mainland, which has
464 population size $10N$. The shape of the distribution of L'_c is impacted by the migration rate

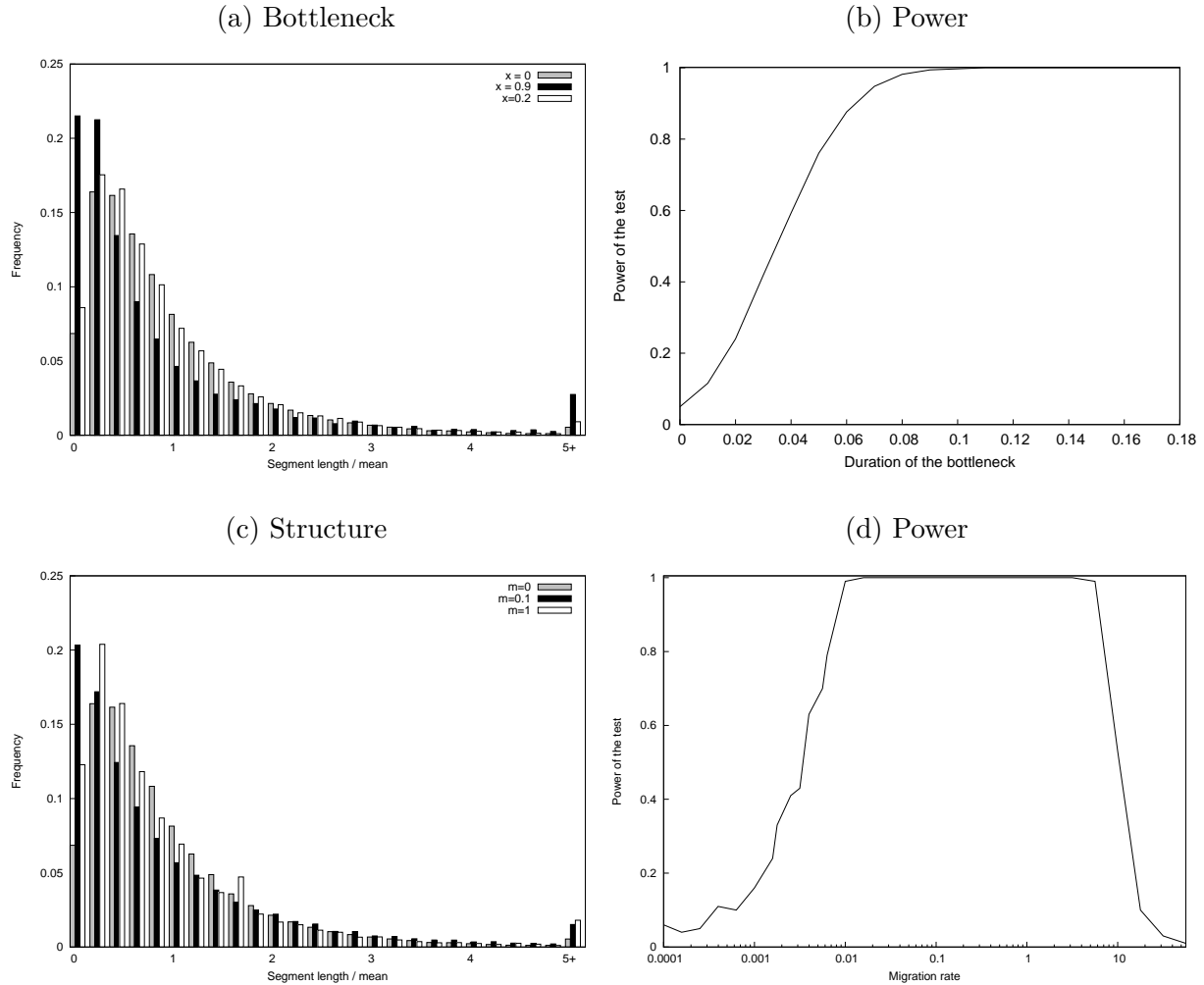


Figure 9: Distribution of L'_c and power of the f_c test under two alternative scenarios. (a) In the bottleneck scenario, the population size is constant equal to N except during the period $[1 - x, 1]$ (measured in units of N generations backwards from the present) during which it equals $N/10$, with $x = \{0, 0.2, 0.9\}$ ($x = 1$ corresponds to decline without recovery). (b) Power of the f_c test on the bottleneck population as a function of the duration of the bottleneck, $x \in [0, 0.18]$. (c) In the island-mainland scenario, the population size is constant equal to N (island) and receives migrants at individual rate $m = \{0, 0.1, 1\}$ from a population of size $10N$ (mainland). (d) Power of f_c test as a function of m , the migration rate from the mainland to the island.

465 and the ratio of population sizes between the island and the mainland (data not shown).
466 When the migration rate gets too small ($m < 0.001$) or too large ($m > 10$), the distribution
467 of L'_c is the same as under H_0 (Fig 9c). For intermediate values (*i.e.* $0.001 < m < 10$),
468 an excess of short and long blocks will be observed (Fig 9d). However, the shape of the
469 distribution of L'_c is visually different from the one of a declining population, that is, the
470 excess of small blocks is higher than under a declining scenario. For a value of $m = 0.1$,
471 the proportion of blocks between 0 and 0.1 times the mean is significantly higher than the
472 proportion of blocks between 0.1 and 0.2 times the mean, which is not observed for the
473 distribution of block lengths under decline. This suggests that the distribution could be
474 used to differentiate between the effects of demography and structure.

475 8 Discussion

476 We have explored the impact of demography, more specifically recent population decline,
477 on the pattern of recombination in a sample of n genomes, where $n \gg 2$. We have shown
478 that the distribution of the distances between recombination breakpoints (MRF block
479 lengths) is strongly affected by the demography. More specifically, a decline will result in
480 an overdispersion of the distribution, that is, a relative excess of short and long blocks.
481 As most recombination breakpoints are difficult, and sometimes impossible, to detect in a
482 sequence alignment, we have proposed to restrict ourselves to the ones that can be detected
483 using the four-gamete test. These detectable breakpoints delineate blocks in full linkage
484 disequilibrium that we named MLD blocks.

485 Although different from the distribution of MRF blocks, the distribution of MLD block
486 lengths is also overdispersed when the population has been declining recently. Using simple
487 tests based on an excess of small and long blocks (f and f_c), one can detect declines for a
488 wide range of different dates and strengths.

489 Surprisingly, the f_c test based on MLD blocks has more power for very recent declines
490 ($\tau \approx 0.01$) than the f test based on MRF blocks. The past demography of the population
491 impacts the distribution of the length L of MRF blocks but also the fraction p of MRF
492 breakpoints that also correspond to MLD breakpoints. When recombination occurs at
493 distant times when only $k \ll n$ ancestor lineages are present (*i.e.* the most ancient times of
494 the tree), it rarely produces incompatibilities detectable with the four-gamete test (never
495 when $k = 2$). For a declining population, these ancient lineages have longer branches

496 than the ones of a constant population scenario, so that recombination events occur more
497 frequently in these lineages. This results in a smaller p for declining populations and thus
498 in more numerous (ancient, small) MRF blocks per MLD block. The relative abundance of
499 long recent MLD blocks becomes thus more important in the distribution. This effect fades
500 away for distant declines. In summary, the effect of recent declines on the L_c distribution
501 is the result of both a change in the L distribution and a change in the fraction p of
502 breakpoints detected, which can explain the difference in power between the f and the f_c
503 tests.

504 We also show that using the f_c statistic, the decline can still be detected even if the
505 population has recently recovered its original size (bottleneck scenario). Finally, we showed
506 that local sampling of a small deme with constant size also leads to rejection of H_0 for f_c
507 but that the distribution of block lengths seems distorted in a way that can help distinguish
508 the two scenarios. We leave this for future work.

509 One interesting advantage of using the f and f_c tests are their efficiency in computing
510 time, such that they can scale up to a very large sample of long genomes. For example, the
511 chopping of the entire human chromosome 1 (1,636,975 SNPs) for a sample of 10 unphased
512 genomes takes 16 seconds on a laptop (with an Intel Core i7 processor running macOS High
513 Sierra). This very short computing time is an interesting asset of this test compared to
514 other methods based on variations (e.g. in SNP density) induced by recombination events
515 (Li and Durbin, 2011; Palamara et al., 2012; MacLeod et al., 2013; Harris and Nielsen,
516 2013; Browning and Browning, 2015). The main choice that influences the computation
517 time of the chopping algorithm is the number of sites considered for the chopping window.
518 An increase in the chopping window size will increase the number of sites to test for
519 incompatibility.

520 In the theoretical assessment of the f_c test, we have made the assumption that the
521 recombination rate is constant along the genome and that entire genomes are aligned. Let
522 us discuss the limits of these assumptions.

523 First, the recombination rate is known to vary along the genome, especially in regions
524 of high recombination known as recombination hot-spots. It could be possible to integrate
525 these variations via the knowledge of the recombination map. Indeed, if the recombination
526 rate is twice higher in a given region of the genome, MRF blocks will be twice smaller, so

527 we can correct this distortion by multiplying all MRF block lengths by 2.

528 Another issue of the test based on MLD block lengths is the need of whole genome
529 data. For normalisation of the MLD block distribution, the average length of a block is
530 needed. If the whole distribution of MLD block is not available, it can compromise the
531 estimate of the average length, and so can compromise the test based on the normalised
532 distribution. The test requires genome data with good SNP quality for all the individuals.

533 The f_c test is a genome-wide approach that can detect population decline that started
534 even very recently, down to orders of $\tau = 0.01N_0$, where N_0 is the current (effective)
535 population size. This corresponds to very recent times, in particular when considering en-
536 dangered populations. For example, there are approximately 600 mature mountain Gorilla
537 individuals alive (IUCN Red List of 31 July 2018). Assuming that the current effective
538 population size is a third of the mature individuals, $N_0 \approx 200$, the f_c test will detect
539 decline as recent as $0.01 * 200 = 2$ generations ago. Great apes populations (Bonobos,
540 Chimpanzees, Orangutans) have been sequenced (Prado-Martinez et al., 2013) and are
541 actively re-sequenced (Gordon et al., 2016). The coverage used to sequence the data cur-
542 rently available is not high enough to apply our test. To infer MLD blocks, the sequenced
543 DNA tracts need to be uninterrupted. Using these whole-genome data in higher quality,
544 we will be able to confirm their decline thanks to the f_c test.

545 Giant pandas have a ‘vulnerable’ conservation status in the Red List. Recently they
546 have seen their population increased (around 500 mature individuals, from IUCN website
547 2019). Applying the f_c test on some genomes of theirs (Zhao et al., 2013) sequenced in
548 higher quality, will give some precise information on their demography. As the test is
549 influenced by duration and strength of a bottleneck, the strength and the date of the
550 increase in the population size impact the result of the test. Applying the f_c test to
551 mammals with approximately known demography will be interesting to verify the method.
552 However, the real asset of this test is its possible application to a much wider range of
553 organisms. Whole-genome data start to become more and more common for non-model,
554 non-vertebrate organisms like honeybee (Wallberg et al., 2014), as well as organisms with
555 no conservation status such as mimicry butterflies (Zhang et al., 2017).

556 The chopping algorithm detects incompatibilities among trees along the genome. All
557 the sites in a MLD block are compatible with one topology. We developed the algorithm

558 to detect a recent change in population size. However, its use is not limited to population
559 demography. Conflicting genealogies are also present in phylogenetic inference (Maddison,
560 1997). This algorithm could be used to partition the genome according to compatible trees
561 before estimating the trees.

562 Recombination and mutation events leave a joint imprint on genomes which depends
563 notably on the demography of the population. Their frequency and locations carry infor-
564 mation about the past history of this population (decline, bottleneck, structure...). Using
565 MLD breakpoints to chop genomes gives insights into this history and may be used to
566 gain further information on other aspects impacting the frequency of recombination events
567 through time and along the genome (e.g. hitch-hiking due to selection).

568 **Acknowledgements**

569 E.K. is funded by the PhD program ‘Interfaces pour le Vivant’ of Sorbonne Université.
570 G.A., A.L. and E.K. thank the *Center for Interdisciplinary Research in Biology* and the
571 *Fondation François Sommer* for funding.

572 **References**

- 573 Adams, A. M. and Hudson, R. R. (2004). Maximum-likelihood estimation of demographic
574 parameters using the frequency spectrum of unlinked single-nucleotide polymorphisms.
575 *Genetics*, 168(3):1699–712.
- 576 Alonso-Blanco, C., Andrade, J., Becker, C., Bemm, F., Bergelson, J., Borgwardt, K. M.,
577 Cao, J., Chae, E., Dezwaan, T. M., Ding, W., Ecker, J. R., Exposito-Alonso, M., Farlow,
578 A., Fitz, J., Gan, X., Grimm, D. G., Hancock, A. M., Henz, S. R., Holm, S., Horton,
579 M., Jarsulic, M., Kerstetter, R. A., Korte, A., Korte, P., Lanz, C., Lee, C.-R., Meng,
580 D., Michael, T. P., Mott, R., Mulyati, N. W., Nägele, T., Nagler, M., Nizhynska, V.,
581 Nordborg, M., Novikova, P. Y., Picó, F. X., Platzer, A., Rabanal, F. A., Rodriguez, A.,
582 Rowan, B. A., Salomé, P. A., Schmid, K. J., Schmitz, R. J., Seren, Ü., Sperone, F. G.,
583 Sudkamp, M., Svardal, H., Tanzer, M. M., Todd, D., Volchenboum, S. L., Wang, C.,
584 Wang, G., Wang, X., Weckwerth, W., Weigel, D., and Zhou, X. (2016). 1,135 genomes

585 reveal the global pattern of polymorphism in *arabidopsis thaliana*. *Cell*, 166(2):481 –
586 491.

587 Ball, F. and Stefanov, V. T. (2005). Evaluation of identity-by-descent probabilities for
588 half-sibs on continuous genome. *Mathematical Biosciences*, 196(2):215 – 225.

589 Barnosky, A. D., Matzke, N., Tomiya, S., Wogan, G. O. U., Swartz, B., Quental, T. B.,
590 Marshall, C., McGuire, J. L., Lindsey, E. L., Maguire, K. C., Mersey, B., and Ferrer, E. A.
591 (2011). Has the earth’s sixth mass extinction already arrived? *Nature*, 471(7336):51–7.

592 Beichman, A. C., Huerta-Sanchez, E., and Lohmueller, K. E. (2018). Using genomic data
593 to infer historic population dynamics of nonmodel organisms. *Annual Review of Ecology,
594 Evolution, and Systematics*, 49(1):433–456.

595 Browning, B. L. and Browning, S. R. (2013). Improving the accuracy and efficiency of
596 identity-by-descent detection in population data. *Genetics*, 194(2):459–71.

597 Browning, S. R. and Browning, B. L. (2010). High-resolution detection of identity by
598 descent in unrelated individuals. *Am J Hum Genet*, 86(4):526 – 539.

599 Browning, S. R. and Browning, B. L. (2015). Accurate non-parametric estimation of
600 recent effective population size from segments of identity by descent. *Am J Hum Genet*,
601 97(3):404–18.

602 Carmi, S., Wilton, P. R., Wakeley, J., and Pe’er, I. (2014). A renewal theory approach to
603 ibd sharing. *Theoretical Population Biology*, 97:35–48.

604 Ceballos, G., Ehrlich, P. R., and Dirzo, R. (2017). Biological annihilation via the ongoing
605 sixth mass extinction signaled by vertebrate population losses and declines. *Proceedings
606 of the National Academy of Sciences*, 114(30):E6089–E6096.

607 Chapman, N. and Thompson, E. (2003). A model for the length of tracts of identity by
608 descent in finite random mating populations. *Theor Popul Biol*, 64(2):141 – 150.

609 Chikhi, L., Rodríguez, W., Grusea, S., Santos, P., Boitard, S., and Mazet, O. (2018). The
610 iicr (inverse instantaneous coalescence rate) as a summary of genomic diversity: insights
611 into demographic inference and model choice. *Heredity*, 120(1):13–24.

- 612 Díez-del Molino, D., Sánchez-Barreiro, F., Barnes, I., Gilbert, M. T. P., and Dalén, L.
613 (2018). Quantifying temporal genomic erosion in endangered species. *Trends in Ecology*
614 *and Evolution*, 33(3):176–185.
- 615 Donnelly, K. P. (1983). The probability that related individuals share some section of
616 genome identical by descent. *Theoretical Population Biology*, 23(1):34 – 63.
- 617 Excoffier, L., Dupanloup, I., Huerta-Sánchez, E., Sousa, V. C., and Foll, M. (2013). Robust
618 demographic inference from genomic and snp data. *PLoS Genet*, 9(10):e1003905.
- 619 Fu, Y. X. (1995). Statistical properties of segregating sites. *Theor Popul Biol*, 48(2):172–97.
- 620 Gibbs, R. A., Boerwinkle, E., Doddapaneni, H., Han, Y., Korchina, V., Kovar, C., Lee, S.,
621 Muzny, D., Reid, J. G., and et al. (2015). A global reference for human genetic variation.
622 *Nature*, 526(7571):68–74.
- 623 Gordon, D., Huddleston, J., Chaisson, M. J. P., Hill, C. M., Kronenberg, Z. N., Munson,
624 K. M., Malig, M., Raja, A., Fiddes, I., Hillier, L. W., Dunn, C., Baker, C., Armstrong,
625 J., Diekhans, M., Paten, B., Shendure, J., Wilson, R. K., Haussler, D., Chin, C.-S.,
626 and Eichler, E. E. (2016). Long-read sequence assembly of the gorilla genome. *Science*,
627 352(6281).
- 628 Griffiths, R. and Marjoram, P. (1997). An ancestral recombination graph. In *Progress in*
629 *population genetics and human evolution*, pages 257 – 270. Springer.
- 630 Grusea, S., Rodríguez, W., Pinchon, D., Chikhi, L., Boitard, S., and Mazet, O. (2019).
631 Coalescence times for three genes provide sufficient information to distinguish population
632 structure from population size changes. *J Math Biol*, 78(1-2):189–224.
- 633 Gusev, A., Lowe, J. K., Stoffel, M., Daly, M. J., Altshuler, D., Breslow, J. L., Fried-
634 man, J. M., and Pe’er, I. (2009). Whole population, genome-wide mapping of hidden
635 relatedness. *Genome research*, 19(2):318–26.
- 636 Gutenkunst, R. N., Hernandez, R. D., Williamson, S. H., and Bustamante, C. D. (2009).
637 Inferring the joint demographic history of multiple populations from multidimensional
638 snp frequency data. *PLoS Genet*, 5(10):e1000695.

- 639 Harris, K. and Nielsen, R. (2013). Inferring demographic history from a spectrum of shared
640 haplotype lengths. *PLoS Genet*, 9(6):e1003521.
- 641 Hayes, B. J., Visscher, P. M., McPartlan, H. C., and Goddard, M. E. (2003). Novel
642 multilocus measure of linkage disequilibrium to estimate past effective population size.
643 *Genome research*, 13(4):635–43.
- 644 Hey, J. and Wakeley, J. (1997). A coalescent estimator of the population recombination
645 rate. *Genetics*, 145(3):833–46.
- 646 Hill, W. G. and Robertson, A. (1968). Linkage disequilibrium in finite populations. *Theor*
647 *Appl Genet*, 38(6):226–31.
- 648 Hollenbeck, C. M., Portnoy, D. S., and Gold, J. R. (2016). A method for detecting recent
649 changes in contemporary effective population size from linkage disequilibrium at linked
650 and unlinked loci. *Heredity*, 117(4):207–16.
- 651 Hudson, R. R. and Kaplan, N. L. (1985). Statistical properties of the number of recombina-
652 tion events in the history of a sample of dna sequences. *Genetics*, 111(1):147–64.
- 653 Kelleher, J., Etheridge, A. M., and McVean, G. (2016). Efficient coalescent simulation and
654 genealogical analysis for large sample sizes. *PLoS Comput Biol*, 12(5):e1004842.
- 655 Kingman, J. (1982). The coalescent. *Stochastic Processes and their Applications*, 13(3):235–
656 248.
- 657 Lapiere, M., Lambert, A., and Achaz, G. (2017). Accuracy of demographic inferences from
658 the site frequency spectrum: The case of the yoruba population. *Genetics*, 206(1):439–
659 449.
- 660 Lewontin, R. C. and ichi Kojima, K. (1960). The evolutionary dynamics of complex poly-
661 morphisms. *Evolution*, 14(4):458–472.
- 662 Li, H. and Durbin, R. (2011). Inference of human population history from individual
663 whole-genome sequences. *Nature*, 475(7357):493–496.
- 664 MacLeod, I. M., Larkin, D. M., Lewin, H. A., Hayes, B. J., and Goddard, M. E. (2013).
665 Inferring demography from runs of homozygosity in whole-genome sequence, with cor-
666 rection for sequence errors. *Mol Biol Evol*, 30(9):2209–23.

- 667 MacLeod, I. M., Meuwissen, T. H. E., Hayes, B. J., and Goddard, M. E. (2009). A novel
668 predictor of multilocus haplotype homozygosity: comparison with existing predictors.
669 *Genetics research*, 91(6):413–26.
- 670 Maddison, W. P. (1997). Gene trees in species trees. *Systematic Biology*, 46(3):523–536.
- 671 Marjoram, P. and Wall, J. D. (2006). Fast "coalescent" simulation. *BMC Genet*, 7:16.
- 672 Marth, G. T., Czabarka, E., Murvai, J., and Sherry, S. T. (2004). The allele frequency spec-
673 trum in genome-wide human variation data reveals signals of differential demographic
674 history in three large world populations. *Genetics*, 166(1):351–72.
- 675 Mazet, O., Rodríguez, W., and Chikhi, L. (2015). Demographic inference using genetic data
676 from a single individual: Separating population size variation from population structure.
677 *Theor Popul Biol*, 104:46–58.
- 678 Mazet, O., Rodríguez, W., Grusea, S., Boitard, S., and Chikhi, L. (2016). On the impor-
679 tance of being structured: instantaneous coalescence rates and human evolution—lessons
680 for ancestral population size inference? *Heredity*, 116(4):362–71.
- 681 McVean, G. A. T. and Cardin, N. J. (2005). Approximating the coalescent with recombina-
682 tion. *Philos Trans R Soc Lond B Biol Sci*, 360(1459):1387–93.
- 683 Palamara, P. F., Lencz, T., Darvasi, A., and Pe'er, I. (2012). Length distributions of
684 identity by descent reveal fine-scale demographic history. *Am J Hum Genet*, 91(5):809–
685 22.
- 686 Patin, E., Siddle, K. J., Laval, G., Quach, H., Harmant, C., Becker, N., Froment, A.,
687 Régnauld, B., Lemée, L., Gravel, S., and et al. (2014). The impact of agricultural emer-
688 gence on the genetic history of african rainforest hunter-gatherers and agriculturalists.
689 *Nature Communications*, 5(1).
- 690 Prado-Martinez, J., Sudmant, P. H., Kidd, J. M., Li, H., Kelley, J. L., Lorente-Galdos,
691 B., Veeramah, K. R., Woerner, A. E., O'Connor, T. D., Santpere, G., and et al. (2013).
692 Great ape genetic diversity and population history. *Nature*, 499(7459):471–475.
- 693 Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M. A., Bender, D., Maller,
694 J., Sklar, P., de Bakker, P. I., Daly, M. J., and Sham, P. C. (2007). Plink: A tool set

695 for whole-genome association and population-based linkage analyses. *Am J Hum Genet*,
696 81(3):559 – 575.

697 Régnier, C., Achaz, G., Lambert, A., Cowie, R. H., Bouchet, P., and Fontaine, B. (2015).
698 Mass extinction in poorly known taxa. *Proceedings of the National Academy of Sciences*,
699 112(25):7761–6.

700 Ringbauer, H., Coop, G., and Barton, N. H. (2017). Inferring recent demography from
701 isolation by distance of long shared sequence blocks. *Genetics*, 205(3):1335–1351.

702 Rodrigues, A. S. L., Pilgrim, J. D., Lamoreux, J. F., Hoffmann, M., and Brooks, T. M.
703 (2006). The value of the iucn red list for conservation. *Trends Ecol Evol*, 21(2):71–6.

704 Rodríguez, W., Mazet, O., Grusea, S., Arredondo, A., Corujo, J. M., Boitard, S., and
705 Chikhi, L. (2018). The iicr and the non-stationary structured coalescent: towards demo-
706 graphic inference with arbitrary changes in population structure. *Heredity*, 121(6):663–
707 678.

708 Sánchez-Bayo, F. and Wyckhuys, K. A. (2019). Worldwide decline of the entomofauna: A
709 review of its drivers. *Biological Conservation*, 232:8–27.

710 Schiffels, S. and Durbin, R. (2014). Inferring human population size and separation history
711 from multiple genome sequences. *Nat Genet*, 46(8):919–25.

712 Sheehan, S., Harris, K., and Song, Y. S. (2013). Estimating variable effective population
713 sizes from multiple genomes: A sequentially markov conditional sampling distribution
714 approach. *Genetics*, 194(3):647–662.

715 Stam, P. (1980). The distribution of the fraction of the genome identical by descent in
716 finite random mating populations. *Genetical Research*, 35(2):131–155.

717 Stefanov, V. T. (2000). Distribution of genome shared identical by descent by two individ-
718 uals in grandparent-type relationship. *Genetics*, 156(3):1403–10.

719 Terhorst, J., Kamm, J. A., and Song, Y. S. (2017). Robust and scalable inference of
720 population history from hundreds of unphased whole genomes. *Nat Genet*, 49(2):303–
721 309.

- 722 Tired, M. and Hospital, F. (2017). Blocks of chromosomes identical by descent in a popu-
723 lation: Models and predictions. *PLoS one*, 12(11):e0187416.
- 724 van der Valk, T., Díez-del Molino, D., Marques-Bonet, T., Guschanski, K., and Dalén, L.
725 (2019). Historical genomes reveal the genomic consequences of recent population decline
726 in eastern gorillas. *Current Biology*, 29(1):165–170.e6.
- 727 Wallberg, A., Han, F., Wellhagen, G., Dahle, B., Kawata, M., Haddad, N., Simões, Z. L. P.,
728 Allsopp, M. H., Kandemir, I., De la Rúa, P., Pirk, C. W., and Webster, M. T. (2014).
729 A worldwide survey of genome sequence variation provides insight into the evolutionary
730 history of the honeybee *apis mellifera*. *Nat Genet*, 46(10):1081–8.
- 731 Watterson, G. A. (1975). On the number of segregating sites in genetical models without
732 recombination. *Theor Popul Biol*, 7(2):256–76.
- 733 Wiuf, C. and Hein, J. (1999). The ancestry of a sample of sequences subject to recombina-
734 tion. *Genetics*, 151(3):1217–28.
- 735 Wright, S. (1931). Evolution in mendelian populations. *Genetics*, 16(2):97–159.
- 736 Zhang, W., Westerman, E., Nitzany, E., Palmer, S., and Kronforst, M. R. (2017). Tracing
737 the origin and evolution of supergene mimicry in butterflies. *Nature Communications*,
738 8(1).
- 739 Zhao, S., Zheng, P., Dong, S., Zhan, X., Wu, Q., Guo, X., Hu, Y., He, W., Zhang, S.,
740 Fan, W., Zhu, L., Li, D., Zhang, X., Chen, Q., Zhang, H., Zhang, Z., Jin, X., Zhang, J.,
741 Yang, H., Wang, J., Wang, J., and Wei, F. (2013). Whole-genome sequencing of giant
742 pandas provides insights into demographic history and local adaptation. *Nat Genet*,
743 45(1):67–71.