



**HAL**  
open science

# Heterogeneity in effective size across the genome: effects on the inverse instantaneous coalescence rate (IICR) and implications for demographic inference under linked selection

Simon Boitard, Armando Arredondo, Lounès Chikhi, Olivier Mazet

## ► To cite this version:

Simon Boitard, Armando Arredondo, Lounès Chikhi, Olivier Mazet. Heterogeneity in effective size across the genome: effects on the inverse instantaneous coalescence rate (IICR) and implications for demographic inference under linked selection. *Genetics*, 2022, 220 (3), pp.iyac008. 10.1093/genetics/iyac008 . hal-03649481

**HAL Id: hal-03649481**

**<https://hal.inrae.fr/hal-03649481>**

Submitted on 22 Apr 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 4.0 International License

1 Heterogeneity in effective size across the  
2 genome: effects on the Inverse Instantaneous  
3 Coalescence Rate (IICR) and implications for  
4 demographic inference under linked selection

5 Simon Boitard\*, Armando Arredondo†, Camille Noûs‡,

6 Lounès Chikhi§,\*\*, Olivier Mazet†

7 \*: CBGP, Université de Montpellier, CIRAD, INRAE, Institut Agro, IRD,  
8 Montpellier, France.

9 †: Université de Toulouse, Institut National des Sciences Appliquées, In-  
10 stitut de Mathématiques de Toulouse, Toulouse, France.

11 ‡: Laboratoire Cogitamus, Toulouse, France.

12 §: Instituto Gulbenkian de Ciência, Oeiras, Portugal

13 \*\*: Laboratoire Évolution & Diversité Biologique (EDB UMR 5174), CNRS,  
14 IRD, UPS, Université de Toulouse Midi-Pyrénées, Toulouse, France

15 **Running title:**

16 The IICR under linked selection

17 **Keywords:**

18 demographic inference, linked selection, effective population size, coales-  
19 cence times, population structure, drosophila melanogaster, humans

20 **Corresponding author:**

21 Simon Boitard

22 CBGP, 755 avenue du Campus Agropolis, CS 30016, 34988 Montferrier sur

23 Lez cedex, France

24 [simon.boitard@inrae.fr](mailto:simon.boitard@inrae.fr)

25 0033 4 99 62 33 36

## 26 Abstract

27 The relative contribution of selection and neutrality in shaping species genetic diversity is  
28 one of the most central and controversial questions in evolutionary theory. Genomic data  
29 provide growing evidence that linked selection, i.e. the modification of genetic diversity  
30 at neutral sites through linkage with selected sites, might be pervasive over the genome.  
31 Several studies proposed that linked selection could be modelled as first approximation  
32 by a local reduction (e.g. purifying selection, selective sweeps) or increase (e.g. balanc-  
33 ing selection) of effective population size ( $N_e$ ). At the genome-wide scale, this leads to  
34 variations of  $N_e$  from one region to another, reflecting the heterogeneity of selective con-  
35 straints and recombination rates between regions. We investigate here the consequences of  
36 such genomic variations of  $N_e$  on the genome-wide distribution of coalescence times. The  
37 underlying motivation concerns the impact of linked selection on demographic inference,  
38 because the distribution of coalescence times is at the heart of several important demo-  
39 graphic inference approaches. Using the concept of Inverse Instantaneous Coalescence  
40 Rate, we demonstrate that in a panmictic population, linked selection always results in a  
41 spurious apparent decrease of  $N_e$  along time. Balancing selection has a particularly large  
42 effect, even when it concerns a very small part of the genome. We also study more gen-  
43 eral models including genuine population size changes, population structure or transient  
44 selection and find that the effect of linked selection can be significantly reduced by that of  
45 population structure. The models and conclusions presented here are also relevant to the  
46 study of other biological processes generating apparent variations of  $N_e$  along the genome.

## 47 Introduction

48 One of the greatest challenges of evolutionary biology is to understand how natural se-  
49 lection, mutation, recombination and genetic drift have shaped and are still shaping the  
50 patterns of genomic diversity of species living today (Charlesworth, 2010, Lewontin, 1974,  
51 Walsh and Lynch, 2018). In the last decade genomic data have become increasingly avail-  
52 able for both model and non-model species. It is expected that by analysing these genomic  
53 data we will be able to better understand the respective roles of the different evolutionary  
54 forces (Charlesworth, 2010, Lewontin, 1974). In particular, it is believed that we will be  
55 able to identify the regions that have been shaped by selection, and those that may be  
56 more neutral (Johri et al., 2020, Pouyet et al., 2018). The relative importance of selection  
57 and neutrality in generating the genomic patterns of diversity we see today has been at  
58 the heart of many evolutionary debates and controversies over the last decades (Kimura,  
59 1983, Lewontin, 1974, Ohta, 1992) and recent studies suggest that it still is (Comeron,  
60 2017, Jensen et al., 2019, Kern and Hahn, 2018).

61 The concept of effective size ( $N_e$ ) is central to these debates (Charlesworth, 2009)  
62 because selection is expected to be more efficient when  $N_e$  is large, and genetic drift to  
63 be the main driver of evolutionary change when  $N_e$  is small (Ohta, 1992). For instance,  
64 Charlesworth (2009) notes that an autosomal locus under positive selection will behave  
65 neutrally when  $s < 1/4N_e$ , where  $s$  is the selection intensity at this locus. At the same  
66 time it is commonly assumed that selection will itself imply a variation of  $N_e$  across the  
67 genome (Charlesworth, 2009, Gossmann et al., 2011, Jiménez-Mena et al., 2016b). For  
68 instance, Gossmann et al. (2011) write that “*The effective population size is expected to*  
69 *vary across the genome as a consequence of genetic hitchhiking (Smith and Haigh, 1974)*  
70 *and background selection (Charlesworth et al., 1993)*”. They add that “*The action of both*

71 *positive and negative natural selection, is expected to reduce the effective population size*  
72 *leading to lower levels of genetic diversity and reduced effectiveness of selection.” They*  
73 also stress that “*The evidence that there is variation in  $N_e$  within a genome comes from*  
74 *three sources. First, it has been shown that levels of neutral genetic diversity are correlated*  
75 *to rates of recombination in *Drosophila* [...], humans [...], and some plant species...”*. In  
76 his 2009 review on the concept of  $N_e$  Charlesworth (2009) made a similar comment: “ *$N_e$*   
77 *may also vary across different locations in the genome of a species [...] because of the*  
78 *effects of selection at one site in the genome on the behaviour of variants at nearby sites”*.  
79 More recently, Jiménez-Mena et al. (2016a) stated that “*recent studies [...] suggest that*  
80 *different segments of the genome might undergo different rates of genetic drift, potentially*  
81 ***challenging the idea that a single  $N_e$  can account for the evolution of the***  
82 ***genome”*** (emphasis ours).

83 Under these explicit or implicit modelling frameworks, genomic regions with limited  
84 genetic diversity are thus seen as regions of low  $N_e$  as a result of selective sweeps (Smith  
85 and Haigh, 1974) or background selection (Charlesworth et al., 1993), whereas regions  
86 with very high levels of genetic diversity may be seen as regions of large  $N_e$  and could  
87 be explained by balancing selection (Charlesworth, 2009) (see also Hill and Robertson  
88 (1966)). Following that rationale, Jiménez-Mena et al. (2016b) suggested that different  
89 species might thus differ in the statistical distribution of  $N_e$  across the genome and they  
90 presented such distributions for eleven species.

91 Given the central role played by the  $N_e$  concept to detect, identify, and even *conceptual-*  
92 *ize* selection, it may be important, perhaps even enlightening, to explore the consequences  
93 of the ideas presented above with the concept of IICR (inverse instantaneous coalescence  
94 rate) recently introduced by Mazet et al. (2016). Indeed, the IICR is equivalent to the  
95 past temporal trajectory of  $N_e$ , previously defined as the coalescent  $N_e$  (Sjödín et al.,

2005), in a panmictic population under neutrality, and it is the quantity estimated by the popular PSMC method of Li and Durbin (2011). The IICR was first defined by Mazet et al. (2016) for a sample size of two and its properties were studied under several models of population structure (Chikhi et al., 2018, Grusea et al., 2018, Rodríguez et al., 2018). It can also be used for demographic inference under neutrality and models of population structure (Arredondo et al., 2021, Chikhi et al., 2018). These studies showed that the IICR will significantly change over time when populations are structured, even when population size is actually constant. They also outlined that the IICR not only depends on the model of population structure but also on the sampling scheme, which questions the notion that an  $N_e$  can be easily associated to (or is a property of) the model of interest when the model is structured (Chikhi et al., 2018, Rodríguez et al., 2018). The reason for this dependency is that the IICR is by definition a function of the distribution of coalescence times for two genes ( $T_2$ ), which is itself a function of both the evolutionary model and the location (in time and space) of the sampled genes.

One important assumption of the IICR studies mentioned above is that this distribution of  $T_2$  is homogeneous along the genome. The IICR, as defined and computed in previous studies, is thus a genomic average assuming that all loci follow a single Wright-Fisher model, with or without population structure, but with the same number of haploid genes. Whichever definition of  $N_e$  one assumes, the underlying model assumes that  $N_e$  is constant along the genome. If we now assume that  $N_e$  varies across the genome as a consequence of selection (even as an approximation) then the variance of coalescence times should be different from that expected under a standard Wright-Fisher model, and the IICR should be a function of the underlying distribution of the  $N_e$  values across the sampled genes. Genomic regions under different selection regimes might then exhibit specific signatures leading to differing IICR curves for each region. Alternatively, these regions

121 might not be easy to identify but they might still influence the average genomic IICR  
122 estimated from sequenced genomes. In the present study we thus wish to explore ideas  
123 related to drift, selection and patterns of genomic diversity by studying the consequences  
124 of this putative genomic variation of  $N_e$  on the IICR.

125 We first study the IICR under panmixia and constant population size but assuming  
126 that  $N_e$  varies across the genome as a result of recurrent selection, using hypothetical  
127 distributions of  $N_e$  and distributions inferred from genomic data. We then generalise the  
128 model to integrate temporal population size variations, population structure or transient  
129 selection effects. Finally, we compare IICR predictions with PSMC estimations obtained  
130 from simulated data under a model including variations of  $N_e$  along the genome. Alto-  
131 gether, we advocate the use of the IICR as a concept that may help clarify what  $N_e$  means  
132 and as one way, among others, to improve our understanding of the recent and ancient  
133 evolutionary history of species.

## 134 **The IICR under panmixia with several classes of (con-** 135 **stant size) $N_e$ along the genome**

### 136 **Methods: model description**

137 We assume that the genome can be divided in  $K$  distinct classes, each of them charac-  
138 terized by a different  $N_e$  that is constant over time. To model these differences of  $N_e$ ,  
139 we consider that each class  $i$  ( $i = 1 \dots K$ ) evolves under a constant size Wright-Fisher  
140 (WF) model (i.e. panmictic with non-overlapping generations) with diploid population  
141 size  $\lambda_i N$  ( $2 \lambda_i N$  haploids), for some reference population size  $N$  corresponding to the  
142 actual number of diploids. Note that  $2N$  represents an actual number of haploid genomes



143 and that under the WF model, there is no ambiguity and  $N$  represents the  $N_e$  under  
144 neutrality. Thus,  $\lambda_i$  reflects the ratio of effective population size  $N_e$  in class  $i$  relative to  
145  $N$  and for convenience we may sometimes refer to  $\lambda_i$  as *the* effective population size in  
146 class  $i$ . Assuming that  $N$  is large (i.e. that all  $\lambda_i N$  are large), we rescale time by units of  
147  $2N$  generations and study the pairwise coalescence time resulting from this model. For  
148 two sequences sampled in the present (at time  $t = 0$ ) for a locus from the  $i^{\text{th}}$  class of the  
149 genome, we know from standard coalescent theory that the coalescence time  $T_2^i$  follows  
150 an exponential distribution with parameter  $\mu_i = \frac{1}{\lambda_i}$ , whose probability density function  
151 (pdf) is

$$f_i(t) = \mu_i e^{-\mu_i t}, i = 1 \dots K.$$

152 Denoting by  $a_i$  the proportion of the genome corresponding to class  $i$ , the pdf of the  
153 coalescence time  $T_2$  at a random locus is thus

$$f(t) = \sum_{i=1}^K a_i f_i(t) = \sum_{i=1}^K a_i \mu_i e^{-\mu_i t}. \quad (1)$$

154 One may also see this distribution as the one we would obtain if we were able to sample  
155 a large number of independent coalescence times along the genome while covering each  
156 class  $i$  according to its true proportion  $a_i$ . In the next section we study the properties of  
157 the IICR under this model.

## 158 **Results: IICR expression and main properties under panmixia**

159 The IICR is a theoretical function that is intrinsically related to the expected distribution  
160 of coalescence times. Denoting  $F$  the cumulative distribution function of  $T_2$  for a given  
161 evolutionary model and sampling scheme, and  $f(t) = F'(t)$  its pdf, the IICR of a sample

162 of size 2 is defined Mazet et al. (2016) as:

$$\text{IICR}(t) = \frac{R(t)}{f(t)}$$

163 where

$$R(t) = \mathbb{P}(T_2 \geq t) = 1 - F(t).$$

164 This theoretical quantity can be evaluated for any coalescent model by simulating a large  
165 number of independent  $T_2$  values and computing their empirical distribution (Chikhi et al.,  
166 2018). For a large class of models, it can also be obtained exactly using analytical or  
167 numerical approaches (Rodríguez et al., 2018). When analyzing a pair of real sequences,  
168 the evolutionary model that generated these sequences is unknown but the associated  
169 IICR can be estimated by SMC approaches like PSMC or MSMC (Schiffels and Durbin,  
170 2013), which exploit the correlation structure of polymorphic sites along the genome to  
171 infer local coalescence times and their genome-wide distribution.

172 For our model with  $K$  different  $\lambda_i$ , we have from equation (1):

$$\text{IICR}(t) = -\frac{R(t)}{R'(t)} = \frac{\sum_{i=1}^K a_i e^{-\mu_i t}}{\sum_{i=1}^K a_i \mu_i e^{-\mu_i t}}. \quad (2)$$

173 It is straightforward to see that the IICR is not constant as soon as there are at least  
174 two different values of  $\lambda_i$  with non null proportion  $a_i$  across the genome. To be more  
175 specific, we prove in the Supplementary Material that the IICR defined in formula (2) is  
176 *always increasing* from  $t = 0$  to  $t = +\infty$  (i.e. backward in time). Thus, in a stationary  
177 panmictic population, the existence of at least two distinct  $N_e$  across the genome ( $\lambda_i, i >$   
178 1) is sufficient to infer a decreasing IICR (forward in time). In this situation, classical  
179 interpretations of PSMC plots under panmixia will lead to the wrong conclusion that

180 the population size decreased through time. Alternatively, this signal could be (also  
181 wrongly) interpreted as the presence of population structure, since population structure  
182 can generate similar changes in the IICR (Mazet et al., 2016).

183 The magnitude of the IICR decrease can also be deduced from formula (2). Indeed,  
184 the value of the IICR at present is

$$\text{IICR}(0) = \frac{1}{\sum_{i=1}^K a_i \mu_i} = \frac{1}{\sum_{i=1}^K \frac{a_i}{\lambda_i}} \quad (3)$$

185 and the limit value when  $t \rightarrow +\infty$  is equal to

$$\frac{1}{\mu_{i_0}} = \lambda_{i_0} = \max_{i=1\dots K}(\lambda_i). \quad (4)$$

186 The present time value  $\text{IICR}(0)$  is thus necessarily between the smallest and largest  $\lambda_i$ ,  
187 as it is the harmonic mean of the  $\lambda_i$ s weighted by their respective proportions  $a_i$ . The  
188 asymptotic value  $\text{IICR}(+\infty)$  is always the largest  $\lambda_i$  found in the genome, *independent*  
189 of its proportion. In other words, even if a minute proportion of the genome has a high  
190  $\lambda_i$  due to balancing selection, under panmixia the IICR will necessarily plateau to this  
191 value in the ancient past. One intuitive explanation for the IICR growing (backward  
192 in time) towards the largest  $\lambda_i$  is that the genes that are characterized by a large  $N_e$   
193 have much larger coalescence times than the rest of the genome. They thus contribute  
194 proportionately more to the most ancient part of the IICR curve.

## 195 **Results: a two-class panmictic model**

196 These properties can be observed in Figure 1 where we represent the simplest case with  
197  $K = 2$  classes of genomic regions. In this figure we present the IICRs for  $\lambda_1 = 0.1$  and  
198  $\lambda_2 = 1$ , for proportions of  $\lambda_2$  (represented by the parameter  $a_2$ ) varying from 0 to 1.

199 Consistent with the choice made in most studies inferring past population size changes,  
200 time is plotted in log10 scale in this Figure and all others shown in the main text.

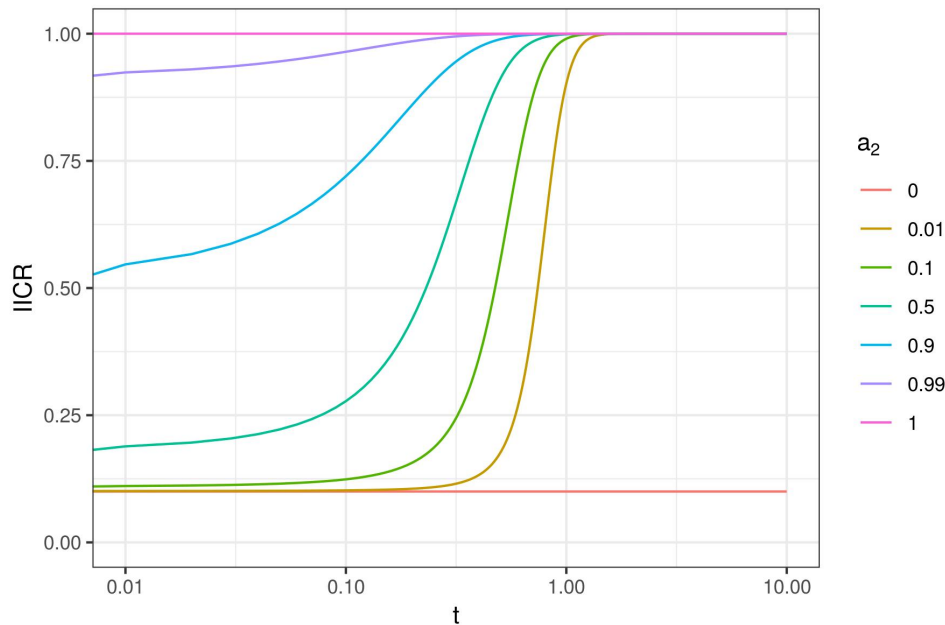


Figure 1: IICR curves for a panmictic model with  $K = 2$  classes of genomic regions with constant size. Genomic regions of class  $i$  ( $i = 1, 2$ ) have a constant population size  $\lambda_i N$ , with  $\lambda_1 = 0.1$  and  $\lambda_2 = 1$ . Their frequencies are  $a_1$  and  $a_2$ , respectively, with  $a_1 + a_2 = 1$ . The IICR curves are represented for  $a_2$  values (representing neutrality, see main text) varying between zero and one. Time is plotted in log10 scale.

201 To simplify the interpretation of our results, we consider (by convention) throughout  
202 this manuscript that  $\lambda_i = 1$  corresponds to the neutral regions of the genome, whether  
203  $a_i$ , their relative proportion in the genome, is large or not. We thus do not necessarily  
204 consider that most of the genome is neutral in that sense. In this setting and in Figure 1,  
205 where  $\lambda_1 = 0.1$  and  $\lambda_2 = 1$ ,  $a_1$  can be interpreted as the fraction of the genome showing  
206 reduced  $N_e$  by a multiplicative factor  $\lambda_1 = 0.1$  as a consequence of positive or background

207 selection.

208 Figure 1 shows that for small values of  $a_2$  (i.e. when most of the genome is under  $N_e$ -  
209 reducing selection) the IICR is S-shaped, slowly increasing backward from  $\lambda_1 = 0.1$  in the  
210 recent past to a plateau at  $\lambda_2 = 1$  in the ancient past. For increasing  $a_2$  values the IICR  
211 curves are becoming flatter as their left-most section flattens upward. Consistent with the  
212 properties outlined in previous section, these curves start (in recent times) at increasing  
213 IICR values above  $\lambda_1 = 0.1$  when the value of  $a_2$  increases, but the curves always reach  
214 the same ancient plateau at  $\lambda_2 = 1$ . However, and this is an important point, this plateau  
215 is reached earlier as  $a_2$  increases. When  $a_2=1$ , only the plateau remains and the IICR is  
216 flat at  $\lambda_2 = 1$  and when  $a_2 = 0$ , it is a flat at  $\lambda_1 = 0.1$ . Thus, when there is only one  $\lambda_i$   
217 over the genome, the IICR is constant over time and equal to that value, as expected for  
218 a population with constant size  $\lambda_i N$  (Li and Durbin, 2011, Mazet et al., 2016).

219 If we now assume that the only type of selection present in the genome increases the  
220 effective size by an order of magnitude, with  $a_1$  and  $a_2$  corresponding to  $\lambda_1 = 1$  and  $\lambda_2 =$   
221 10, we obtain exactly the same figure with the only difference that it is rescaled (Figure  
222 S1). This figure now shows that even if most of the genome is neutral, tiny amounts of  
223  $N_e$  increasing selection strongly influence the IICR, as it always grows backward towards  
224 the plateau corresponding to the largest of the two  $\lambda_i$  values.

225 Altogether Figures 1 and S1 suggest that there is a strong asymmetry between selection  
226 reducing (background and positive) or increasing (balancing)  $N_e$  in the genome in the way  
227 they affect IICR shapes. Balancing selection generates an ancient and high plateau at the  
228 level of  $\lambda_2$ , even for small proportions of  $a_2$  (Figure S1), whereas positive and background  
229 selection generate a recent and relatively more modest decrease of the IICR for small  
230 values of  $a_1$ , even assuming, as in Figure 1, that these generate a ten-fold decrease in  $N_e$   
231 (Figure 1).

## 232 **Results: a three-class panmictic model**

233 To further explore the influence of both types of selection (reducing and increasing  $N_e$ ),  
234 we considered a model with 3 classes such that  $\lambda_1 < 1$ ,  $\lambda_2 = 1$  and  $\lambda_3 > 1$  (Figure 2). In  
235 this Figure we set the three  $\lambda_i$  as  $(\lambda_1, \lambda_2, \lambda_3) = (0.1, 1, 3)$ . As above,  $\lambda_1 < 1$  corresponds  
236 to genomic regions under positive or background selection,  $\lambda_2 = 1$  corresponds to the  
237 neutral part of the genome and  $\lambda_3 = 3$  to genomic regions under balancing selection. In  
238 the left panel, we considered a fixed small proportion of balancing selection ( $a_3 = 0.01$ ),  
239 and allowed the proportions of neutral and positive or background selection to vary ( $a_1$   
240 varied from 0 to 0.8, and thus  $a_2$  from 0.99 to 0.19). In the right panel, we considered a  
241 fixed and large proportion of positive or background selection ( $a_1 = 0.5$ ) and varied the  
242 proportion of regions under balancing selection ( $a_3$  from 0 to 0.1), and thus the proportion  
243 of neutral regions too ( $a_2$  between 0.5 and 0.4).

244 Figure 2 shows similarities with Figure 1. Specifically, both figures suggest that regions  
245 reducing  $N_e$  impact the IICR curves in the recent past whereas regions increasing  $N_e$   
246 impact the IICR in the ancient past. This is worth stressing given that our model assumes  
247 here that  $N_e$  is reduced (in class 1) or increased (in class 3) in a stationary way throughout  
248 the genealogical history of the sampled genes (see the sections on transient selection for  
249 a different assumption). Also, small proportions of balancing selection seem to generate  
250 much bigger changes than small proportions of positive or background selection, as shown  
251 by the comparison of the IICRs obtained for  $a_1 = 0.01$  vs  $a_1 = 0$  on one hand (left panel)  
252 and for  $a_3 = 0.01$  vs  $a_3 = 0$  on the other hand (right panel).

253 There are however differences between Figure 2 and Figure 1. The simple fact that we  
254 consider both  $N_e$ -reducing and  $N_e$ -increasing forms of selection generates complex IICR  
255 curves, in which both forms of selection directly or indirectly impact the whole IICR  
256 curves. When neutral regions are frequent enough ( $a_1 \leq 0.5$  and  $a_3 \leq 0.01$ ), the IICR

257 exhibits a plateau or a flattening at  $\lambda_2$  in its middle section, but for larger values of either  
258  $a_1$  (left panel,  $a_1 = 0.8$ ) or  $a_3$  (right panel,  $a_3 = 0.1$ ) the proportion of neutral genomic  
259 regions decreases and the IICR curve only exhibits a short inflexion corresponding to  
260  $\lambda_2 = 1$  before increasing backwards towards  $\lambda_3$ . An interesting pattern related to this  
261 intermediate plateau is observed on the left panel when  $a_3$  is fixed: the IICR in the  
262 ancient past increases more and quicker (backward in time) for  $a_1 = 0.8$  than for lower  
263 values of  $a_1$ , although  $a_1$  models the proportion of low  $N_e$  regions in the region. This  
264 counterintuitive result likely comes from the fact that the proportion of neutral regions  
265 decreases when  $a_1$  increases, so that the IICR becomes more similar to that of a two class  
266 model with only  $\lambda_1$  and  $\lambda_3$ , directly increasing to  $\lambda_3$ .

267 Despite this complex interplay, Figure 2 provides some insights about our capacity  
268 to detect or quantify either type of selection based on the IICR. The left panel suggests  
269 that the IICR includes relevant information about the proportion of the genome under  
270 positive or background selection: for large values of  $a_1$ , there is a quick decline of the IICR  
271 (forward in time) followed by a low plateau around  $\lambda_1$ , whereas lower  $a_1$  values see a more  
272 recent and gradual decrease of the IICR without any clear recent plateau. However, this  
273 distinction is far less visible when plotting on a natural scale (Figure S2), in which case  
274  $a_1$  values as different as 0.1 and 0.5 lead to quite similar IICRs. Besides, results on the  
275 importance of  $a_1$  are likely exaggerated by the small value of  $\lambda_1$  used in Figure 2, which  
276 implies a 10-fold reduction of  $N_e$ . In comparison, our choice of  $\lambda_3$  only implies a 3-fold  
277 increase of  $N_e$  in Figure 2.

278 While the value of  $\lambda_3$  (more generally of the highest  $\lambda_i$ ) determines the plateau of the  
279 IICR, the proportion of this class ( $a_3$ ) appears to determine to a large extent the speed of  
280 convergence (backward) to this ancient plateau (right panel). For the smallest  $a_3$  values  
281 (0.1 or 0.01%), this ancient plateau is not reached within the figure (for  $t \leq 10$ ) whereas

282 a plateau corresponding to the neutral regions ( $\lambda_2 = 1$ ) is observed for quite long periods.  
283 For the largest  $a_3$  values considered here (1 or 10%), the convergence backward to the  
284 ancient plateau is so fast that the IICR does not exhibit the middle plateau around the  
285 neutral value, as already mentioned.

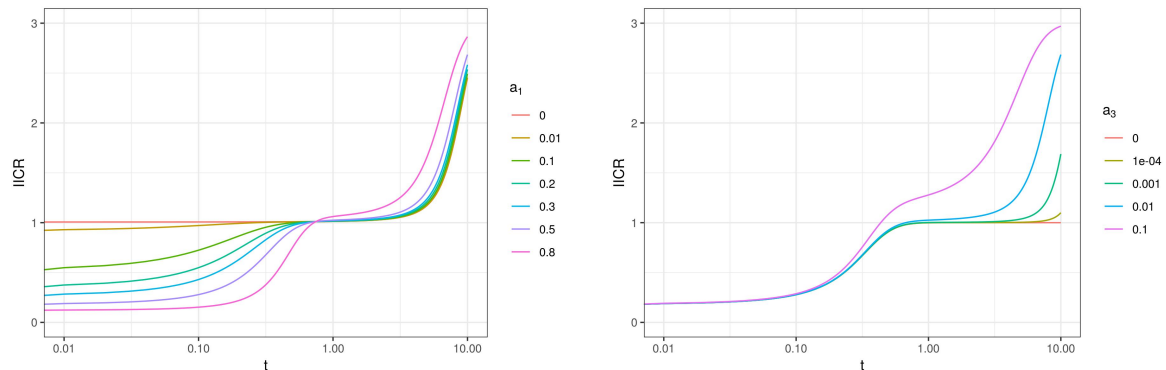


Figure 2: IICR for a panmictic model with  $K = 3$   $\lambda_i$  values such that  $\lambda_1 < 1$ ,  $\lambda_2 = 1$  and  $\lambda_3 > 1$ . The first class (or type) of genomic regions ( $\lambda_1 < 1$ ) is meant to represent regions of the genome under positive or negative selection and is modelled by a constant population size  $\lambda_1 N$  with  $\lambda_1 = 0.1$ . Genomic regions of class 2 are meant to represent neutrality and they have a constant population size  $\lambda_2 N$  where  $\lambda_2 = 1$ . Regions of class 3 are meant to represent genomic regions under balancing selection, they have a constant population size  $\lambda_3 N$  with  $\lambda_3 = 3$ . Left panel: the frequency of class 3 is fixed at  $a_3 = 0.01$  and the frequencies of classes 1 and 2 are allowed to vary. The frequency  $a_1$  is given by the legend. Right panel: the frequency of class 1 is fixed at  $a_1 = 0.5$  and the frequency of classes 2 and 3 are allowed to vary. The frequency  $a_3$  is given by the legend.

286 In any case, these results suggest that if selection can be seen as reducing or increasing  
287  $N_e$  in a panmictic population, the strongest effect on the IICR seems to be dispropor-  
288 tionately the result of the largest  $N_e$ , even though it may in practice affect ancient parts



289 of the IICR curves that may not be easily reconstructed from real data. PSMC curves  
290 obtained from real data show a sharp decrease (forward in time) in the very ancient past  
291 in several species, including humans and Neanderthals. While this ancient decrease is  
292 usually ignored or interpreted as a statistical artefact resulting from the very low number  
293 of coalescence events dating back to this period, Figure 2 suggests that it is possibly due  
294 to divergent alleles maintained by balancing selection.

### 295 **Methods: distributions of $N_e$ inferred from real data**

296 The above examples highlighted important and partly unexpected properties of the IICR  
297 when  $N_e$  is variable along the genome. However, they relied on a very small number of  
298 classes with arbitrary  $\lambda_i$  and  $a_i$  values. It is thus not clear to which extent they inform us  
299 on the impact of linked selection in real species, where the combined variations of gene  
300 density, selection form or intensity and local recombination rate generate complex  $N_e$   
301 distributions. In this section we consider two model species for which variation in  $N_e$  has  
302 been documented or estimated, the fruit fly *Drosophila melanogaster* and humans (Figure  
303 3).

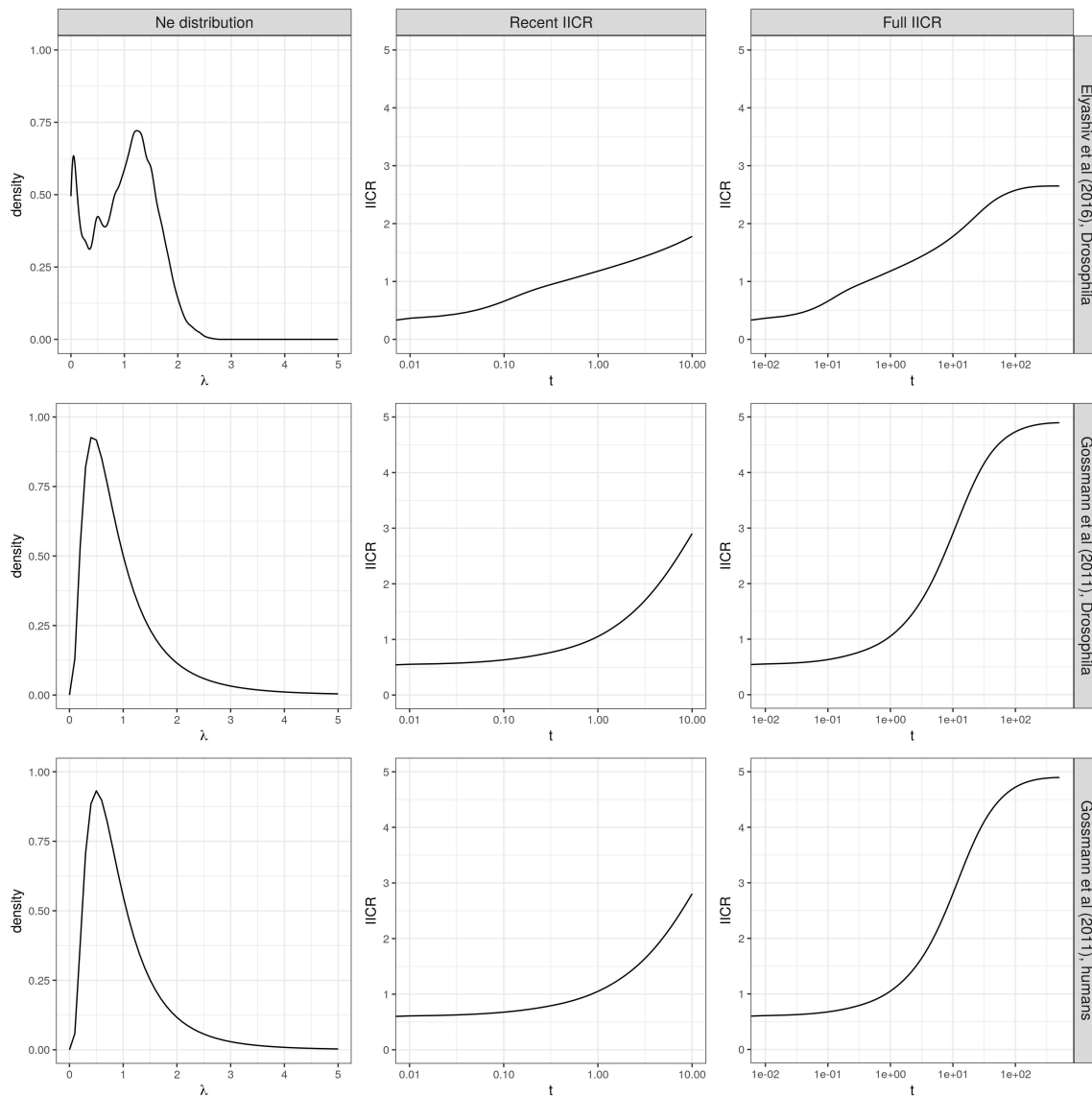


Figure 3: IICRs for panmictic models with large numbers of classes. This figure represents genome-wide distributions of  $\lambda_i$  (left panels) and the associated IICRs until  $t = 10$  (middle panels) or  $t = 500$  (right panels). Top panels: IICR for *Drosophila melanogaster* (Raleigh, North Carolina population) based on the  $N_e$  distribution estimated by Elyashiv et al. (2016). Middle panels: IICR for *D. melanogaster* (Zimbabwe population) based on the  $N_e$  distribution estimated by Gossmann et al. (2011) assuming a lognormal distribution. To make the two IICRs comparable, the distribution estimated by Elyashiv et al. (2016) (top left) was re-scaled to have an average of one, as assumed in the analysis of Gossmann et al. (2011) (middle left). Bottom panels: IICR for humans (Yoruba population) based on the  $N_e$  distribution estimated by Gossmann et al. (2011) assuming a lognormal distribution.

304 In the case of *Drosophila melanogaster*, we compared two different distributions of  $\lambda_i$   
305 over the genome, obtained by Gossmann et al. (2011) and Elyashiv et al. (2016). These  
306 two methods combine polymorphism data from the focal species and divergence data with  
307 closely related species, but they are based on very different approaches: the method of  
308 Elyashiv et al. (2016) explicitly models selection and its impact on the pairwise coalescence  
309 rate in each genomic region, while the method of Gossmann et al. (2011) assumes a log-  
310 normal distribution of  $N_e$  over the genome and estimates its scale parameter from a large  
311 number of loci. For each of these two methods, the distribution obtained for *Drosophila*  
312 *melanogaster* was converted into a discrete distribution of  $\lambda_i$  values with  $K = 25$  and  
313 the associated IICR was computed using formula (2) (see the Supplementary Material for  
314 more details). As a comparison with another species, we also considered the distribution  
315 obtained by Gossmann et al. (2011) for humans.

## 316 **Results: distributions of $N_e$ inferred from real data**

317 The distribution of  $\lambda$  inferred by Elyashiv et al. (2016) for *Drosophila* differed from the  
318 other two on two aspects (Figure 3). First, it had a lower support (up to  $\lambda_i = 2.5$ , versus  
319  $\lambda_i = 5$  for the others). This implied a smaller plateau of the IICR (as expected from  
320 equation (4)), but this effect was mainly visible at very ancient times (back to  $t = 500$ ,  
321 right column) for which the IICR is unlikely to be observed from real data. Second, it had  
322 a mode for very low  $\lambda_i$  values, which probably resulted from the inclusion of regions with  
323 very low recombination where the impact of linked selection is substantial. This mode  
324 had a limited effect on the IICR (see Figure S3 for an IICR obtained after filtering out  $\lambda$   
325 values below 0.25 from the distribution).

326 Despite the differences between the species and the methods used to estimate the  
327 variation in  $N_e$ , we obtained rather similar IICRs between  $t = 0$  and  $t = 10$  (middle

328 column). The magnitude of the decrease observed in these IICRs was also comparable  
329 to that expected from Figure 2 for small values of  $a_1$  (e.g.  $a_1 = 0.1$ , top right panel).  
330 Consequently, a long term 5 fold IICR decrease (from  $t = 10$  to  $t = 0$  forward in time)  
331 could realistically be the result, in both humans and *Drosophila melanogaster*, of a mod-  
332 erate proportion of loci with very small  $N_e$  (Figure 2,  $a_1 = 0.1$ , Figure 3, top) or from a  
333 larger proportion of loci with only slightly decreased  $N_e$  (Figure 3, middle and bottom),  
334 all as a consequence of linked selection. Obviously, this conclusion can only be seen as  
335 a first order approximation, given that neither the estimation of the  $N_e$  distribution by  
336 Elyashiv et al. (2016) or Gossmann et al. (2011), nor the computation of the resulting  
337 IICR, account for population demography or structure. Models including these aspects  
338 when computing the IICR are considered in the next section.

## 339 **Generalisation to more complex models**

### 340 **Methods: extended model**

341 We can generalise equation (2) to more complex models by still assuming that the genome  
342 is divided into  $K$  groups of loci each characterized by a different coalescence rate history.  
343 However, instead of describing this history by assuming panmixia and constant popula-  
344 tion size ( $\lambda_i N$ ), we can study different demographic models with departures from these  
345 assumptions, including models with panmixia and population size changes, models with  
346 population structure and models with transient (rather than recurrent) selection. In this  
347 more general framework, let us denote  $f_i(t)$  the *pdf* of the coalescence time  $T_2^i$  in the  $i$ -th  
348 class and  $a_i$  the proportion of the genome in this class. The IICR is:

$$\text{IICR}(t) = \frac{\sum_{i=1}^K a_i R_i(t)}{\sum_{i=1}^K a_i f_i(t)}. \quad (5)$$

349 where  $f_i(t) = -R'_i(t)$ .

## 350 **Results: panmixia and population size changes**

351 One first potential application of this general framework is to study how linked selection  
352 interferes with genuine temporal variations of the population size. For instance, a natural  
353 question would be to know whether the spurious signal of recent population size decline  
354 arising from positive or background selection is strong enough to mask a genuine recent  
355 population expansion. To answer this question, we considered a simple extension of the  
356 two-class model studied in Figure 1 ( $K = 2$ ,  $\lambda_1 = 0.1$  and  $\lambda_2 = 1$ ), where the population  
357 sizes in the two classes are multiplied by the same factor at a given time  $T$  before present.  
358 This expansion factor was set either to 5 in order to mimic the magnitude of (opposite)  
359 linked selection effects (Figure 4), or to 100 to mimic the very strong recent expansion  
360 that may be observed in some species including humans (Figure S4. The IICR of this  
361 model was computed by inserting known analytical expressions for the pdf of  $T_2^i$  in each  
362 class  $i$  (e.g. (Mazet et al., 2015)) into formula (5). Note that the same approach could  
363 be applied to arbitrary complex demographic and selective scenarios, as long as the same  
364 temporal variations are applied to all classes.

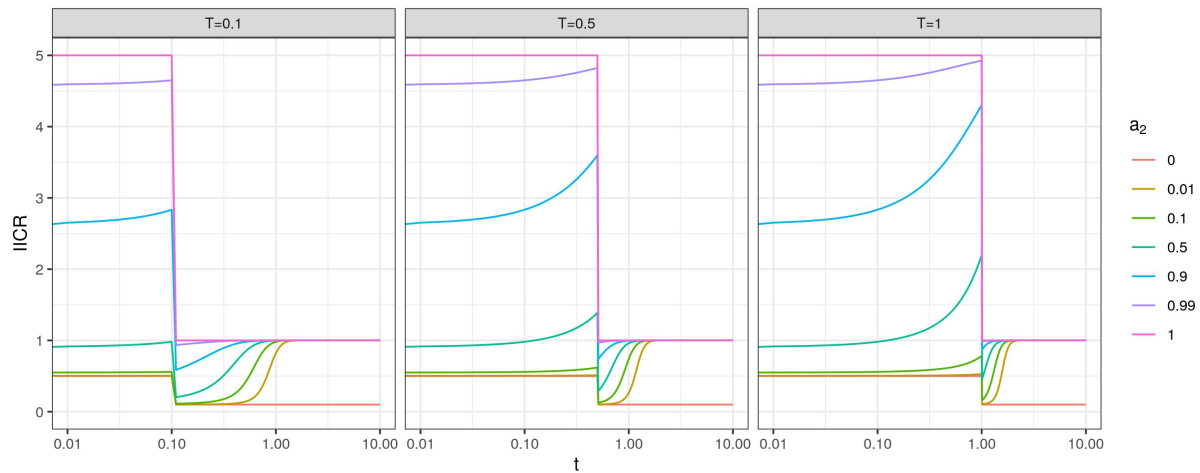


Figure 4: IICR curves for a panmictic model with a recent 5 fold expansion and  $K = 2$  classes of genomic regions. Regions of class 1 and 2 have an ancestral population size  $2N\lambda_1$  and  $2N\lambda_2$  and a recent population size  $10N\lambda_1$  and  $10N\lambda_2$ , with  $\lambda_1 = 0.1$  and  $\lambda_2 = 1$ . Each panel corresponds to a different expansion time, indicated in the panel header. Frequencies  $a_1$  and  $a_2$  of the 2 classes are given by the legend ( $a_1 + a_2 = 1$ ).

365 In the specific scenario considered here, we found that a strong proportion of selec-  
366 tion in the genome could mask a genuine 5 fold expansion or even lead to the opposite  
367 conclusion of a population size decline (Figure 4). When 50% of the genome was under  
368 selection, the IICR showed transient temporal variations around the expansion time  $T$   
369 (whose magnitude depended on  $T$ ) but could at first approximation be interpreted as  
370 a constant population size history. When 90% of the genome was under selection, the  
371 overall pattern was that of a two fold decline. In contrast, smaller proportions of selection  
372 (10% of the genome or less) did not strongly affect the signal of population expansion. For  
373 stronger expansion events (100 fold, Figure S4), the IICR showed a significant increase for  
374 all values of  $a_1$  and  $T$ , but the IICR increase was much weaker than the true population

375 size expansion: around 15 fold for  $a_1 = 0.5$  and 10 fold for  $a_1 = 0.9$ . These results confirm  
376 that linked selection can significantly bias population size change inference, even in the  
377 presence of clear genuine demographic events.

## 378 **Results: stationary population structure**

379 One other important extension of the models considered above is to account for population  
380 structure when modelling each genomic class. To illustrate this idea, we first considered  
381 a model with  $K = 2$ ,  $\lambda_1 = 0.1$  and  $\lambda_2 = 1$  as in Figure 1. Here we assumed that these  
382 two classes evolved under a n-island model with the same number of demes ( $n = 10$ ),  
383 the difference in  $N_e$  being modelled through the use of different deme sizes in the two  
384 classes ( $\lambda_1 N$  and  $\lambda_2 N$ ) We further assumed that selection did not affect migration, so  
385 that the *per* generation migration rate  $m$  was the same for the two classes. In other  
386 words, selection reducing  $N_e$  is assumed to operate after migration and thus only affects  
387 coalescence rates, but not migration rates, of the two genomic regions. This implies that  
388 the scaled migration rate  $M = 2Nm$  is identical in the two classes (time scale is still  $2N$   
389 here, but  $\lambda_i N$  now refers to deme diploid size rather than to the entire population size).  
390 One way of seeing this is by considering that there are  $2N$  haploid genomes in each deme  
391 with scaled migration rate  $2Nm$  and that selection acts on the different genomic regions  
392 by changing drift by a factor  $\lambda_i$ .

393 As already mentioned and exploited in previous studies on the IICR (Grusea et al.,  
394 2018, Mazet et al., 2016, Rodríguez et al., 2018), the distribution of coalescence times un-  
395 der a symmetrical n-island model can be derived analytically (Herbots, 1994). Extending  
396 these derivations to a model with general deme size  $\lambda_i N$ , instead of  $N$  in previous studies,  
397 we can show (see the Supplementary Material) that in this case

$$f_i(t) = p_i e^{-\alpha_i t} + \left(\frac{1}{\lambda_i} - p_i\right) e^{-\beta_i t} \quad (6)$$

398 with

$$\alpha_i = \frac{1}{2} \left( \frac{1}{\lambda_i} + n\gamma + \sqrt{\left(\frac{1}{\lambda_i} + n\gamma\right)^2 - \frac{4}{\lambda_i}\gamma} \right),$$
$$\beta_i = \frac{1}{2} \left( \frac{1}{\lambda_i} + n\gamma - \sqrt{\left(\frac{1}{\lambda_i} + n\gamma\right)^2 - \frac{4}{\lambda_i}\gamma} \right),$$
$$\gamma = \frac{M}{n-1}$$

399 and

$$p_i = \frac{\gamma - \alpha_i}{\lambda_i(\beta_i - \alpha_i)}.$$

400 Setting  $\lambda_i = 1$  for all  $i$  recovers the results of Mazet et al. (2016). The IICR of an n-island  
401 model with two classes of deme size can be obtained by computing  $f_i(t)$  with each  $\lambda_i$   
402 using Equation (6) and inserting the results into Equation (5).

403 IICR curves obtained for this two class n-island model are shown in Figure 5 for  
404 different values of the scaled migration rate. For  $M = 5$ , they are similar to those shown  
405 in Figure 1. This was expected given that an n-island model with high migration ( $M \gg 1$ )  
406 should behave in a way that is similar to a panmictic model with population size  $Nn$ ,  
407 except in the recent past where the IICR of the n-island still reflects local deme size  
408 (Mazet et al., 2016). For lower migration rates, the two extreme models with  $a_2 = 0$   
409 (red curve) or  $a_2 = 1$  (violet) show that a higher plateau of the IICR is observed as  $M$   
410 decreases, which was again expected (Mazet et al., 2016).



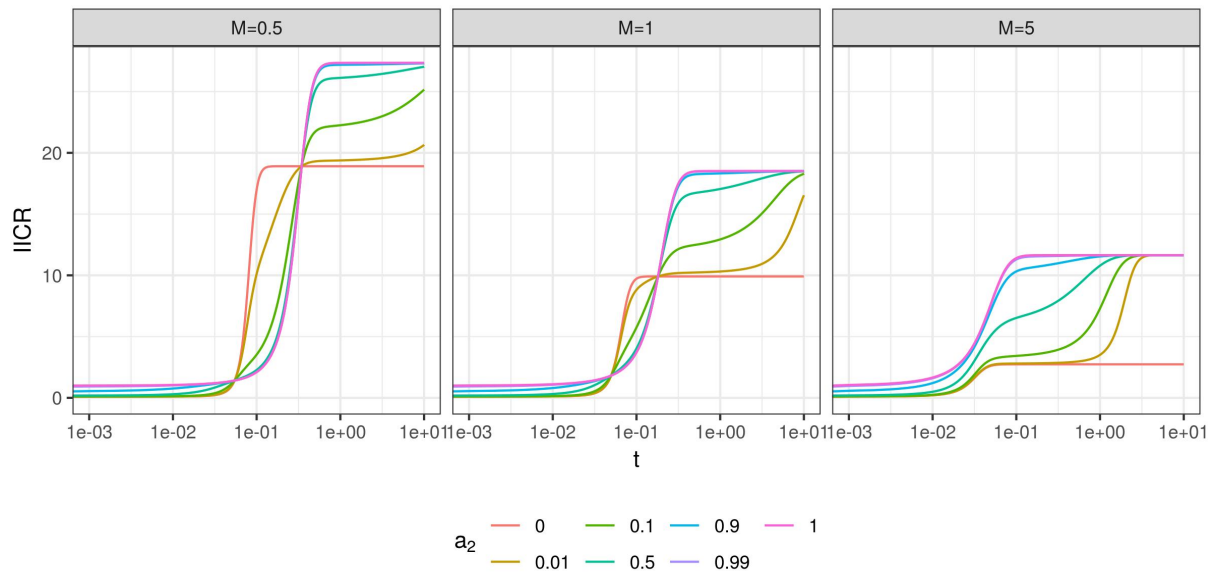


Figure 5: IICR curves for a symmetrical  $n$ -island model with  $n = 10$  demes and  $K = 2$  classes of genomic regions. Regions of class 1 and 2 have a constant deme size  $2N\lambda_1$  and  $2N\lambda_2$  with  $\lambda_1 = 0.1$  and  $\lambda_2 = 1$ . Scaled migration rate  $M = 4Nm$  is the same for the two classes, each panel corresponding to a different value of this parameter. Frequencies  $a_1$  and  $a_2$  of the 2 classes are given by the legend (having in mind that  $a_1 + a_2 = 1$ ). For comparison with panmictic models (in particular those in Figure 1), time is scaled by the meta-population size  $2Nn$  rather than by the deme size  $2N$  as in Equation (6).

411 For lower migration rates ( $M \leq 1$  in Figure 5), models with rather large values of  
 412  $a_1$  are hard to distinguish from the model with  $a_1=0$  (no selection). For instance, the  
 413 IICR with  $a_2 = a_1 = 0.5$  is not very different from that with  $a_2 = 1$ , in contrast to  
 414 Figure 1 where panmixia was assumed. This suggests that population structure may tend  
 415 to mask the effect of positive or negative selection even when a quite important part of  
 416 the genome is under selection. On the other hand, the IICR with  $a_2 = 0.01$  is more

417 similar to that with  $a_2 = 0$  than under panmixia. This suggests that, in the presence of  
418 population structure, models with pervasive selection (99% of the genome with  $\lambda = 0.1$ )  
419 may be interpreted as neutral models with small effective size (100% of the genome with  
420  $\lambda = 0.1$ ).

421 Another interesting observation from Figure 5 is the existence of a time window where  
422 the IICR is lower when  $a_2$ , corresponding to the largest  $N_e$ , is largest, i.e. the IICR  
423 is lower for models with a smaller part of their genome under selection reducing  $N_e$ .  
424 This time window occurs in the recent past and is wider for lower migration rates. This  
425 counterintuitive result illustrates the limits of interpreting the IICR as a trajectory of  
426 effective size, as already outlined for several other demographic scenarios (Chikhi et al.,  
427 2018, Mazet et al., 2016). Outside this period, the IICR curves seem to always reach  
428 higher values when  $a_2$  is larger. This is in particular the case for  $t$  close to 0, which is  
429 expected analytically (Equation (3)).

## 430 **Results: non stationary population structure**

431 To check whether these conclusions may still hold for more realistic evolutionary scenarios,  
432 we next assume that each genomic class evolves under the non stationary n-island model  
433 estimated by Arredondo et al. (2021) to fit the observed PSMC of a modern human from  
434 Karitiana (Li and Durbin, 2011). This model includes 11 islands with symmetric migration  
435 and (diploid) deme size 1,380 and it assumes that these islands go through 4 changes of  
436 connectivity in the past:  $M \approx 0.9$  ( $m \approx 1.6e-4$ ) from present to 24,437 generations before  
437 present (BP),  $M \approx 17.7$  ( $m \approx 3.2e-3$ ) from 24,437 to 82,969 generations BP,  $M \approx 2.5$   
438 ( $m \approx 4.5e-4$ ) from 82,969 to 107,338 generations BP,  $M \approx 0.7$  ( $m \approx 1.3e-4$ ) from 107,338  
439 to 179,666 generations BP and  $M \approx 1.1$  ( $m \approx 2e-4$ ) in more ancient times. We define  $K$   
440 classes of genomic regions: one neutral region with deme size  $N$  and  $K - 1$  other regions

441 under selection with deme size  $\lambda_i N$ , for  $\lambda_i$  either smaller or larger than 1. Results are  
442 shown in Figure 6, where two different options are considered to model the heterogeneity  
443 of effective size along the genome: (i) the hypothetical three class model of Figure 2 with  
444 one class corresponding to positive or negative selection and one other corresponding to  
445 balancing selection (top panels), and (ii) the 25 class model of Figure 3 estimated from  
446 Gossmann et al. (2011)’s analysis of human real data (bottom panel).

447 We find that large values of  $a_1$  could have a significant impact on the IICR in the  
448 period ranging from 10,000 to 30,000 generations ago (corresponding to 200-300,000 to  
449 600-900,000 years ago). For instance with  $a_1 = 0.8$ , the IICR is around 17 in the most  
450 recent hump and around 5 in the most recent “valley”, versus 22 and 12 without selection  
451 (top left panel). However, this effect is very moderate when considering the  $\lambda_i$  distribu-  
452 tion estimated by Gossmann et al. (2011) (bottom panel). Much more dramatic is the  
453 effect observed in the ancient past above 100,000 generations ( $\approx$  2-3 million years) before  
454 present, where the IICR with selection is significantly larger than the neutral IICR. This  
455 difference is driven by the part of the genome with large effective size (i.e. under balancing  
456 selection) and is found (with varying magnitude) in all scenarios.

457 While the neutral model considered here was estimated without accounting for se-  
458 lection and may thus be itself a biased representation of the true neutral history, the  
459 results shown in Figure 6 provide a first approximation of the impact of linked selection  
460 on demographic inference in a realistic scenario.

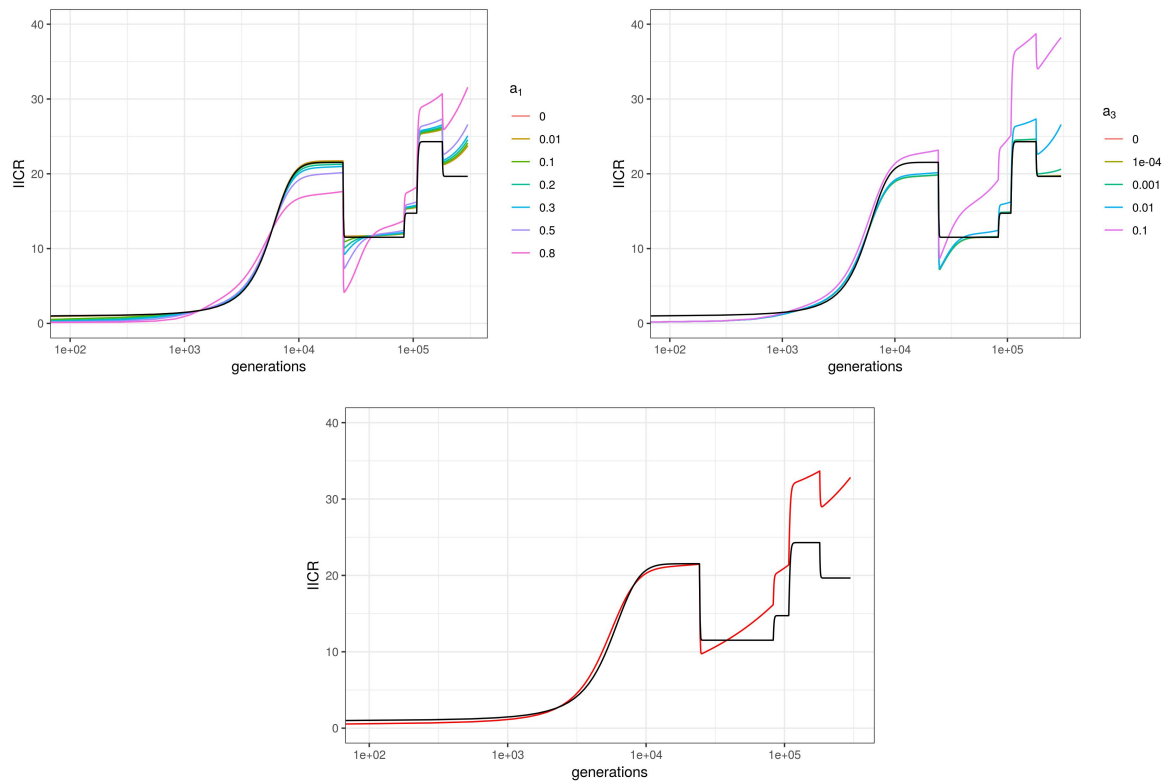


Figure 6: IICRs for demographic models combining population structure and linked selection in humans. The neutral part of the genome evolves under the non stationary n-island model estimated by Arredondo et al. (2021) to fit the observed PSMC of a modern human from Karitiana (Li and Durbin, 2011). This model includes 11 islands with (diploid) deme size  $N = 1380$ , whose connectivity varied along time according to a 3 step process (see the text for details). To account for selection, this neutral class only represents a fraction of the genome and other classes with lower or higher  $N_e$  are also considered. The number of these classes, their proportions and deme sizes (relative to the neutral class) are taken either from Figure 2 (top, where  $a_3$  is fixed to 0.01 in the left panel, and  $a_1$  fixed to 0.5 in the right one) or from Figure 3 (bottom, red line). The black curve on all panels depicts the IICR for this demographic scenario but without selection. Time is shown in generations and in log10 scale.

## 461 **Methods: modelling transient selection**

462 We finally apply this general framework to model the transient effect of recent selec-  
463 tive sweeps, rather than the effect of recurrent positive, negative or balancing selection  
464 considered until now. For this analysis we consider a panmictic population. A similar  
465 question was tackled by Schrider et al. (2016), who showed in their Figure 5 the estima-  
466 tions obtained when applying the PSMC to a 15Mb genomic region that experienced one  
467 or several recent selective sweeps. We focus here on a scenario similar to theirs, with one  
468 single selective sweep and approximate the resulting IICR using a model with different  
469 classes of  $\lambda_i$  that are time-dependent. In contrast to the model considered in Figure 4,  
470 these temporal variations differ between classes, because they depend on the distance to  
471 the selected site. Although this model is built based on the expected variations of effec-  
472 tive size (or coalescence rate) in a 15Mb region, we note that it also applies to a whole  
473 genome having experienced on average one recent selective sweep per 15 Mb region. In  
474 other words, our aim here is not to switch from the analysis of global to local IICRs, but  
475 rather to explore the local and implicitly global effects in a relatively realistic example.

476 To approximate the IICR resulting from a recent selective sweep, we assume that the  
477 effect of this sweep can be modelled by a reduction of effective population size that is  
478 limited both in time (from the emergence of the derived favorable allele to its eventual  
479 fixation in the population) and in "genomic space" (i.e. in a genomic neighborhood of  
480 this selected variant). More precisely, we consider that the region affected by the sweep  
481 on one side of the selected locus is of size

$$L = -\log(0.05) \frac{\alpha}{8Nr \log(\alpha)}$$

482 with  $N$  the diploid population size,  $r$  the per site recombination rate and  $\alpha = 2Ns$

483 the scaled selection intensity ( $s$  being the fitness advantage of homozygotes carrying the  
484 selected mutation). This quantity corresponds to the distance in base pairs (bp) from  
485 the selected site such that heterozygosity is reduced by only 5% at the end of the sweep  
486 (Walsh and Lynch, 2018, chap. 8). To capture the fact that the reduction of effective  
487 size caused by the sweep depends on the physical distance to the selected site, we further  
488 divide this affected region in 10 classes of size  $2\frac{L}{10}$  with increasing distance from the sweep,  
489 where the factor two results from the sweep extending on both sides of the selected site.

490 Modelling the selective sweep under the classical “star-like” hypothesis (Nielsen et al.,  
491 2005), we approximate (see the Supplementary Material) the average coalescence rate  
492 during the sweep as

$$\mu_{sweep} = (1 - q)^2 \frac{1}{\tau} + q^2 \frac{1}{2N}$$

493 where

$$\tau = 8N \log(\alpha) / \alpha$$

494 is the duration of the sweep (in generations) and

$$q = 1 - e^{-4drN \log(\alpha) / \alpha}$$

495 is the per lineage probability of recombination between the selected site and the genomic  
496 class. Thus, the relative effective population size in a given genomic class affected by the  
497 sweep is equal to 1 before and after the sweep and to

$$\lambda_{sweep} = \frac{1 / \mu_{sweep}}{2N}$$

498 during the  $\tau$  generations of the sweep. A neutral class with  $\lambda = 1$  at all times is also  
499 included to account for positions within the 15Mb segment but with physical distance to  
500 the selected site greater than  $L$ .

## 501 **Results: transient selection**

502 As shown in Figure 7, top panel, the resulting IICR for  $\alpha = 200$  (corresponding to  
503  $s = 0.01$  for  $N = 10,000$ ) is very close to that of a neutral scenario. The IICR for  
504  $\alpha = 1000$  (corresponding to  $s = 0.05$  for  $N = 10,000$ ) shows a reduction of about one half  
505 at sweep time, similar to the average PSMC plot in Figure 6B of Schrider et al. (2016).  
506 The IICR for  $\alpha = 10000$  (corresponding to  $s = 0.5$  for  $N = 10,000$  or to  $s = 0.05$  for  
507  $N = 100,000$ ) shows a much stronger decline, down to almost zero. However, the IICR  
508 decline in our analysis is very localized in time, while the PSMC decline in (Schrider et al.,  
509 2016) extends for a longer period. Another important difference is that the PSMC plot  
510 in the simulations of Schrider et al. (2016) not only recovers the neutral value after the  
511 sweep but increases up to more than twice this value in the recent past. To understand  
512 these differences, we simulated coalescence times along a 15Mb region under the same  
513 sweep scenario, with  $\alpha = 1000$ , using the software *msms* (Ewing and Hermisson, 2010)  
514 and estimated the resulting empirical IICR as in Chikhi et al. (2018).

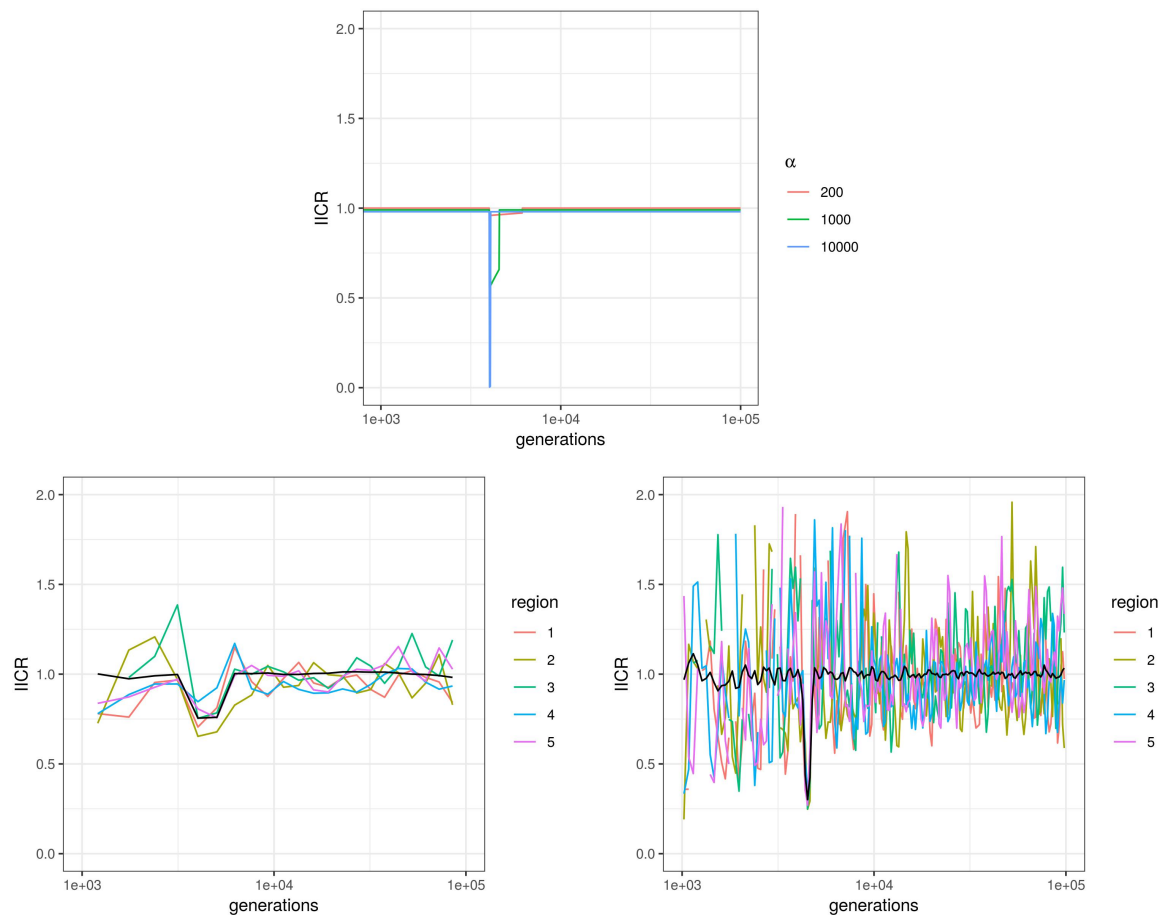


Figure 7: IICRs for a 15Mb region experiencing a single recent selective sweep. Parameter values were chosen to reproduce those in Figure 5 of Schrider et al. (2016):  $N = 10000$  (diploid size),  $r = 10^{-8}$  (per site recombination rate) and  $t_0 = 4000$  generations before present (time where the derived allele got fixed). Times are given in generations and are shown in log10 scale. Top: Expected IICRs when modelling selection using a panmictic model with  $K = 11$  classes of regions. Class 11 represents the neutral part of the region (unaffected by the sweep), with relative population size  $\lambda_{11} = 1$ . Class  $j$  ( $1 \leq j \leq 10$ ) represents a part of the region affected by the sweep, with a given physical distance from the selected site (which increases with  $j$ ). Relative population size is equal to  $\lambda_j = 1$  before and after the sweep and is decreased during the sweep to match the larger coalescence rate (see the text for more details). The proportion of each selected class  $j \geq 10$  is  $L/5$ , where  $L$  is the size of the region affected by the sweep on either side of the selected site. Scaled selection intensity  $\alpha = 2Ns$  was equal to 200, 1000 or 10000 (see the legend). Bottom: Empirical IICRs based on coalescence times simulated with the software *msms*, for  $\alpha = 1000$ . Two hundreds independent 15Mb regions were simulated. Colored lines show the IICRs for 5 of these regions (taken at random) and thus represent typical local IICRs. Black lines show the IICRs obtained when merging coalescence times from all regions, they thus correspond to genome-wide IICRs obtained for a 3Gb genome ( $200 \times 15\text{Mb}$ ) with one selective sweep every 15Mb. The number of time windows considered (i.e. of distinct estimated IICR values) was equal to 25 (left) or 200 (right) and the length of these windows was increasing exponentially backward in time, as in the PSMC approach.



515 Similar to PSMC estimations, these empirical IICR estimations depend on the number  
516 of time windows considered, the assumption being that  $N_e$  is constant within each time  
517 window but may vary between time windows. In the bottom left panel of Figure 7, we  
518 consider 25 time windows, which corresponds to the order of magnitude used in most  
519 PSMC studies. The resulting IICR, averaged over 200 replicates, is transiently reduced  
520 around the sweep time and shows no increase above 1 in the recent past, similar to our  
521 theoretical prediction (top panel). However, the reduction of  $N_e$  is both longer and of  
522 lower magnitude than in our prediction, as in the PSMC plots of Schrider et al. (2016).  
523 In the bottom right panel, we consider 200 time windows and obtain an average IICR  
524 in which the magnitude and duration of the decrease is much more consistent with our  
525 theoretical prediction. IICRs from single replicates also correctly capture this reduction  
526 around the sweep time but are very noisy outside this period as a side effect of the  
527 finer time discretization. Altogether, these results show that modelling selective sweeps  
528 by local transient changes of population size leads to a reasonable approximation of the  
529 IICR (or equivalently of the genome-wide distribution of  $T_2$ ) but that discretizing time  
530 using a limited number of time windows may lead to soften the true sweep signature by  
531 an averaging effect. They also outline that some aspects of a PSMC estimation, as the  
532 recent expansion following the sweep in the study of Schrider et al. (2016), cannot be  
533 predicted by the IICR, whatever method is used to compute the IICR. The next section  
534 explores in more details the link between IICR predictions and PSMC estimations.

## 535 **IICR predictions and PSMC estimations**

536 The models and results presented so far allow to predict the effect of linked selection on  
537 the IICR, or equivalently on the genome-wide distribution of pairwise coalescence times.

538 However, coalescence times are not directly observed from real data so the IICR is in  
539 practice estimated from methods like PSMC or MSMC. When population size history  
540 is homogeneous along the genome (i.e.  $K = 1$  class), PSMC generally provides a very  
541 good estimation of the IICR (Mazet et al., 2016) (taking apart considerations relative the  
542 amount or the quality of the data). But when population size history is heterogeneous  
543 along the genome, as considered here to approximate the effects of selection, the answer  
544 may depend on the scale (10kb? 100kb? 1Mb?) at which this heterogeneity is detectable.  
545 In other words, for a fixed proportion of genomic positions with reduced effective size due  
546 to linked selection, PSMC results may depend on the spatial clustering of these positions  
547 along the genome, while the IICR does not.

548 To explore this question, we tested whether genomic data including genome-wide het-  
549 erogeneity of  $N_e$  at different scales could generate PSMC plots consistent with our IICR  
550 predictions. To do this we carried out a limited number of additional simulations in  
551 which, using the genomic sizes  $\lambda_1 = 0.1$  and  $\lambda_2 = 1$ , we varied the lengths  $L_1$  and  $L_2$  of  
552 contiguous DNA chunks belonging to a given class, while keeping constant the propor-  
553 tions  $a_1$  and  $a_2 = 1 - a_1$  at which these classes are represented. The lengths  $L_2$  for the  
554 chunks of class 2 were chosen to be  $10^6$ ,  $10^5$  and  $10^4$  base pairs, and the lengths for the  
555 chunks of class 1 followed from the proportions  $a_1$  and  $a_2$ . We tested three values for  
556 the frequency  $a_1$  (0.5, 0.9 and 0.99), and for each combination of  $a_1$  and  $L_1$  we simulated  
557 two independent genomes of length  $10^9$  base pairs, where the two size classes were evenly  
558 spaced in the form  $(L_1, L_2, L_1, L_2, \dots, L_1, L_2)$ . We found that PSMC estimations fit well  
559 IICR predictions for large chunks ( $L_2 = 10^6$  and  $10^5$ ), but may highlight more complex  
560 and unpredicted patterns for smaller ones (Figure 8).

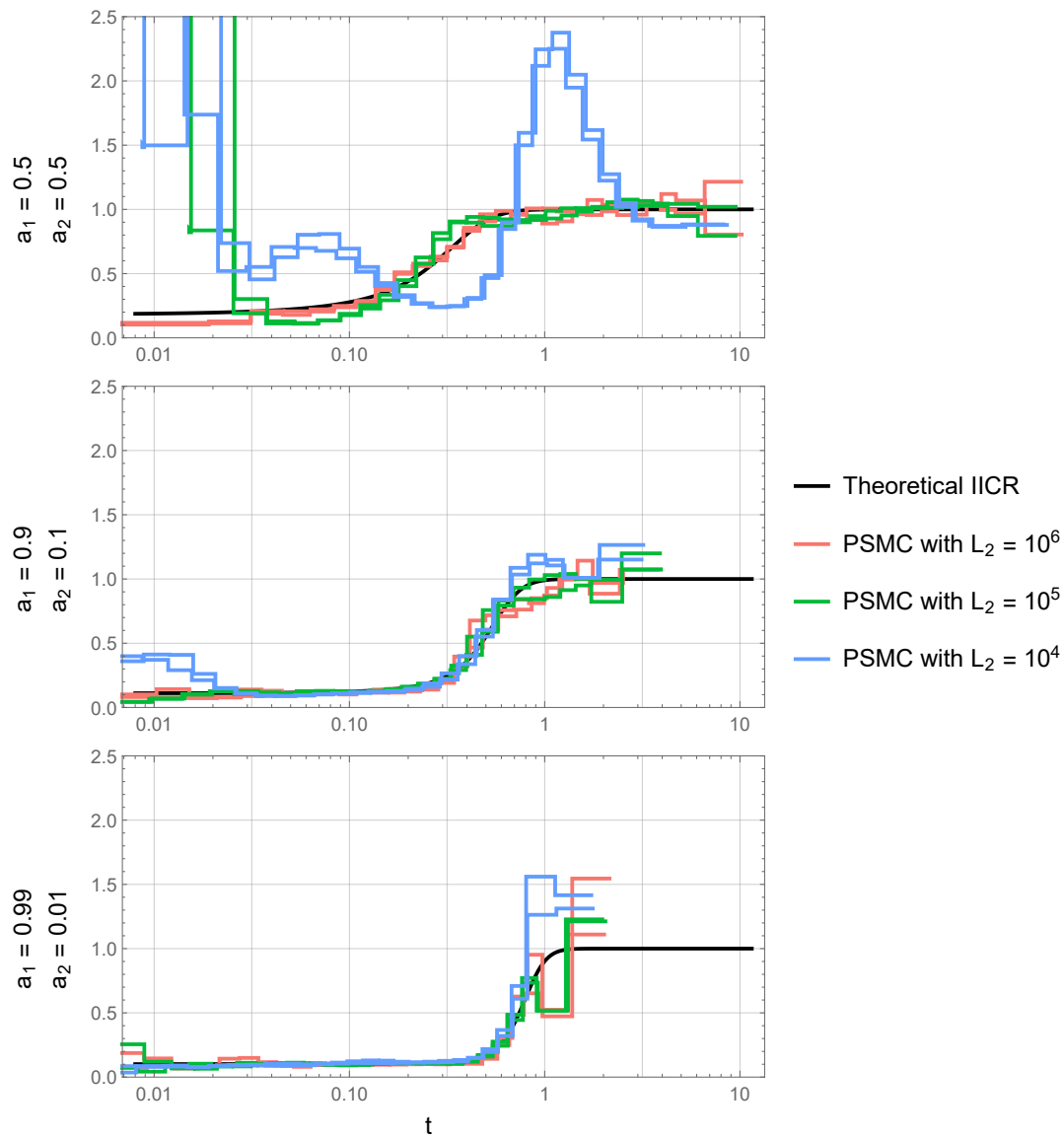


Figure 8: Comparison between theoretical IICR and inferred PSMC. For each frequency distribution  $(a_1, a_2)$  of the two size classes  $\lambda_1 = 0.1$  and  $\lambda_2 = 1$  we show the corresponding theoretical IICR (black) and two independent PSMC simulations for three values of the chunk length  $L_2$ . In each case,  $L_1 = \frac{a_1}{a_2} L_2$ . The simulated sequence has a total length of  $10^9$  bp and the two class chunks are evenly alternated in the form  $(L_1, L_2, L_1, \dots, L_2)$ . Population size was equal to 10000.

## 561 Discussion

### 562 Effects of linked selection on the IICR

563 A classical assumption in population genetics considers that linked selection can be mod-  
564 elled as a first approximation by a local change in effective population size (Hill and  
565 Robertson, 1966). Background selection and selective sweeps, which tend to reduce ge-  
566 netic diversity locally (Charlesworth et al., 1993, Smith and Haigh, 1974), are then seen  
567 as resulting in lower  $N_e$  values, whereas genomic regions under balancing selection are in  
568 contrast interpreted in terms of higher  $N_e$  values. In both cases, the impact of selection on  
569 genetic diversity or  $N_e$  is stronger for regions with lower recombination or higher selective  
570 constraints (number of selected sites, selection intensity) (Charlesworth, 2009). At the  
571 genome-wide level, linked selection appears thus to generate an apparent heterogeneity of  
572  $N_e$  among genomic regions, reflecting the variations of the mode (increasing or decreasing  
573  $N_e$ ) and the intensity of linked selection (Gossmann et al., 2011, Jiménez-Mena et al.,  
574 2016a). Following this simplifying assumption, we described in this study the distribu-  
575 tion of the coalescence time between two sequences ( $T_2$ ) for models including variable  
576 classes of  $N_e$  along the genome. More precisely, we characterized the IICR (Mazet et al.,  
577 2016) of such models, a quantity that is equivalent to the  $T_2$  distribution and corresponds  
578 to the graphical output of the popular PSMC approach (Li and Durbin, 2011), which is  
579 generally interpreted as the past temporal trajectory of  $N_e$  of the population or species  
580 under study. This analysis allowed us to predict the expected effects of linked selection  
581 on PSMC or related demographic inference approaches (Schiffels and Durbin, 2013).

582 One of the main conclusions of our work is that, under panmixia and constant popula-  
583 tion size, the existence of several classes of  $N_e$  (induced by linked selection) *always* results  
584 in a spurious signal of population size decline: the IICR of such models is a decreasing

585 function (forward in time) whose highest value (reached in the ancient past) corresponds  
586 to the largest genomic  $N_e$  and lowest value (reached in the most recent past) to the har-  
587 monic mean of genomic  $N_e$  values weighted by their relative proportion in the genome  
588 (Figure 1, Equation 3). Specifically, we found that selection reducing  $N_e$  (background  
589 selection or sweeps) has a stronger effect on the IICR in the recent past, while selection  
590 increasing  $N_e$  (balancing selection) mainly influences the IICR in the intermediate and  
591 ancient past (Figure 2). There is a striking asymmetry between the two forms of selection:  
592 because the IICR plateau is determined by the class with the largest  $N_e$  independently  
593 of the proportion of this class, even a minute proportion of balancing selection can have  
594 a large effect on the IICR, whereas higher proportions of background selection or sweeps  
595 are necessary to generate significant and detectable effects on the IICR (Figure 2). Com-  
596 bining the two forms of selection by considering  $N_e$  distributions inferred from real data  
597 (Elyashiv et al., 2016, Gossmann et al., 2011) we found that linked selection is expected  
598 to cause a long term apparent five-fold decrease of the IICR in organisms such as humans  
599 or *Drosophila melanogaster* (Figure 3). However, we stress that these results assumed  
600 panmixia and constant population size.

601 Another important conclusion of our work is indeed that the effects of linked selection  
602 on the IICR mentioned above may be largely hidden by those of population structure.  
603 Considering a symmetrical  $n$ -island model, we observed for instance that even when a  
604 large proportion of the genome is influenced by selection reducing  $N_e$  the effect on the  
605 IICR could be difficult to see for models with reduced migration rates between islands  
606 (Figure 5). Focusing on humans we also considered a simple but reasonable demographic  
607 scenario of variable population structure (Arredondo et al., 2021) together with a realistic  
608 genomic  $N_e$  distribution for this species (Gossmann et al., 2011). We found that the  
609 largest and most visible effect of linked selection on the IICR was an ancient population

610 size decline related to the presence of balancing selection (Figure 6, bottom).

611 Such ancient declines are indeed observed in PSMC plots inferred in humans and a  
612 number of other species, but a further complication is that these patterns may also arise  
613 due to the low number of informative coalescence events available to PSMC in this ancient  
614 time period. PSMC analyses of genomic data simulated under realistic demographic sce-  
615 narios, with and without balancing selection, will be necessary to investigate whether these  
616 ancient signatures of balancing selection can be disentangled from statistical artifacts. As  
617 a simple test we simulated genomic data under the demographic model of Figure 6 with  
618 a single genomic  $N_e$  (i.e. no selection). We applied PSMC to these data and found no  
619 ancient decrease in the estimated trajectory compared to the expected IICR (Figure S5).  
620 These admittedly limited results suggest that the PSMC is not necessarily *statistically*  
621 biased in the ancient past, and that the signals observed in several species including hu-  
622 mans and chimpanzees might be due to balancing selection or other forms of selection  
623 maintaining high levels of diversity over very long periods. One possible strategy to limit  
624 the influence of regions submitted to such forms of selection would be to first detect them  
625 and filter them out from the PSMC analysis. For the demographic scenario of Figure 6,  
626 we found that this would reduce the biases observed in the ancient past without affecting  
627 significantly other parts of the IICR (Figure S6).

## 628 **The intriguing signature of background selection on the IICR**

629 The framework developed in this study makes no particular distinction between posi-  
630 tive and background selection, which are both modelled as leading to a reduction of  $N_e$ .  
631 Thus, one possible interpretation of our results would be that ignoring background selec-  
632 tion leads to infer spurious population declines. This conclusion is at odds with several  
633 previous studies, which concluded that unaccounted background selection may actually

634 lead to a spurious signature of recent population expansion. For instance, Zeng and  
635 Charlesworth (2011) and Walczak et al. (2012) developed theoretical approximations of  
636 the genealogical process at a neutral locus linked to a site under negative selection and  
637 showed that this process shared many properties with that of an expanding population.  
638 The former study accounted for intra-locus recombination, whereas the latter ignored  
639 it. Several recent studies have applied demographic inference methods to genomic data  
640 simulated with and without background selection (Ewing and Jensen, 2016, Johri et al.,  
641 2021, Lapierre et al., 2016, Pouyet et al., 2018) and observed a signal of recent popula-  
642 tion expansion in the scenarios including selection. Finally, Johri et al. (2020) analyzed  
643 real data from an African population of *Drosophila melanogaster* with a new ABC demo-  
644 graphic inference approach accounting for background selection. They estimated that the  
645 size of this population has been relatively constant for a few millions generations, while  
646 several previous studies on this or other related populations, which ignored background  
647 selection, estimated a strong recent population size increase, e.g. (Arguello et al., 2019,  
648 Kapopoulou et al., 2018).

649 Two main reasons may resolve this apparent paradox between these previous results  
650 and ours. First, we assume that linked selection can be modelled by a local change of  
651  $N_e$  without any temporal dynamics (except in Figure 7 and related text, whose focus is  
652 specifically on recent selective sweeps). In particular, our results do not hold for demo-  
653 graphic inference approaches based on the Site Frequency Spectrum (SFS), because weak  
654 background selection is expected to produce an excess of low frequency alleles, in partic-  
655 ular singletons, which cannot be mimicked by just assuming a smaller  $N_e$ . Such an excess  
656 of rare alleles is also a classical signature of expanding populations, which may explain  
657 the conclusions of several of the studies mentioned above (Ewing and Jensen, 2016, Johri  
658 et al., 2020, Lapierre et al., 2016, Pouyet et al., 2018).

659 Second, even when focusing on pairwise statistics such as heterozygosity or  $T_2$ , the  
660 signature of population decline predicted by the IICR can only be observed if the data  
661 considered exhibit some heterogeneity in  $N_e$ . As it can easily be seen from Figure 1,  
662 panmictic models with either no ( $a_2 = 1$ ) or only ( $a_2 = 0$ ) selection do not show declining  
663 but constant IICRs. Consequently, a decline signature is not necessarily expected when  
664 analyzing a single locus under selection as in Zeng and Charlesworth (2011) or Walczak  
665 et al. (2012). It is also not necessarily expected when analyzing genome-wide data with  
666 homogeneous selective constraints along the genome. For instance, Johri et al. (2021)  
667 simulated genome-wide sequences including background selection by considering a regular  
668 alternance of functional (selected) and intergenic (neutral) regions of fixed and relatively  
669 small sizes: depending on the scenario, the size of a single 'unit' including one functional  
670 and one intergenic region ranged from  $\approx 13$  to 55 kb. The PSMC analyses of these  
671 sequences suggested a population under constant size or slight recent expansion. We  
672 believe that some of the results obtained by these (and possibly other) authors could  
673 be due to the fact that the data simulated with this approach do not exhibit enough  
674 heterogeneity in population sizes among (short) sliding windows over the genome. Such a  
675 regularity is at odds with observations made in different organisms (Elyashiv et al., 2016,  
676 Gossmann et al., 2011).

## 677 **IICR predictions and PSMC estimations**

678 Understanding the difference between our results and those of Johri et al. (2021) also  
679 leads to the fundamental question of the link between a PSMC curve and the IICR. The  
680 results obtained in Figure 8 suggest that the IICRs computed in this study are good  
681 predictors of PSMC outputs when variations of  $N_e$  occur at a relatively large scale (100  
682 kb or more), but not always when these variations occur at a smaller scale. This may



683 explain the discrepancy between our predictions and the PSMC results in the scenario  
684 simulated by Johri et al. (2021), where the heterogeneity of  $N_e$  was detectable only at  
685 very small scale ( $\leq 55\text{kb}$ ).

686 The recent selective sweep scenario considered in Figure 7 provides another example  
687 of potential differences between PSMC estimations and IICR predictions in the case of  
688 genomic heterogeneity. Simulating *genome sequences* in a single 15Mb region experiencing  
689 one recent selective sweep, Schrider et al. (2016) found that PSMC applied to these  
690 sequences would infer a bottleneck around the time of the sweep completion, generally  
691 followed by a more recent expansion exceeding the 'neutral' effective size. Simulating  
692 *coalescence times* under the same selective sweep scenario and estimating the IICR from  
693 these simulated values, we observed a similar bottleneck but no recent expansion. This  
694 difference likely results from the fact that short coalescence times are mostly clustered  
695 around the selected site in the real data, while for IICR estimation only their proportion  
696 over the 15Mb region matters. Approximating the IICR under a selective sweep through  
697 a model with several classes of time-dependent  $N_e$ , we managed to reproduce the main  
698 characteristics of the IICR of this scenario, but this is not exactly similar to the PSMC  
699 that would be estimated in this scenario.

700 Overall, these results suggest that assessing potential PSMC biases in a given species  
701 may require specific simulations based on precise genomic annotations (positions and  
702 lengths of genes, local recombination rates ...). As an alternative to such specific studies,  
703 we provide here a quick and flexible approach to predict the distribution of coalescence  
704 times in the presence of linked selection, which is to some extent also representative of  
705 expected PSMC outputs.

## 706 **Perspectives for demographic inference**

707 The above discussion illustrates that the effects of linked selection on demographic infer-  
708 ence are complex, as they not only depend on the type and intensity of linked selection  
709 but also on the inference approach applied (SFS or  $T_2$  based for instance) or the scale  
710 at which selection constraints vary along the genome. If the future confirms that linked  
711 selection is pervasive in the genome as claimed for several model species (Elyashiv et al.,  
712 2016, Pouyet et al., 2018) new demographic inference approaches accounting for linked  
713 selection and population structure will be needed. One way of achieving this objective is  
714 to jointly estimate demographic and selection parameters, as proposed in two recent stud-  
715 ies relying on simulation based approaches, deep learning (Sheehan and Song, 2016) and  
716 Approximate Bayesian Computation (ABC) (Johri et al., 2020). These studies focused  
717 on relatively simple models, considering panmictic populations with a single population  
718 size change and only some types of selection (background selection in one study, sweeps  
719 and balancing selection in the other). To integrate more complex demographic scenarios,  
720 several recent studies considered demographic models including two classes of  $N_e$  along  
721 the genome, one for neutral loci and one for loci under linked selection. The proportion  
722 of the two classes and the ratio of  $N_e$  between them were estimated together with other  
723 parameters of the demographic model, using either ABC (Rougemont and Bernatchez,  
724 2018, Roux et al., 2016) or a modification (Rougemont et al., 2020, Rougeux et al., 2017)  
725 of the diffusion approach implemented in the software  $\partial a \partial i$  (Gutenkunst et al., 2009). Our  
726 study suggests that a similar inference approach, accounting for linked selection through  
727 variable classes of  $N_e$  along the genome, could be developed based on the IICR. An IICR-  
728 based inference framework was recently proposed for the estimation of non stationary  
729  $n$ -island models and provided very encouraging results (Arredondo et al., 2021). Given  
730 the strong impact of linked selection on the IICR under panmixia, we believe that a simi-

731 lar approach could allow to jointly infer parameters related to demographic history and to  
732 the  $N_e$  distribution. However, the results obtained under models of population structure  
733 suggest that it may be necessary to use the IICR in addition to other summaries of ge-  
734 nomic diversity to overcome identifiability issues. Also, we should stress that separating  
735 the effects of population size change, selection and population structure is likely to be one  
736 of the major challenges of population genetics in the future.

### 737 **Pros and cons of an IICR approach**

738 Whether the objective is to predict potential effects of linked selection or to estimate linked  
739 selection parameters from real data, two nice features of an IICR-based approach such as  
740 the one considered here are flexibility and speed of computation. This approach allows  
741 to simultaneously include different forms of selection and to combine linked selection  
742 with arbitrary complex demographic models. The examples considered here included for  
743 instance panmictic models with temporal variations of the population size (Figure 4)  
744 and n-island models with temporal variations of the migration rate (Figure 6). We also  
745 considered different distributions of  $\lambda_i$ , some of them including a large number of classes.  
746 More general models could be considered, for instance including other forms of structure  
747 or combining population structure and temporal population size variations. In the case  
748 of structured models, variable migration rates along the genome may be considered: we  
749 could either decrease  $M$  in the linked selection class(es) to account for possible effects of  
750 selection on migration success or introduce new classes with lower  $M$  values in order to  
751 model possible barriers to gene flow (Roux et al., 2016). As outlined in Figure 7, transient  
752 selection can be modelled by including population size changes in a subset of classes, and  
753 this approach could also be extended to model more complex fluctuating selection effects.  
754 Whatever the complexity of the demographic model and the  $N_e$  distribution considered,

755 the associated IICR can be computed exactly in a very small time using the rate matrix  
756 approach described in Rodríguez et al. (2018) or Arredondo et al. (2021), which allows to  
757 efficiently explore a very large number of scenarios or parameter values.

758 We should also stress that apparent variations of  $N_e$  along the genome may result  
759 from other biological processes than linked selection. The models presented here, and  
760 the general conclusion that heterogeneity in  $N_e$  is expected to generate population size  
761 decline patterns, also apply to these other biological processes. For instance, genome-  
762 wide variations of the mutation rate may have similar effects on the data than genome-  
763 wide variations of  $N_e$ , because high mutation rates and large population sizes both lead  
764 to increase the number of polymorphic sites in a region. Consistent with our results,  
765 Sellinger et al. (2021) showed that applying SMC methods to genomic sequences that were  
766 simulated with local variations of the mutation rate leads to infer spurious population size  
767 declines. Actually, a direct consequence of  $N_e$  heterogeneity is to increase the variance of  
768 coalescence times along the genome (see the Supplementary Materials for a proof of this  
769 statement under panmixia). Inference methods like the PSMC, which do not account for  
770 *genomic* variations of  $N_e$ , try to explain this additional variance using *temporal* variations  
771 of  $N_e$ , more precisely population size declines.

772 The main limitation of the IICR approach described in this study is that it focuses on  
773 pairs of sequences. It provides information that is complementary to that provided by the  
774 SFS, as we have noted elsewhere (Arredondo et al., 2021, Chikhi et al., 2018) For instance,  
775 some effects of weak background selection or selective sweeps may be visible on the SFS  
776 but not on the IICR. Currently we have mainly focused on the IICR as defined for a pair  
777 of sequences, but extensions to multiple sequences might provide additional information  
778 on the distribution of higher order coalescence times ( $T_3, T_4, \dots$ ), hence allowing a finer  
779 characterization of selective and neutral processes.

## 780 **Closing comments**

781 We have used the IICR as a way to explore important ideas that are central to population  
782 genetics such as the notion of effective size (see also Chikhi et al. (2018), Mazet et al.  
783 (2016) for discussions on these questions), drift and selection. We wished to re-open  
784 discussions regarding the influence of selective and neutral processes on genetic diversity,  
785 some of them general and theoretical, others more specific and practical: Can selection be  
786 modelled as a genomic variation in  $N_e$ ? What are the limits of such an approximation?  
787 Can linked selection, and more generally  $N_e$  variation along the genome, be detected in real  
788 genomes by applying the PSMC method of (Li and Durbin, 2011) or related approaches?  
789 These are exciting questions to ask and the recent years have shown that they are at the  
790 heart of modern population genetics.

## 791 **Data availability statement**

792 Code used to generate the exact and simulated IICRs shown in this study can be found  
793 at [https://github.com/sboitard/IICR\\_selection](https://github.com/sboitard/IICR_selection).

## 794 **Acknowledgements**

795 Armando Arredondo was funded by the Université Fédérale Toulouse Midi Pyrénées  
796 (UFTMiP) and the Région Occitanie (formerly Midi-Pyrénées) with PhD grant No.  
797 31I2017M248. Lounès Chikhi was funded by Fundação para a Ciência e Tecnologia (ref.  
798 PTDC-BIA-EVL/30815/2017). Olivier Mazet and Lounès Chikhi were funded by the  
799 2015–2016 BiodivERsA COFUND call for research proposals, with the national funders  
800 ANR (ANR-16-EBI3-0014) and the Fundação para a Ciência e Tecnologia ref. Bio-

801 diversa/0003/2015 and PT-DLR (01LC1617A). This work was also supported by the  
802 LABEX entitled TULIP (ANR-10-LABX-41 and ANR-11-IDEX-0002-02) as well as the  
803 LIA BEEG-B (Laboratoire International Associé-Bioinformatics, Ecology, Evolution, Ge-  
804 nomics and Behaviour). We acknowledge an Investissement d’Avenir grant of the Agence  
805 Nationale de la Recherche (CEBA: ANR-10-LABX-25-01).

## 806 **Supplementary Material**

### 807 **Monotony of the IICR in a panmictic model with several classes** 808 **of constant $N_e$**

809 We consider here the first model introduced in this study, where a proportion  $a_i$  of  
810 the genome evolves under a Wright-Fisher model with constant population size  $\lambda_i N$   
811 ( $i=1, \dots, K$ ). The IICR under this model is given by equation (2). To characterize the  
812 dynamics of the IICR over time, we study the derivative of the IICR as a function of time  
813 (backward from present):

$$\text{IICR}'(t) = \frac{R(t)R''(t) - R'(t)^2}{R'(t)^2}$$

814 which has the sign of

$$\begin{aligned}
 R(t)R''(t) - R'(t)^2 &= \sum_{i=1}^K a_i e^{-\mu_i t} \sum_{j=1}^K a_j \mu_j^2 e^{-\mu_j t} - \sum_{i=1}^K a_i \mu_i e^{-\mu_i t} \sum_{j=1}^K a_j \mu_j e^{-\mu_j t} \\
 &= \sum_{i=1}^K \sum_{j \neq i} a_i e^{-\mu_i t} a_j e^{-\mu_j t} \mu_j^2 - \sum_{i=1}^K \sum_{j \neq i} a_i e^{-\mu_i t} a_j e^{-\mu_j t} \mu_i \mu_j \\
 &= \sum_{i=1}^K \sum_{j > i} a_i e^{-\mu_i t} a_j e^{-\mu_j t} (\mu_i^2 + \mu_j^2 - \mu_i \mu_j - \mu_j \mu_i) \\
 &= \sum_{i=1}^K \sum_{j > i} a_i e^{-\mu_i t} a_j e^{-\mu_j t} (\mu_i - \mu_j)^2
 \end{aligned}$$

815 This quantity is always positive so we can conclude that the IICR is *always increasing*  
 816 from  $t = 0$  to  $t = +\infty$  (i.e. backward in time).

817 **Variance of  $T_2$  in a panmictic model with several classes of con-**  
 818 **stant  $N_e$**

819 We consider here the same model as in previous section. For a given position in the  
 820 genome, let us denote  $T_2$  the pairwise coalescence time (in  $2N$  units) and  $X$  the genomic  
 821 class.  $X$  is a stochastic variable that is equal to  $i$  with probability  $a_i$ , and the distribution  
 822 of  $T_2$  conditional on  $X = i$  is an exponential distribution with parameter  $\mu_i = \frac{1}{\lambda_i}$ . In  
 823 particular, we have  $\mathbb{E}[T_2^i | X = i] = \lambda_i$  and  $Var(T_2^i | X = i) = \lambda_i^2$ . From these

824 assumptions, we can deduce that

$$\begin{aligned}\mathbb{E}[T_2] &= \mathbb{E}[\mathbb{E}[T_2 | X]] \\ &= \sum_i a_i \mathbb{E}[T_2 | X = i] \\ &= \sum_i a_i \lambda_i\end{aligned}$$

825 and

$$\begin{aligned}\text{Var}(T_2) &= \text{Var}(\mathbb{E}[T_2 | X]) + \mathbb{E}[\text{Var}(T_2 | X)] \\ &= \left( \sum_i a_i \lambda_i^2 - \left( \sum_i a_i \lambda_i \right)^2 \right) + \sum_i a_i \lambda_i^2 \\ &= 2 \sum_i a_i \lambda_i^2 - \left( \sum_i a_i \lambda_i \right)^2\end{aligned}$$

826 where the derivation from the first to the second line follows from the fact that (i)  $\mathbb{E}[T_2 | X]$   
827 is a stochastic variable equal to  $\lambda_i$  with probability  $a_i$  and (ii)  $\text{Var}(T_2 | X)$  is a stochastic  
828 variable equal to  $\lambda_i^2$  with probability  $a_i$ .

In comparison, the variance of  $T_2$  in a model with a single class of  $N_e$  and the same expected value of  $T_2$  is

$$\text{Var}(T_2^{\text{const}}) = \left( \sum_i a_i \lambda_i \right)^2$$

829 Thus, we have

$$\begin{aligned}\text{Var}(T_2) \geq \text{Var}(T_2^{\text{const}}) &\iff 2 \sum_i a_i \lambda_i^2 - \left( \sum_i a_i \lambda_i \right)^2 \geq \left( \sum_i a_i \lambda_i \right)^2 \\ &\iff \sum_i a_i \lambda_i^2 \geq \left( \sum_i a_i \lambda_i \right)^2 \\ &\iff \left( \sum_i a_i \right) \left( \sum_i a_i \lambda_i^2 \right) \geq \left( \sum_i \sqrt{a_i} \sqrt{a_i} \lambda_i \right)^2\end{aligned}$$



830 which is always true from the Cauchy Schwartz inequality.

831 Let us denote  $R = \frac{Var(T_2)}{Var(T_2^{const})}$  the ratio of the two variances, which is thus always  
832 larger than 1. We observed that this ratio generally increased with the proportion of the  
833 genome associated to the smallest  $\lambda_i$ . For instance, in the two class model of Figure 1  
834 with  $\lambda_1 = 0.1$  and  $\lambda_2 = 1$ ,  $R$  was equal to 1.08 for  $a_1 = 0.1$ , 1.53 for  $a_1 = 0.5$  and 2.24 for  
835  $a_1 = 0.1$ . In the three class model of Figure 2 with  $\lambda_1 = 0.1$ ,  $\lambda_2 = 1$ ,  $\lambda_3 = 3$  and  $a_3 = 0.01$   
836 (left panel),  $R$  was equal to 1.13 for  $a_1 = 0.1$ , 1.61 for  $a_1 = 0.5$  and 2.75 for  $a_1 = 0.1$ .

### 837 **Estimation of the distribution of $N_e$ in drosophila and humans**

838 Two different distributions of  $\lambda_i$  over the genome were obtained for *Drosophila melanogaster*.  
839 The first one was taken from the study of Elyashiv et al. (2016), who developed a method  
840 for inferring the distribution of fitness effects in different classes of functional annota-  
841 tions (UTRs, codons ...) for both beneficial and deleterious mutations. This method  
842 requires polymorphism data from the focal species, divergence data with closely related  
843 species and precise recombination and annotation maps allowing to assess the selection  
844 constraints acting on each position in the genome. A by-product of their analysis is  
845 that an estimation of  $N_e$  can be obtained for sliding windows along the genome. Inter-  
846 estingly, these  $N_e$  values resulting from the strength of linked selection in each genomic  
847 region are defined as the inverse of the coalescence rate between two sequences and all  
848 computations rely on heterozygosity values observed between pairs of individuals. This  
849 suggests that the  $N_e$  estimates should be directly comparable with our  $\lambda_i$  values, which  
850 also correspond to the inverse of pairwise coalescence rates. The values of  $N_e$  estimated  
851 by Elyashiv et al. (2016) for 1Mb sliding windows in *Drosophila melanogaster*, based on  
852 162 inbred lines derived from the Raleigh, North Carolina population, were downloaded at  
853 <https://github.com/sellalab/LinkedSelectionMaps>. Their distribution (top left panel) was

854 converted into a discrete distribution of  $\lambda_i$  values with  $K = 25$  classes using the *hist()*  
855 function of R. The IICR resulting from this distribution is shown in the top middle and  
856 right panels.

857 The second distribution used for this species was that estimated by Gossmann et al.  
858 (2011) for a Zimbabwe population. While these authors also used polymorphism and  
859 divergence data, they focused on exons and did not aim at modelling the distribution  
860 of fitness effects. They assumed a log-normal distribution of  $N_e$  with mean value of 1  
861 and estimated the scale parameter of this distribution from the observed data at several  
862 independent genes in the genome. Using the parameter obtained by this approach for  
863 *Drosophila melanogaster* and no recombination within genes (Table 1 of their study), we  
864 randomly sampled 100,000 values of  $N_e$  (or  $\lambda$ ) under the log-normal distribution (middle  
865 left panel). A discrete distribution of the  $\lambda_i$ 's and the associated IICR were then computed  
866 as explained above, filtering out large  $\lambda$  values (we arbitrarily excluded values above  
867 five). Indeed, it is not clear whether such large values would be realistic or statistical  
868 artifacts resulting from the use of a continuous distribution estimated mainly from smaller  
869  $\lambda$  values. Also, they represent less than 0.6% of the distribution. As a comparison with  
870 another species, we also applied this second approach with the scale parameter inferred by  
871 Gossmann et al. (2011) for humans based on data from the Yoruba population (bottom  
872 panels).

## 873 **Derivation of the pdf of $T_2$ in a $n$ -island model**

874 We derive here the pdf density of  $T_2$ , the coalescence time of two lineages sampled in the  
875 same deme (resp. different deme), in an  $n$ -island model. We follow the identity by descent  
876 approach used in Durrett's process (Durrett, 2008, p. 150). The size of each deme is  $\lambda N$ ,  
877 the probability of each lineage to migrate from a deme to another each generation is  $m$ ,

878 and the per locus mutation rate is  $u$ . Define the rescaled mutation and migration rates  
879 by  $\theta = 4Nu$  and  $M = 4Nm$ . Note that two lineages coalesce at rate  $c = \frac{1}{\lambda}$  when they  
880 are in the same deme, migrate at rate  $2m \cdot 2N = M$  and experience mutations at rate  
881  $2u \cdot 2N = \theta$ .

882 Let  $p_s(\theta)$  and  $p_d(\theta)$  be the probabilities that two lineages are identical by descent  
883 when they are chosen in the same or different demes. Following back two lineages from  
884 the same deme, three different events can occur: a coalescence with probability  $\frac{c}{c+\theta+M}$ ,  
885 a migration with probability  $\frac{M}{c+\theta+M}$  and a mutation with probability  $\frac{\theta}{c+\theta+M}$ . If lineages  
886 are in different demes, the only possible events are mutation, with probability  $\frac{\theta}{\theta+M}$  and  
887 migration. In this second case lineages arrive in the same deme with probability  $\frac{1}{n-1}$  and  
888 stay in different ones with probability  $\frac{n-2}{n-1}$ . Hence we have the two coupled equations:

$$p_s(\theta) = \frac{c}{c+M+\theta} \cdot 1 + \frac{M}{c+M+\theta} \cdot p_d(\theta),$$

889 and

$$p_d(\theta) = \frac{M/(n-1)}{M+\theta} \cdot p_s(\theta) + \frac{M(n-2)/(n-1)}{M+\theta} \cdot p_d(\theta).$$

890 The second equation gives

$$\begin{aligned} \left(1 - \frac{M(n-2)}{(n-1)(M+\theta)}\right) p_d(\theta) &= \frac{M}{(n-1)(M+\theta)} p_s(\theta) \\ \Leftrightarrow \frac{\theta(n-1)+M}{(n-1)(M+\theta)} p_d(\theta) &= \frac{M}{(n-1)(M+\theta)} p_s(\theta) \\ \Leftrightarrow p_d(\theta) &= \frac{M}{\theta(n-1)+M} p_s(\theta). \end{aligned}$$

891 We then inject in the first equation:

$$p_s(\theta) = \frac{c}{c + M + \theta} + \frac{M}{c + M + \theta} \frac{M}{\theta(n-1) + M} p_s(\theta)$$

892 hence

$$p_s(\theta) \left( 1 - \frac{M^2}{(c + M + \theta)(\theta(n-1) + M)} \right) = \frac{c}{c + M + \theta}$$

893 and since

$$(c + M + \theta)(\theta(n-1) + M) - M^2 = \theta^2(n-1) + \theta(c(n-1) + Mn) + cM,$$

894 we get

$$p_s(\theta) = \frac{c(\theta(n-1) + M)}{\theta^2(n-1) + \theta(c(n-1) + Mn) + cM} = \frac{c(\theta + \gamma)}{\theta^2 + \theta(c + n\gamma) + c\gamma}$$

895 and

$$p_d(\theta) = \frac{cM}{\theta^2(n-1) + \theta(c(n-1) + Mn) + cM} = \frac{c\gamma}{\theta^2 + \theta(c + n\gamma) + c\gamma}$$

896 with

$$\gamma = \frac{M}{n-1}.$$

897 Let's now note that the probability  $p_s(\theta)$  that two lineages has reached their common  
898 ancestor without undergoing any mutation is also the expected value  $\mathbb{E}(e^{\theta T_2})$ . In other  
899 words,  $p_s$  is the Laplace transform of  $T_2$ . It can be inverted by looking for the roots of  
900  $\theta^2 + \theta(c + n\gamma) + c\gamma$ . Let  $\Delta = (c + n\gamma)^2 - 4c\gamma$ , then

$$p_s(\theta) = \frac{c(\theta + \gamma)}{(\theta + \alpha)(\theta + \beta)} = \frac{a}{\theta + \alpha} + \frac{b}{\theta + \beta}$$

901 with

$$\alpha = \frac{1}{2} \left( c + n\gamma + \sqrt{\Delta} \right),$$

$$\beta = \frac{1}{2} \left( c + n\gamma - \sqrt{\Delta} \right),$$

$$a = \frac{c(\gamma - \alpha)}{\beta - \alpha}$$

902 and

$$b = \frac{c(\gamma - \beta)}{\alpha - \beta} = c - a.$$

903 Hence the probability density function of  $T_2$  is:

$$f_{T_2}(t) = ae^{-\alpha t} + (c - a)e^{-\beta t}.$$

904 Note that  $-\alpha$  and  $-\beta$  are the non zero eigenvalues of the Q-matrix,  $-\beta$  being the  
905 closest to 0, and we have the relationships  $\alpha + \beta = c + n\gamma$  and  $\alpha\beta = c\gamma$ . Note also that we  
906 could similarly obtain the pdf distribution of the coalescence time of two lineages sampled  
907 in different demes, as  $p_d$  is its Laplace transform as well.

## 908 **Approximation of the coalescence rate in a selective sweep sce-** 909 **nario**

910 Assuming a selective sweep scenario with scaled selection intensity  $\alpha$ , we consider here  
911 the genealogy at a neutral locus located  $d$  bp away from the selected site. This process  
912 can be modelled using a structured coalescent where lineages are either in the 'derived' or  
913 'ancestral' background, depending on which allele at the selected locus they are associated  
914 with (to avoid any confusion, we remind here that this structure is a modelling facility and  
915 has nothing to do with the island structure considered in some sections of the main text).  
916 In this framework, ancestral recombination events creating or breaking the association  
917 with the derived allele can be seen as migration events from one background to the other  
918 (Kaplan et al., 1988). In the case of a complete selective sweep, lineages sampled at  
919 present all belong to the derived background, because the derived allele is then fixed in  
920 the population. Following previous studies on this topic, e.g. (Nielsen et al., 2005), we  
921 further assume a "star-like" model where these lineages can either (i) escape this derived  
922 background through recombination and stay in the ancestral background until the end of  
923 the sweep phase (i.e. at the time when the derived allele appeared, as we go backward in  
924 time) or (ii) coalesce all together at the end of the sweep phase. Actually, we slightly relax  
925 this second hypothesis and simply assume that their average coalescence time corresponds  
926 to the end of the sweep phase. The probability for each lineage to escape the sweep is  
927 approximately

$$q = 1 - e^{-4drN \log(\alpha)/\alpha}$$

928 where  $r$  is the recombination rate per generation and per bp. Because lineages can only  
929 coalesce if they are in the same background (derived with probability  $(1 - q)^2$  or ancestral

930 with probability  $q^2$ ), we assume that the average coalescence rate during the sweep is

$$\mu_{sweep} = (1 - q)^2 \frac{1}{\tau} + q^2 \frac{1}{2N}$$

931 where

$$\tau = 8N \log(\alpha) / \alpha$$

932 is the duration of the sweep (in generations). In this formula,  $\frac{1}{\tau}$  approximates the average  
933 coalescence rate for two lineages not escaping the sweep, which follows from our assump-  
934 tion that the average coalescence time is  $\tau$ , and  $\frac{1}{2N}$  is the standard neutral coalescence  
935 rate which applies to two lineages having escaped the sweep.

936 **Supplementary figures**

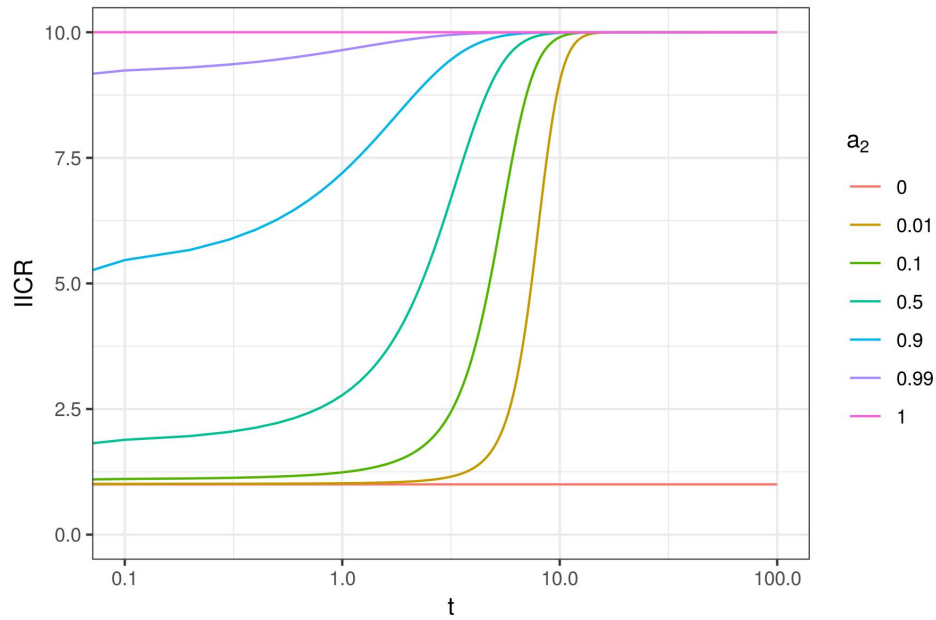


Figure S1: IICR curves for a panmictic model with  $K = 2$  classes of genomic regions with constant size. Same as Figure 1 with  $\lambda_1 = 1$ ,  $\lambda_2 = 10$  and time from 0 to 100 (in log10 scale)



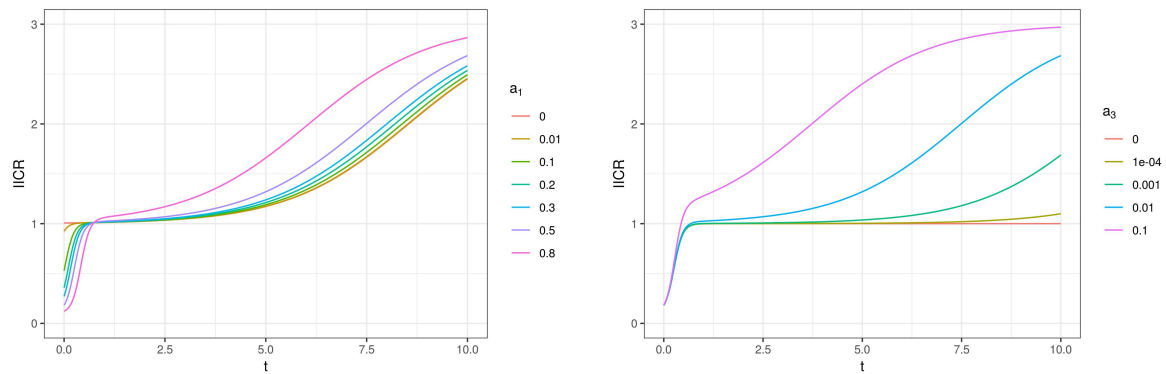


Figure S2: IICR for a panmictic model with  $K = 3$   $\lambda_i$  values such that  $\lambda_1 < 1$ ,  $\lambda_2 = 1$  and  $\lambda_3 > 1$ . Same as Figure 2 except that time is plotted in natural scale.

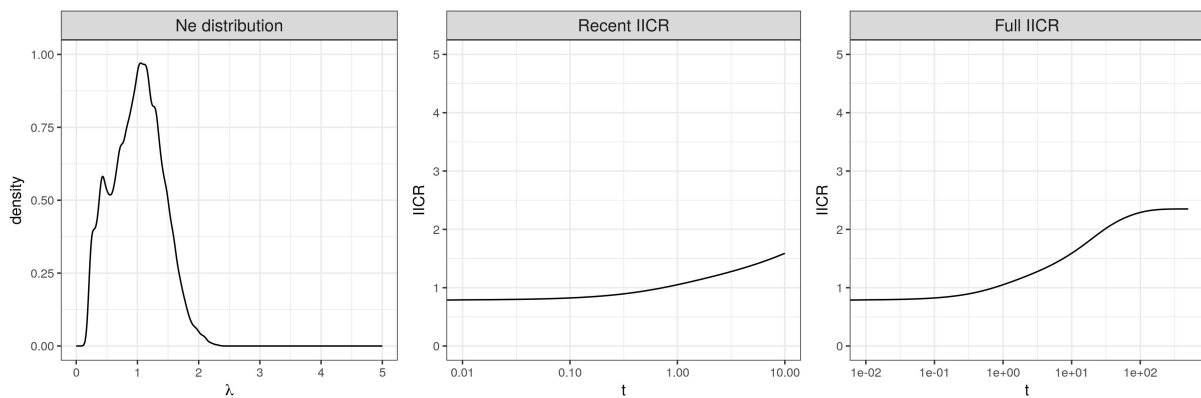


Figure S3: IICR obtained when removing low  $N_e$  values from the distribution estimated by Elyashiv et al. (2016). This truncated distribution (rescaled to have a mean of 1 as the others) is shown on the left panel. The associated IICR is shown until  $t = 10$  (middle panel) or  $t = 500$  (right panel), in log10 scale.

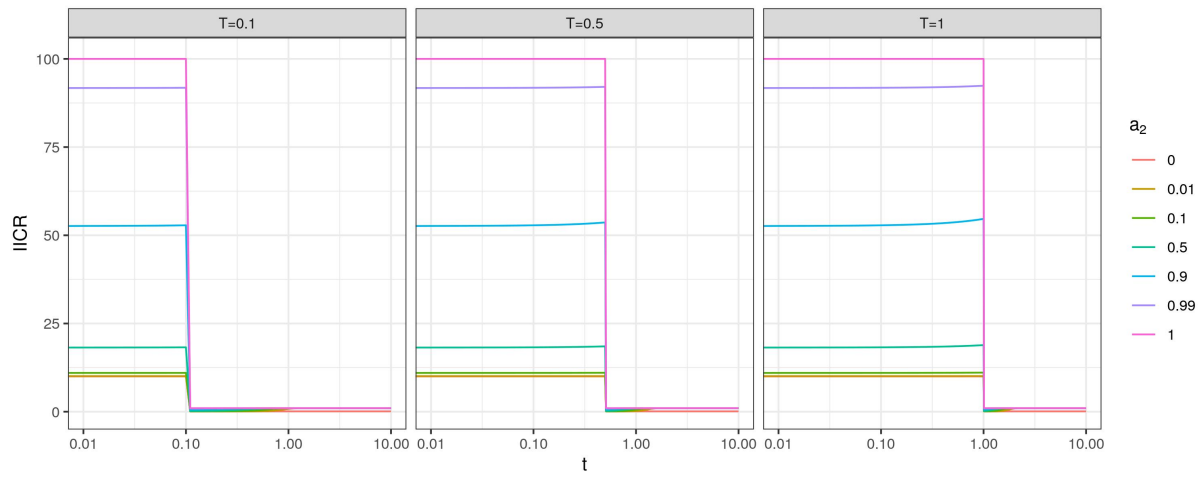


Figure S4: IICR curves for a panmictic model with a recent 100 fold expansion and  $K = 2$  classes of genomic regions. Same as Figure 4 with a stronger population expansion (100 fold vs 5 fold).

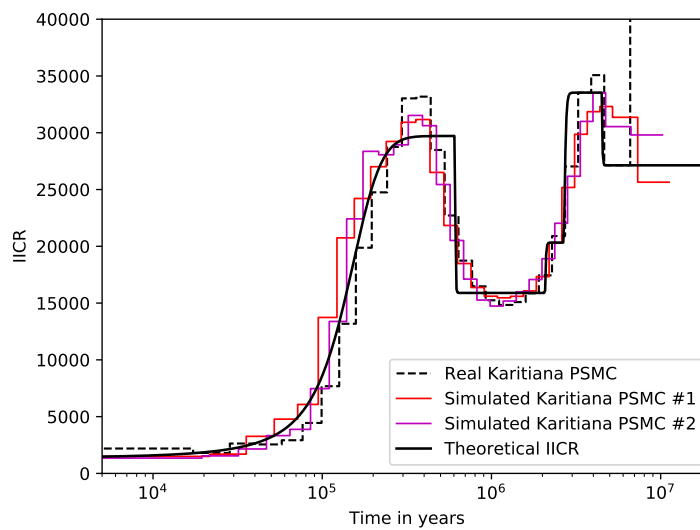


Figure S5: PSMC curves of simulated data under a non-stationary  $n$ -island model. We show in black the exact IICR corresponding to an inferred  $n$ -island model for a Karitiana individual in Arredondo et al. (2021). In color, we show various PSMC curves obtained by independently simulating genomic sequences under this structured model. The real PSMC curve for this Karitiana individual is represented by the dashed plot (Prado-Martinez et al., 2013). The horizontal axis is the time in years, with a generation time of 25 years. The vertical axis is the diploid population size. Times and population sizes were scaled assuming a mutation rate  $\mu=1.25e-8$ .

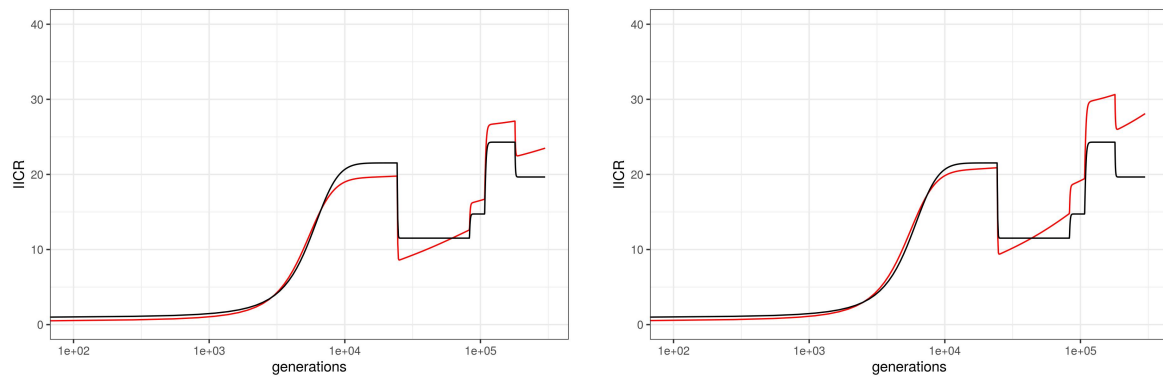


Figure S6: IICRs for demographic models combining population structure and linked selection in humans. Same as Figure 6, bottom panel, except that  $\lambda$  values greater than 2 (left) or 3 (right) were filtered out from the distribution in order to mimic a situation where loci under balancing selection could be detected and removed before computing the IICR. The resulting truncated distribution was rescaled.

## References

- 937
- 938 Arguello, J. R., Laurent, S., and Clark, A. G. (2019). Demographic history of the human  
939 commensal *Drosophila melanogaster*. *Genome biology and evolution*, 11(3):844–854.
- 940 Arredondo, A., Mourato, B., Nguyen, K., Boitard, S., Valcarce, W. R. R., Noûs, C.,  
941 Mazet, O., and Chikhi, L. (2021). Inferring number of populations and changes in  
942 connectivity under the n-island model. *Heredity*.
- 943 Charlesworth, B. (2009). Effective population size and patterns of molecular evolution  
944 and variation. *Nature Reviews Genetics*, 10(3):195.
- 945 Charlesworth, B. (2010). *Elements of evolutionary genetics*. Roberts Publishers.
- 946 Charlesworth, B., Morgan, M., and Charlesworth, D. (1993). The effect of deleterious  
947 mutations on neutral molecular variation. *Genetics*, 134(4):1289–1303.
- 948 Chikhi, L., Rodriguez, W., Grusea, S., Santos, P., Boitard, S., and Mazet, O. (2018).  
949 The IICR (inverse instantaneous coalescence rate) as a summary of genomic diversity:  
950 insights into demographic inference and model choice. *Heredity*, 120:13–24.
- 951 Comeron, J. M. (2017). Background selection as null hypothesis in population genomics:  
952 insights and challenges from *Drosophila* studies. *Philosophical Transactions of the Royal  
953 Society B: Biological Sciences*, 372(1736):20160471.
- 954 Durrett, R. (2008). *Probability models for DNA sequence evolution*. Springer.
- 955 Elyashiv, E., Sattath, S., Hu, T. T., Strutsovsky, A., McVicker, G., Andolfatto, P., Coop,  
956 G., and Sella, G. (2016). A genomic map of the effects of linked selection in *Drosophila*.  
957 *PLoS genetics*, 12(8):e1006130.

- 958 Ewing, G. and Hermisson, J. (2010). Msms: a coalescent simulation program including  
959 recombination, demographic structure and selection at a single locus. *Bioinformatics*,  
960 26(16):2064–2065.
- 961 Ewing, G. B. and Jensen, J. D. (2016). The consequences of not accounting for background  
962 selection in demographic inference. *Molecular ecology*, 25(1):135–141.
- 963 Gossmann, T. I., Woolfit, M., and Eyre-Walker, A. (2011). Quantifying the variation in  
964 the effective population size within a genome. *Genetics*, 189(4):1389–1402.
- 965 Grusea, S., Rodriguez, W., Boitard, S., Chikhi, L., and Mazet, O. (2018). Coalescence  
966 times for three genes are sufficient to detect population structure. *Journal of Mathe-*  
967 *matical Biology*, xxx(x):xxx–xxx.
- 968 Gutenkunst, R. N., Hernandez, R. D., Williamson, S. H., and Bustamante, C. D. (2009).  
969 Inferring the joint demographic history of multiple populations from multidimensional  
970 SNP frequency data. *PLoS Genetics*, 5(10):e1000695.
- 971 Herbots, H. M. J. D. (1994). *Stochastic models in population genetics: genealogy and*  
972 *genetic differentiation in structured populations*. PhD thesis.
- 973 Hill, W. G. and Robertson, A. (1966). The effect of linkage on limits to artificial selection.  
974 *Genetical Research*, 8(3):269–294.
- 975 Jensen, J. D., Payseur, B. A., Stephan, W., Aquadro, C. F., Lynch, M., Charlesworth,  
976 D., and Charlesworth, B. (2019). The importance of the neutral theory in 1968 and 50  
977 years on: a response to kern and hahn 2018. *Evolution*, 73(1):111–114.
- 978 Jiménez-Mena, B., Bataillon, T., et al. (2016a). Heterogeneity in effective population

- 979 size and its implications in conservation genetics and animal breeding. *Conservation*  
980 *genetics resources*, 8(1):35–41.
- 981 Jiménez-Mena, B., Tataru, P., Brøndum, R. F., Sahana, G., Guldbbrandtsen, B., and  
982 Bataillon, T. (2016b). One size fits all? direct evidence for the heterogeneity of genetic  
983 drift throughout the genome. *Biology letters*, 12(7):20160426.
- 984 Johri, P., Charlesworth, B., and Jensen, J. D. (2020). Toward an evolutionarily appro-  
985 priate null model: Jointly inferring demography and purifying selection. *Genetics*,  
986 215(1):173–192.
- 987 Johri, P., Riall, K., Becher, H., Excoffier, L., Charlesworth, B., and Jensen, J. D. (2021).  
988 The impact of purifying and background selection on the inference of population history:  
989 problems and prospects. *Molecular Biology and Evolution*. msab050.
- 990 Kaplan, N. L., Darden, T., and Hudson, R. R. (1988). The coalescent process in models  
991 with selection. *Genetics*, 120(3):819–829.
- 992 Kapopoulou, A., Pfeifer, S. P., Jensen, J. D., and Laurent, S. (2018). The demographic  
993 history of african drosophila melanogaster. *Genome biology and evolution*, 10(9):2338–  
994 2342.
- 995 Kern, A. D. and Hahn, M. W. (2018). The neutral theory in light of natural selection.  
996 *Molecular biology and evolution*, 35(6):1366–1371.
- 997 Kimura, M. (1983). *The neutral theory of molecular evolution*. Cambridge University  
998 Press.
- 999 Lapiere, M., Blin, C., Lambert, A., Achaz, G., and Rocha, E. P. (2016). The impact of

- 1000 selection, gene conversion, and biased sampling on the assessment of microbial demog-  
1001 raphy. *Molecular biology and evolution*, 33(7):1711–1725.
- 1002 Lewontin, R. C. (1974). *The genetic basis of evolutionary change*, volume 560. Columbia  
1003 University Press New York.
- 1004 Li, H. and Durbin, R. (2011). Inference of human population history from individual  
1005 whole-genome sequences. *Nature*, 475(7357):493–496.
- 1006 Mazet, O., Rodríguez, W., and Chikhi, L. (2015). Demographic inference using genetic  
1007 data from a single individual: Separating population size variation from population  
1008 structure. *Theoretical Population Biology*, 104:46–58.
- 1009 Mazet, O., Rodriguez, W., Grusea, S., Boitard, S., and Chikhi, L. (2016). On the im-  
1010 portance of being structured: instantaneous coalescence rates and human evolution—  
1011 lessons for ancestral population size inference. *Heredity*, 116(4):362–371.
- 1012 Nielsen, R., Williamson, S., Kim, Y., Hubisz, M. J., Clark, A. G., and Bustamante,  
1013 C. (2005). Genomic scans for selective sweeps using snp data. *Genome research*,  
1014 15(11):1566–1575.
- 1015 Ohta, T. (1992). The nearly neutral theory of molecular evolution. *Annual review of*  
1016 *ecology and systematics*, 23(1):263–286.
- 1017 Pouyet, F., Aeschbacher, S., Thiéry, A., and Excoffier, L. (2018). Background selec-  
1018 tion and biased gene conversion affect more than 95% of the human genome and bias  
1019 demographic inferences. *Elife*, 7:e36317.
- 1020 Prado-Martinez, J., Sudmant, P. H., Kidd, J. M., Li, H., Kelley, J. L., Lorente-Galdos,



- 1021 B., Veeramah, K. R., Woerner, A. E., O'Connor, T. D., Santpere, G., et al. (2013).  
1022 Great ape genetic diversity and population history. *Nature*, 499(7459):471–475.
- 1023 Rodríguez, W., Mazet, O., Grusea, S., Arredondo, A., Corujo, J. M., Boitard, S., and  
1024 Chikhi, L. (2018). The iicr and the non-stationary structured coalescent: towards demo-  
1025 graphic inference with arbitrary changes in population structure. *Heredity*, 121(6):663.
- 1026 Rougemont, Q. and Bernatchez, L. (2018). The demographic history of atlantic salmon  
1027 (*salmo salar*) across its distribution range reconstructed from approximate bayesian  
1028 computations. *Evolution*, 72(6):1261–1277.
- 1029 Rougemont, Q., Moore, J.-S., Leroy, T., Normandeau, E., Rondeau, E. B., Withler,  
1030 R. E., Van Doornik, D. M., Crane, P. A., Naish, K. A., Garza, J. C., et al. (2020).  
1031 Demographic history shaped geographical patterns of deleterious mutation load in a  
1032 broadly distributed pacific salmon. *PLoS genetics*, 16(8):e1008348.
- 1033 Rougeux, C., Bernatchez, L., and Gagnaire, P.-A. (2017). Modeling the multiple facets  
1034 of speciation-with-gene-flow toward inferring the divergence history of lake whitefish  
1035 species pairs (*coregonus clupeaformis*). *Genome biology and evolution*, 9(8):2057–2074.
- 1036 Roux, C., Fraïsse, C., Romiguier, J., Anciaux, Y., Galtier, N., and Bierne, N. (2016).  
1037 Shedding light on the grey zone of speciation along a continuum of genomic divergence.  
1038 *PLOS Biology*, 14(12):1–22.
- 1039 Schiffels, S. and Durbin, R. (2013). Inferring human population size and separation history  
1040 from multiple genome sequences. *Nature Genetics*, 8(46):919–925.
- 1041 Schrider, D. R., Shanku, A. G., and Kern, A. D. (2016). Effects of linked selective sweeps  
1042 on demographic inference and model selection. *Genetics*, 204(3):1207–1223.

- 1043 Sellinger, T., Abu Awad, D., and Tellier, A. (2021). Limits and convergence properties  
1044 of the sequentially markovian coalescent. *Molecular Ecology Resources*.
- 1045 Sheehan, S. and Song, Y. S. (2016). Deep learning for population genetic inference. *PLoS*  
1046 *computational biology*, 12(3):e1004845.
- 1047 Sjödin, P., Kaj, I., Krone, S., Lascoux, M., and Nordborg, M. (2005). On the meaning  
1048 and existence of an effective population size. *Genetics*, 169(2):1061–1070.
- 1049 Smith, J. M. and Haigh, J. (1974). The hitch-hiking effect of a favourable gene. *Genetics*  
1050 *Research*, 23(1):23–35.
- 1051 Walczak, A. M., Nicolaisen, L. E., Plotkin, J. B., and Desai, M. M. (2012). The structure  
1052 of genealogies in the presence of purifying selection: a fitness-class coalescent. *Genetics*,  
1053 190(2):753–779.
- 1054 Walsh, B. and Lynch, M. (2018). *Evolution and selection of quantitative traits*. Oxford  
1055 University Press.
- 1056 Zeng, K. and Charlesworth, B. (2011). The joint effects of background selection and  
1057 genetic recombination on local gene genealogies. *Genetics*, 189(1):251–266.