



**HAL**  
open science

# Sporadic Occurrence of Recent Selective Sweeps from Standing Variation in Humans as Revealed by an Approximate Bayesian Computation Approach

Guillaume Laval, Etienne Patin, Pierre Bouillier, Lluís Quintana-Murci

► **To cite this version:**

Guillaume Laval, Etienne Patin, Pierre Bouillier, Lluís Quintana-Murci. Sporadic Occurrence of Recent Selective Sweeps from Standing Variation in Humans as Revealed by an Approximate Bayesian Computation Approach. *Genetics*, 2021, 219 (4), 10.1093/genetics/iyab161 . pasteur-03560159

**HAL Id: pasteur-03560159**

**<https://pasteur.hal.science/pasteur-03560159>**

Submitted on 7 Feb 2022

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

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



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

1 **Sporadic Occurrence of Recent Selective Sweeps from Standing Variation in Humans as**  
2 **Revealed by an Approximate Bayesian Computation Approach**

3

4

5

6 Guillaume Laval<sup>1\*</sup>, Etienne Patin<sup>1</sup>, Pierre Bouillier<sup>2</sup> and Lluís Quintana-Murci<sup>1,3</sup>

7

8 <sup>1</sup>Human Evolutionary Genetics Unit, Institut Pasteur, UMR 2000, CNRS, Paris 75015,  
9 France.

10 <sup>2</sup>Department of Systems Biology, Harvard Medical School, Boston, Massachusetts 02115,  
11 USA.

12 <sup>3</sup>Human Genomics and Evolution, Collège de France, 75005 Paris, France

13

14

15 **\*Corresponding author**

16 E-mail: [glaval@pasteur.fr](mailto:glaval@pasteur.fr)

17

18 **Abstract**

19 During their dispersals over the last 100,000 years, modern humans have been exposed to a  
20 large variety of environments, resulting in genetic adaptation. While genome-wide scans for  
21 the footprints of positive Darwinian selection have increased knowledge of genes and  
22 functions potentially involved in human local adaptation, they have globally produced  
23 evidence of a limited contribution of selective sweeps in humans. Conversely, studies based  
24 on machine learning algorithms suggest that recent sweeps from standing variation are  
25 widespread in humans, an observation that has been recently questioned. Here, we sought to  
26 formally quantify the number of recent selective sweeps in humans, by leveraging  
27 approximate Bayesian computation and whole-genome sequence data. Our computer  
28 simulations revealed suitable ABC estimations, regardless of the frequency of the selected  
29 alleles at the onset of selection and the completion of sweeps. Under a model of recent  
30 selection from standing variation, we inferred that an average of 68 (from 56 to 79) and 140  
31 (from 94 to 198) sweeps occurred over the last 100,000 years of human history, in African  
32 and Eurasian populations, respectively. The former estimation is compatible with human  
33 adaptation rates estimated since divergence with chimps, and reveal numbers of sweeps per  
34 generation per site in the range of values estimated in *Drosophila*. Our results confirm the  
35 rarity of selective sweeps in humans and show a low contribution of sweeps from standing  
36 variation to recent human adaptation.

37

## INTRODUCTION

38

39 Evaluating the legacy of positive, Darwinian selection in humans has proved crucial for  
40 increasing our understanding of the genetic architecture of adaptive phenotypes (Fan et al.  
41 2016; Jeong and Di Rienzo 2014; Vitti et al. 2013). Genome-wide scans of selection have  
42 identified large numbers of signals of selective sweeps (Maynard Smith and Haigh 1974;  
43 Pritchard et al. 2010; Stephan et al. 1992) but, in turn, have produced limited evidence of this  
44 mode of selection to recent human adaptation (Hernandez et al. 2011). Conversely, a study  
45 based on the 1000 Genomes (1000G) data (Auton et al. 2015) and a simulation-based machine  
46 learning classifier has suggested that recent sweeps from standing variation ('soft sweeps')  
47 have been pervasive (Schridder and Kern 2017). Yet, the majority of such sweeps have been  
48 proposed to result from mis-classified neutral regions (Harris et al. 2018). Schridder and Kern  
49 (2017) identified ~862 or ~18 sweeps on average per population using the default or a more  
50 stringent probability threshold of 0.9 (Harris et al. 2018), showing the strong dependency  
51 between the number of reported sweeps and the detection thresholds used (Akey 2009;  
52 Pavlidis and Alachiotis 2017; Teshima et al. 2006). This highlights the difficulty to genuinely  
53 assess the real number of sweeps occurred in humans when identifying them individually  
54 (Jensen 2009; Li and Stephan 2006). Li and Stephan (2006) advised against methods counting  
55 the number of detected selection events due to false positives and the relatively low power to  
56 detect weak selection. Alternatively, many studies, including this one, aim to directly estimate  
57 the number of sweeps due to beneficial mutations that have occurred during a given period of  
58 time (e.g., Jensen et al. 2008; Li and Stephan 2006).

59 Inferring the rates of beneficial mutations in various species has been an active area of  
60 research in evolutionary biology. The most convincing results come from studies in  
61 *Drosophila*, where various methods have been implemented to formally estimate the expected  
62 number of selected substitutions per nucleotide site per generation ( $\lambda$ ). These estimates have

63 been calculated from divergence and/or polymorphism data, focusing on beneficial mutations  
64 that have fixed over a high number of generations (Andolfatto 2007; Eyre-Walker 2006;  
65 Jensen et al. 2008; Li and Stephan 2006; Macpherson et al. 2007; Sawyer and Hartl 1992;  
66 Sawyer et al. 2007; Smith and Eyre-Walker 2002). For example, Li and Stephan (2006)  
67 estimated  $\lambda$  considering beneficial mutations that have fixed over the last ~60,000 years in *D.*  
68 *melanogaster* populations (~900,000 generations, assuming 15 generations per year) (Barker  
69 1962; Pool 2015). Here, we sought to assess the extent of recent positive selection considering  
70 beneficial mutations occurring at shorter evolutionary time scales, such as those occurring  
71 after the split of African and Eurasian populations (e.g., in humans ~60,000 years represent  
72 ~2,150 generations assuming a generation interval of 28 years) (Fenner 2005; Moorjani et al.  
73 2016b). We considered both segregating and fixed beneficial mutations (incomplete and  
74 complete sweeps), as many candidate selection targets still segregate in humans (Fan et al.  
75 2016; Jeong and Di Rienzo 2014; Vitti et al. 2013) with few numbers of fixed or nearly-fixed  
76 differences being observed between human populations (Coop et al. 2009; Pritchard et al.  
77 2010). In this study, we leveraged genome-wide polymorphism data and neutrality statistics  
78 known to detect recent selection to assess the number of occurring sweeps ( $X$ ) using an  
79 approximate Bayesian computation (ABC) method (Beaumont et al. 2002; Jensen et al. 2008).

80 The ABC summary statistics we used here have been shown to be efficient in capturing  
81 true signals of selection, as they detect genome-wide excesses of candidate selected SNPs  
82 (i.e., SNPs harboring extreme values for a given neutrality statistic) within and near genes  
83 relative to intergenic regions, e.g., coding SNPs or *cis*-acting eQTLs *vs* intergenic SNPs  
84 (Barreiro et al. 2008; Fagny et al. 2014; Frazer et al. 2007; Jin et al. 2012; Kudaravalli et al.  
85 2009; Schmidt et al. 2019; Voight et al. 2006). Indeed, selective sweeps produce clusters of  
86 candidate SNPs in the vicinity of selection targets, whereas under neutrality candidate SNPs  
87 are more uniformly scattered (Voight et al. 2006). Specifically, we used the odds ratio for

88 selection (OR) (Kudaravalli et al. 2009), a statistic that depends on the ratio between the  
89 percentages of candidate SNPs identified within genic and intergenic regions (Barreiro et al.  
90 2008; Fagny et al. 2014). Under neutrality, false positive candidate SNPs are expected in  
91 genic and intergenic regions at the same proportion (Barreiro et al. 2008; Kudravalli et al.  
92 2009) ( $OR = 1$ , no excess of candidate SNPs). Otherwise, in case of higher rates of candidate  
93 SNPs in genes due to selection ( $OR > 1$ ), we assumed that OR correlates with  $X$  and provides  
94 suitable information to estimate  $X$ .

95 We applied the ABC method to several African, European and East Asian 1000G  
96 populations (Supplemental Material, Table S1), and explored various assumptions about the  
97 nature of the sweeps. To compare our results with the number of reported sweeps from  
98 standing variation (Schridder and Kern 2017), we simulated a single advantageous mutation  
99 per sweep region with a selection coefficient ( $s$ ) that ranged from 0.001 to 0.05 and a specific  
100 time ( $t$ ) that ranged from the present to 3,500 generations ago (or ~100,000 years ago). The  
101 frequency of the selected allele when the sweep begins ( $p_{start}$ ) was similarly ranged from  $1/2N$   
102 to 0.2. We labeled the sweeps from a *de novo* mutation ( $p_{start} = 1/2N$ ) and from standing  
103 variation ( $p_{start} > 1/2N$ ) (Hermisson and Pennings 2005; Innan and Kim 2004; Orr and  
104 Betancourt 2001; Przeworski et al. 2005) SDN and SSV, respectively (Peter et al. 2012).  
105 Although these sweeps correspond to the hard and soft sweeps simulated in Schridder and Kern  
106 (2017), we avoided such a terminology (Hermisson and Pennings 2005) since the nature --  
107 either hard or soft -- of detected sweeps is often ambiguous (Harris et al. 2018; Jensen 2014;  
108 Orr and Betancourt 2001).

109

## MATERIALS AND METHODS

110

### 111 Overview

112 Using ABC, we sought to jointly estimate the number of occurring sweeps  $X$ , the average  
113 strength and the average age of selection,  $S = 1/X \sum_1^X s_i$  and  $T = 1/X \sum_1^X t_i$ , in fifteen 1000G  
114 populations analyzed separately (five African, five European and five East Asian populations,  
115 see the section “1000 Genomes populations analyzed” below). We also considered three  
116 distinct categories of sweeps arbitrarily defined on the basis of the frequency of the selected  
117 allele at the onset of selection. Namely, we estimated with ABC  $X_1$ ,  $X_2$  and  $X_3$ , the number of  
118 sweeps with very low ( $1/2N \leq p_{start} < 0.01$ ), low ( $0.01 \leq p_{start} < 0.1$ ) and intermediate  
119 ( $0.1 \leq p_{start} < 0.2$ ) initial frequencies respectively (in our model  $X = X_1 + X_2 + X_3$ ). ABC  
120 point estimates (posterior means) and the 95% credible intervals (CI) boundaries are obtained  
121 from simulated whole-genome sequence (WGS) data best fitting with the 1000G empirical  
122 WGS data (see the section “ABC, acceptance rules and accuracy evaluation” below).

123 In each analyzed population, the  $X$ s ( $X$ ,  $X_1$ ,  $X_2$  and  $X_3$ ),  $S$  and  $T$  were estimated using  
124 whole-genome sequence data simulated to reproduce the 1000G empirical WGS data. The  
125 simulated and the empirical 1000G WGS data are each summarized by a vector of  $K$  ORs,  $K$   
126 being the number of neutrality statistic used. To compute these ORs, we considered as  
127 candidate SNPs of selection those with the most extreme values for neutrality statistics known  
128 to detect recent selection and previously found to be insensitive to background selection  
129 (BGS). The OR is based on the classic comparison between putatively neutral mutations and  
130 mutations potentially targeted or influenced by selection (Kimura 1977; McDonald and  
131 Kreitman 1991), labeled here ENVs (Evolutionary Neutral Variants) and PSVs (Possibly  
132 Selected Variants) respectively and defined below. Specifically, the OR computed for each  
133 WGS dataset and for each neutrality statistic separately is defined as

$$134 \quad OR = \left[ \frac{P(PSV|I(Candidate=1))}{P(ENV|I(Candidate=1))} \right] \left[ \frac{P(ENV|I(Candidate=0))}{P(PSV|I(Candidate=0))} \right], \quad (\text{equation 1})$$

135 with  $I(Candidate = 1)$  being the indicator function equal to 1 if the SNP is a candidate SNP  
136 of selection or 0 otherwise (Kudaravalli et al. 2009). The candidate SNPs were determined  
137 using genome-wide thresholds defining the 1% most extreme values obtained in simulated or  
138 in empirical 1000G data. The OR is thus expected close to one under neutrality (no sweep on  
139 the genome,  $X=0$ ). Otherwise, since selective sweeps produce clusters of candidate SNPs  
140 around targets of selection (Voight et al. 2006), an OR above one is an indicator of an excess  
141 of candidate SNPs in PSVs due to selection targeting beneficial mutations in this SNP class  
142 ( $X>0$ ). The ENVs are used here as the neutral baseline to control for the rate of false positive  
143 candidate SNPs of selection (Barreiro et al. 2008; Fagny et al. 2014; Kudaravalli et al. 2009).  
144 Finally, as a sanity check, we verified that the selection signals detected by the neutrality  
145 statistics used correspond to well-known examples of selective sweeps. We performed, in  
146 each 1000G population, a classic selection scan based on the same neutrality statistics used to  
147 estimate  $X$  (see the section “Detecting genomic regions enriched in candidate SNPs of  
148 selection” below).

149

### 150 **Odds ratio for selection in the simulated and 1000G populations**

151 For each whole-genome sequence dataset, either simulated or empirical, and each of the  $K$   
152 neutrality statistic used, the OR was computed using all SNPs genome-wide and the logistic  
153 region model set as follows (Kudaravalli et al. 2009):

$$154 \quad \text{Logit}[I(PSV = 1)] = \beta_1 I(Candidate = 1) + \varepsilon \quad (\text{equation 2})$$

155 with  $I(PSV = 1)$  being the indicator function equal to 1 for PSV SNPs or equal to 0  
156 otherwise,  $I(Candidate = 1)$  being the indicator function defined above, and  $\beta_1$  being the  
157 coefficient of the logistic regression (the constant term  $\beta_0$  was omitted). We used the ‘glm’ R  
158 package (family=binomial) to estimate  $\beta_1$  and  $\exp(\beta_1)$  to compute the OR for the effect of  
159 selection according to equation (1). Indeed, if  $\beta_1 > 0$ , this implies an enrichment effect for



160 PSVs among SNPs with selection signals due to the enrichment of SNPs with selection  
161 signals in the PSV SNP class (Kudaravalli et al. 2009).

162 In addition, when analyzing the 1000G populations, we carefully investigated the  
163 sensitivity of our ABC method to various estimations of the ORs. To this end, we also used  
164 the logistic regression-based method proposed by Kudravalli et al. (2009) to estimate the  
165 1000G ORs while controlling for various covariates such as the coverage (or sequencing  
166 quality), a feature of the data which was not simulated in our ABC model. Moreover, while  
167 we performed our simulations using human recombination maps (see the section “Simulating  
168 WGSs”) it is difficult to closely reproduce in genome-wide computer simulations the  
169 recombination pattern observed in humans. Recombination varies both along and across  
170 chromosomes in humans (Coop et al. 2008; Kong et al. 2010; Myers et al. 2005) and can alter  
171 the OR values. Due to the diminished hitchhiking effects of selection in high recombining  
172 regions, which result in reduced clusters of candidate SNPs, some selection signals can be  
173 penalized and ultimately contribute less to the ORs because of their genomic location only.  
174 The OR corrected for these covariates was computed using the logistic region model set as  
175 follows (Kudaravalli et al. 2009):

$$\begin{aligned} 176 \quad \text{Logit}[I(PSV = 1)] &= \beta_1 I(Candidate = 1) + [\beta_2 Cov + \\ 177 \quad \beta_3 Rec + \beta_4 NbSNP + \beta_5 Cov * Rec + \beta_6 Rec * NbSNP + \\ 178 \quad \beta_7 NbSNP * Cov] + \varepsilon, & \quad \text{(equation 3)} \end{aligned}$$

179 with  $I(PSV = 1)$  and  $I(Candidate = 1)$  defined above, and Rec the mean recombination  
180 rate in cM/bp obtained from HapMap recombination maps, Cov the mean coverage and  
181 NbSNP the number of SNPs, all computed with 100 kb sliding windows centered on each  
182 SNP. Finally, the algorithm used to compute the various ORs was set as follows:

183 (i) for each WGS dataset either simulated or empirical

184 (ii) get the definition of PSVs ( $PSV = 1$  or  $PSV = 0$  otherwise) randomly determined when  
185 building the simulated WGS data or according to the 1000G VEP annotations for empirical  
186 WGS data (see below)

187 (iii) for each of the  $K$  neutrality statistics used

188 (iv) get the definition of candidate SNPs ( $Candidate = 1$  or  $Candidate = 0$  otherwise)  
189 determined when building the simulated or empirical WGSs

190 (v) make the logistic regression model without covariates, equation (2)

191 (vi) compute the uncorrected OR merging all chromosomes in the logistic model

192  $OR_1 = \exp(\hat{\beta}_1)$ ,  $\hat{\beta}_1$  obtained from equation (2)

193 (vii) if empirical WGS data then

194 (viii) compute the uncorrected OR averaged across chromosomes

195  $OR_2 = 1/C \sum_c \exp(\hat{\beta}_{1,c})$ ,  $\hat{\beta}_{1,c}$  obtained for the  $c^{\text{th}}$  chromosome from equation (2)

196 (ix) compute covariates

197 (x) make the logistic regression model with covariates, equation (3)

198 (xi) compute the corrected OR merging all chromosomes in the logistic model

199  $OR_3 = \exp(\hat{\beta}_1)$ ,  $\hat{\beta}_1$  obtained from equation (3), which corresponds to that used by  
200 Kudaravalli et al. (2009)

201 (xii) compute the corrected OR averaged across chromosomes

202  $OR_4 = 1/C \sum_c \exp(\hat{\beta}_{1,c})$ ,  $\hat{\beta}_{1,c}$  obtained for the  $c^{\text{th}}$  chromosome from equation (3)

203 back to (iii)

204 back to (i)

205 With this algorithm, the  $10^5$  simulated WGSs used to perform the ABC estimations in  
206 each 1000G population are summarized by a matrix of  $OR_1$  with  $10^5$  rows and  $K$  columns. For  
207 each population, we systematically performed four rounds of ABC estimations using the four  
208 vectors of 1000G ORs of dimension  $(1 \times K)$  computed per population ( $OR_1$ ,  $OR_2$ ,  $OR_3$  and

209 OR<sub>4</sub>, see above). The matrix of simulated OR<sub>1</sub> and each vector of 1000G ORs were used  
210 without any modifications to jointly estimate  $X$ ,  $X_1$ ,  $X_2$ ,  $X_3$ ,  $S$  and  $T$  with a standard ABC  
211 method.

212

### 213 **Definitions of ENVs and PSVs**

214 The ENVs are Evolutionary Neutral Variants. In the simulations, they are neutral SNPs  
215 simulated in absence of selection, while in real data, these are intergenic SNPs far from the  
216 nearest gene and assumed to be unaffected by selection (Barreiro et al. 2008; Schmidt et al.  
217 2019; Voight et al. 2006). Conversely, the PSVs are Possibly Selected Variants; they include  
218 the potential targets of selection, that is mutations prone to alter individual phenotypes. In  
219 addition, we also incorporated in PSVs the SNPs near the potential targets of selection in  
220 order to capture the hitchhiking effects of linked selected variants on neutral polymorphisms.  
221 Indeed, the extent of clusters of candidate SNPs around the targets of selection, which  
222 depends on the intensity and on the age of selection, may provide valuable information for the  
223 estimation of  $X$ ,  $S$  and  $T$ , respectively. In simulations, the PSVs are thus the selected SNPs  
224 and nearby SNPs simulated in the same genomic regions. In real data, the PSVs are  
225 nonsynonymous, regulatory mutations and nearby SNPs, i.e., synonymous mutations, intronic  
226 mutations, which corresponds to the genic SNPs classically used (Barreiro et al. 2008;  
227 Schmidt et al. 2019; Voight et al. 2006). We also considered as PSVs the variants located  
228 upstream/downstream of genes as well as any presumed regulatory sites located in distal  
229 intergenic regions, e.g., SNPs located in transcription factor binding sites or in mature  
230 miRNAs. In our study, we considered that all nonsynonymous or regulatory mutations,  
231 denoted as ‘functional mutations’, are prone to alter phenotypes, as classically assumed in  
232 studies assessing the impact of natural selection (Barreiro et al. 2008; Fagny et al. 2014).

233 In simulated WGS data, the simulated regions containing PSVs, namely the PSV regions,  
234 are randomly defined and contain  $X$  selective sweeps of various intensity, age and frequency  
235 of the selected allele at the onset of selection (the remaining mutations are ENVs, see the  
236 section “Simulating WGSs”). It must be emphasized that PSV regions in simulated WGS data  
237 do not necessarily include a selected variant. Indeed, in the real data, PSV regions contain  
238 functional mutations, but only a fraction of these regions exhibit selection signals and may  
239 thus contain functional selected mutations. The other PSV regions behave as neutral regions  
240 as they only contain functional mutations with no detectable advantageous effects on fitness.

241 For each analyzed population, the PSVs and ENVs are determined using the 1000G VEP  
242 (Ensembl Variant Effect Predictor) annotations. According to our definitions, all mutations in  
243 or near genes, as well as intergenic mutations annotated as “TF\_binding\_site\_variant” or  
244 “mature\_miRNA\_variant” in the VEP files, were set to be PSVs,  $I(PSV = 1)$ . The rest of  
245 intergenic mutations annotated as “intergenic\_variant” were set to be ENVs  $I(PSV = 0)$   
246 (Supplemental Material). With such filters, ~70% of 1000G SNPs were considered to be  
247 PSVs (the remaining 30% are ENVs). We thus systematically reproduced these proportions in  
248 the simulated WGS data. We also dealt with some particular situations: these include the  
249 selected sites that are unknown functional variants annotated as “intergenic\_variant”, small  
250 regulatory regions located far from other PSVs and regulatory variants located in edges of  
251 PSVs tracts. Positive selection targeting such SNPs can bias downward our estimations  
252 because the selection signals (clusters of candidate SNPs) may expand to ENVs. This will  
253 increase the proportion of candidate SNPs of selection in ENVs resulting in lower empirical  
254 ORs, which will be ultimately interpreted by the model as a lower number of sweeps. To  
255 minimize such estimation biases, we annotated as PSVs all SNPs with a genome-wide  
256 significant enrichment in candidate SNPs measured 100 kb around them, since such  
257 enrichments due to clusters of candidate SNPs are indicative of positive selection (Voight et

258 al. 2006). This step has only a marginal effect on the total number of SNPs classified as PSVs  
259 but it can inflate ORs. We thus reproduced this step in the simulated WGS data.

260

## 261 **Simulating WGSs**

262 Computer simulations were performed to reproduce the 1000G WGSs data, i.e., ~3Gb of  
263 DNA sequences for 100 unrelated individuals sampled per population. To do so, we used a  
264 demographic model previously inferred for the analyzed populations. To evaluate inter-  
265 population statistics, such as XP-EHH (Sabeti et al. 2007), we used a three-population  
266 demographic model calibrated to replicate the allele frequency spectrum, population structure  
267 and linkage disequilibrium levels in African (YRI), European (CEU) and East Asian (CHB)  
268 populations (Grossman et al. 2013; Grossman et al. 2010; Schaffner et al. 2005). We  
269 privileged this model because it has been used to detect selective sweeps in the YRI, CEU and  
270 CHB 1000G populations on the basis of simulation-based approaches (Grossman et al. 2013;  
271 Pybus et al. 2015). This model incorporates an African expansion, an Out-of-Africa exodus  
272 ~100 kya (3,500 generations) followed by a bottleneck, a split of Eurasians into European and  
273 Asian populations ~58kya and various migration rates between continents ( $\sim 10^{-5}$  per haploid  
274 genome per generation). An interesting feature is the presence of a second bottleneck in each  
275 non-African population being four times stronger in Asia than in Europe (Pickrell et al. 2009).  
276 We systematically simulated triplets of African, European and East Asian populations. To  
277 simulate the  $X$  selective sweeps per WGS data in the focal population, we set  $X = 0$  in the  
278 two other populations used as reference for the XP-EHH computation and simulated the non-  
279 equilibrium demography of each population using the three-population demographic model  
280 presented above.

281 We used SLiM (Haller and Messer 2017) to simulate 5Mb regions using human  
282 recombination rates drawn from HapMap recombination maps (Frazer et al. 2007). We

283 simulated  $10^4$  neutrally evolving DNA regions and  $4 \times 10^3$  DNA regions with selection for  
284 each sweep scenario envisaged (see the results section). In the latter case, we simulated  
285 selection in the focal population (Africa, Europe or East Asia) by inserting a single  
286 advantageous mutation in the middle of the region. At generation  $t$  being the onset of  
287 selection, the frequency of this mutation  $p_{start}$  was randomly drawn from the allele frequency  
288 spectrum at the generation  $t$ . In all the sweep scenarios investigated, we used initial frequency  
289 of the selected mutation  $p_{start}$ , intensity of selection  $s$  and age of selection  $t$  with the same  
290 ranges of values, i.e., respectively ranged from  $1/2N$  to 0.2, from 0.001 to 0.05 and from  
291 present to 3,500 generations ago. We simulated long DNA regions to avoid premature  
292 truncations because selection signals can extent over mega bases for particularly strong  
293 sweeps, i.e.,  $\sim 2$ Mb in the LCT region (various estimates of  $s$  for rs4988235 range from 0.025  
294 to 0.069) (Chen et al. 2015; Peter et al. 2012; Tishkoff et al. 2007). Because SLiM is a  
295 forward-in-time simulator, the computation time depends on both the effective population size  
296  $N$  and the number of generations considered. We thus rescaled effective population sizes and  
297 times according to  $N/\lambda$  and  $t/\lambda$ , with  $\lambda = 10$ , and used rescaled mutation, recombination and  
298 selection parameters,  $\lambda\mu$ ,  $\lambda r$  and  $\lambda s$  (Haller and Messer 2017; Hoggart et al. 2007).

299 Lastly, simulated genomes were obtained by randomly drawing neutral simulated regions  
300 and regions simulated under the desired sweep scenario, some of which were considered to be  
301 PSVs, the rest being ENVs. The age  $t$  and the intensity  $s$  of selection, used to simulate the  $X$   
302 sweeps per WGS data located in the PSV regions, were randomly drawn from various  
303 distributions depending on the sweep scenario investigated. We assigned to each SNP an  
304 indicator function  $I(PSV)$  equal to 1 for the SNPs in PSV regions, which was subsequently  
305 used to compute the ORs (see above). In summary, the PSVs are simulated as tracts of SNPs  
306 with  $I(PSV = 1)$  randomly spread over the genomes to reproduce the proportions of ENVs  
307 and PSVs observed in the 1000G populations, with numbers of simulated SNPs matching

308 those observed in 1000G populations. In the following, for each ABC estimation performed,  
309 we used  $10^5$  simulated WGS data, namely the ABC simulations, with  $X$  randomly drawn from  
310 a uniform prior distribution,  $U(0, 300)$ . A graphical illustration summarizing the different  
311 steps to simulate WGS data can be found in the Supplemental Material.

312

### 313 **1000 Genomes populations analyzed**

314 Analyses were performed on the 1000 Genomes Project phase 3 data, focusing on African,  
315 European and East Asian populations and excluding populations with diverse continental or  
316 admixed ancestry. We analyzed 1,511 individuals from five African, five European and five  
317 East Asian populations (85 to 113 individuals per population). Phased data (SHAPEIT2)  
318 (Delaneau et al. 2012), ancestral/derived states and VEP annotations were downloaded from  
319 the 1000G Project website.

320

### 321 **Neutrality statistics**

322 For each simulated and 1000G WGS, we computed six neutrality statistics. We used the  
323 haplotype-based neutrality statistics,  $iHS$  (Voight et al. 2006),  $DIND$  (Barreiro et al. 2009)  
324 and  $\Delta iHH$  (Grossman et al. 2010), which compare the haplotypes carrying the derived and  
325 ancestral alleles. We computed two pairwise  $XP-EHHs$  (Sabeti et al. 2007), which compare  
326 the haplotypes carrying the derived allele between the focal population and two other  
327 populations of differing continental origin used as reference (for the 1000G WGS the two  
328 reference populations were chosen from YRI, CEU or CHB). We also used the Fay and Wu's  
329  $H$  (F&W-H) (Fay and Wu 2000), which detects deviations from the neutral allele frequency  
330 spectrum in genomic regions. We used a sliding-windows approach for these computations  
331 (100 kb windows centered on each SNP) (Fagny et al. 2014). The sliding-windows began and  
332 ended 50kb from the edges of the 5Mb simulated regions, to prevent truncation in the 100 kb

333 sliding windows (a similar approach was applied to the 1000G chromosomes). As iHS,  
334 DIND,  $\Delta iHH$  and XP-EHH are sensitive to the inferred ancestral/derived state, we computed  
335 these statistics only when the derived state was determined unambiguously and we  
336 normalized them by DAF bin (Fagny et al. 2014; Voight et al. 2006) (mutations grouped by  
337 DAF bin, from 0 to 1, in increments of 0.025). We applied classic filters by excluding minor  
338 alleles frequencies below 0.05 and we minimized the false-positive discovery by excluding  
339 SNPs with a DAF below 0.2, as the power to detect positive selection has been shown to be  
340 reduced at such low frequencies (Fagny et al. 2014; Voight et al. 2006). Finally, for each  
341 neutrality statistic, we defined the candidate SNPs of selection (top 1% of SNPs genome-  
342 wide). For iHS, DIND,  $\Delta iHH$  and XP-EHHS, we considered extreme values indicative for  
343 selection targeting the derived alleles.

344 Because background selection may generate spurious positive selection signatures  
345 genome-wide (Coop et al. 2009; Hernandez et al. 2011; Pritchard et al. 2010), we excluded  
346 the  $F_{ST}$ , Tajima's D and other neutrality statistics previously found to be affected by BGS  
347 (Zeng et al. 2006). Indeed, the differences in allele frequencies between populations are  
348 expected to be exacerbated in regions affected by BGS, a pattern that can be confounded with  
349 positive selection (Coop et al. 2009; Pickrell et al. 2009; Pritchard et al. 2010). We only  
350 retained Fay and Wu's H and the haplotype-based statistics also previously found to be  
351 insensitive to BGS (Fagny et al. 2014; Zeng et al. 2006), e.g., the haplotype-based statistics  
352 use haplotypes carrying the ancestral alleles as internal controls that should be affected by  
353 BGS to a similar extent.

354

### 355 **ABC, acceptance rules and accuracy evaluation**

356 We used the 'abc' R package and the standard ABC method (Beaumont et al. 2002), in which  
357 posteriors are built from accepted simulated parameters subsequently adjusted by local linear



358 regression (method="Loclinear" in the 'abc' package). The accepted simulated parameters are  
359 those which provide the best fit with empirical data (the 'abc' parameter 'tol' was set to be  
360 equal to 0.005). We used the mean of the posterior distribution as point estimate and  
361 computed the 95% credible interval from these accepted parameters. To test the accuracy of  
362 our estimations, we compared estimated and simulated parameter values,  $\hat{\theta}_i$  and  $\theta_i$   
363 respectively, using classical accuracy indices: the relative error  $rE$  (i.e. difference between  
364 estimated and true values, expressed as a proportion of the true value,  $rE = (\hat{\theta}_i - \theta_i)/\theta_i$ ,  
365  $i = 1, \dots, J$ ), the relative root of the mean square error,  $rRMSE$  (i.e. the root of the MSE  
366 expressed as a proportion of the true value), and the proportion of true values within the 90%  
367 credible interval of estimates,  $90\%COV = \frac{1}{J} \sum_1^J 1(q_1 < \theta_i < q_2)$  where  $1(C)$  is the indicative  
368 function (equal to 1 when  $C$  is true, 0 otherwise) and  $q_1$  and  $q_2$ , the corresponding percentiles  
369 of the posterior distributions.

370

### 371 **ABC estimations obtained simulating BGS in PSVs or changing the PSVs annotation**

372 To check the insensitivity of our estimation to background selection, we simulated  
373 widespread background selection in PSVs: we simulated DNA regions with a fraction of  
374 mutations targeted by purifying selection, determined following the coding region maps, in  
375 which 2/3 of mutations were considered to be nonsynonymous with selection coefficients  
376 randomly drawn from the Gamma distribution of fitness effects determined in Boyko et al.  
377 (2008). These new simulations are used as PSV regions in the ABC simulations and we re-  
378 estimated the parameters accordingly. In addition, for the reasons described above, the SNPs  
379 with a genome-wide significant enrichment in candidate SNPs 100 kb around them were  
380 annotated as PSVs to perform the ABC estimations. To check these estimations, we  
381 performed another round of estimations after excluding from the analysis all ENVs with such  
382 a significant enrichment (some selection signals due to selection targets located close to the

383 edges of PSVs were thus truncated while other selection signals that cover several genes and  
384 intergenic regions were diminished, Table S2C).

385

### 386 **Detecting genomic regions enriched in candidate SNPs of selection**

387 The enrichment of candidate SNPs in a given genomic region was determined by means of a  
388 combined selection score (CSS) computed using the same neutrality statistics used to estimate  
389  $X$ . For each SNP and each of the  $K = 6$  neutrality statistics used in the ABC estimations, we  
390 computed the proportion of candidate SNPs in a 100 kb window around the considered SNP.  
391 We next determined the empirical  $P$ -values ( $P$ ) for these proportions and combined them into  
392 a single combined selection score using a Fisher score,  $CSS = -2 \sum_1^6 \log(P_i)$  (Deschamps et  
393 al. 2016). The rationale behind such a composite approach (Deschamps et al. 2016; Grossman  
394 et al. 2013; Grossman et al. 2010) is that neutrality statistics are more strongly correlated for  
395 positively selected variants than for neutral sites (Grossman et al. 2010). Consequently, false  
396 positives may harbor extreme values for a few neutrality statistics only, whereas SNPs  
397 genuinely selected (or nearby SNPs) should harbor extreme values for several statistics  
398 together, a feature captured by the combined score. Finally, the genomic regions enriched in  
399 candidate SNPs were defined as consecutive SNPs with genome-wide significant CSS values,  
400  $P < 0.01$  (gaps of maximal length of 100kb are allowed). In these enriched regions, the  
401 maximum CSS value was used as a proxy of the strength of selection signal; high CSS values  
402 are expected for SNPs targeted by strong selection.

403 For each population and each enriched genomic region, we computed an overlap score  
404 calculated as the number of populations for which the same region was identified. Overlap  
405 scores were calculated either within or between continents (upper limits of 5 and 10 for  
406 within- and between-continent overlap scores respectively, because we analyzed three  
407 continents of five populations each). Within continents, the overlap scores range from 1 (the

408 considered region was identified in a single population) to 5 (the considered genomic region  
409 was identified in all populations of the same continental origin). Between continents, the  
410 overlap scores range from 0 (the considered region was never identified in another continent)  
411 to 10 (the considered region was identified in all populations of differing continental origin).  
412

### 413 **Data availability**

414 Genome-wide data and SNPs annotations can be downloaded from the 1000 Genome website.  
415 The software used to compute the neutrality statistics can be downloaded from  
416 <https://github.com/h-e-g/selink>. The scripts used to perform the main analysis can be  
417 downloaded from <https://github.com/h-e-g/ABCnumSweeps>. Supplemental material,  
418 including Figures S1 to S26, Tables S1 to S4 (Tables S2-S4 are supplied as Excel files) and  
419 Files S1 and S2 can be found with this article online.  
420

## 421 **RESULTS**

### 422 **Odds Ratios for selection capture the number of occurring sweeps**

423 We estimated the number of selective sweeps using ABC and summary statistics that measure  
424 the enrichment of selection signals in PSVs regions, relative to ENV regions; a graphical  
425 illustration summarizing the rationale behind our approach is shown in Figures 1 and S1.  
426 Because ABC inferences are based on summary statistics that exhibit a monotonic  
427 relationship with the parameter to estimate, we first present the relationships between the odds  
428 ratio for selection (OR) and the number of sweeps in African, European and East Asian  
429 simulated populations. We explored two different scenarios (Figure 2, Supplemental Material,  
430 Figures S2 and S3): in the first, we simulated  $X$  incomplete sweeps per whole-genome  
431 sequence (WGS) dataset (Materials and Methods) by randomly drawing  $s$  and  $t$  from flat  
432 distributions,  $U(0.001, 0.05)$  and  $U(0, 100)$  *kya*, and excluding complete sweeps using a

433 rejection algorithm (Figure 2). The underlying distributions of  $s$  and also  $t$  are therefore  
434 enriched in low values (Figure 2A), since (all other things being equal) complete sweeps tend  
435 to be stronger or older than incomplete sweeps. In a second, more realistic scenario, we  
436 simulated  $X$  sweeps per WGS data while keeping the complete sweeps (Figure S3). The age  
437 of selection was randomly drawn from a flat distribution,  $U(0, 100)$  *kya*, whereas we aimed  
438 to reproduce the excess of mutations with low or moderate effect on fitness (Boyko et al.  
439 2008) with a distribution enriched in low values (a mixture between a Gamma distribution  
440 with 60% of  $s < 0.01$  and a L-shape distribution with 90% of  $s < 0.01$ , Figure S3A).

441 For each scenario, we simulated WGS data with fixed numbers of sweeps,  $X =$   
442  $[0, 50, 100, 150]$ , using human recombination maps and a demographic model previously  
443 inferred for these populations (Materials and Methods). Our simulations clearly show that the  
444 OR well captures the enrichment of candidate SNPs in PSVs owing to the action of selection  
445 (Figures 2B and S3B). PSV regions with more than 1% of candidate SNPs in the vicinity of  
446 the selection targets contribute to the greater proportions of candidate SNPs in PSVs relative  
447 to ENVs ( $OR = 1$  for  $X = 0$  vs  $OR > 1$  otherwise). As assumed, the six ORs increase with  
448 the number of simulated sweeps, resulting in a monotonic relationship between  $X$  and the OR  
449 over the parameter space investigated. In addition, the OR values depend on the frequencies  
450 of alleles at the onset of selection (Figure S2); SSV typically have weaker effects on linked  
451 sites (Pritchard et al. 2010; Przeworski et al. 2005), which reduces the numbers of candidate  
452 SNPs around selection targets and ultimately the OR. However, very similar ORs were  
453 obtained for SDN and SSV with  $p_{start}$  lower than 0.01. This is expected since such sweeps are  
454 difficult to distinguish due to similar signatures when  $p_{start}$  tends to  $1/2N$  (Ferrer-Admetlla et  
455 al. 2014; Jensen 2014; Peter et al. 2012). In light of this, we focused on the estimation of the  
456 numbers of sweeps corresponding to these two sweep models merged ( $p_{start} < 1\%$ ,  $X_1$ ).

457

458 **Odds Ratios for selection are moderately sensitive to demographic assumptions**

459 The ORs values follow the power to detect selection as previously assessed for iHS and XP-  
460 EHH (highest, intermediate and lowest power in the African, European and East Asian  
461 demography respectively) (Pickrell et al. 2009). The lower European and East Asian ORs  
462 relative to those in Africans (Figures 2B and S3B), which reflect lower enrichments of  
463 candidate SNPs in PSVs, are in agreement with a reduced power to detect selection in  
464 bottlenecked relative to expanding populations (Fagny et al. 2014; Grossman et al. 2013;  
465 Gunther and Schmid 2011; Huff et al. 2010; Pickrell et al. 2009). However, little differences  
466 in the behavior of ORs were observed across the populations simulated with contrasting  
467 demographic histories (i.e., African expansion vs Eurasian bottleneck). Whereas the East  
468 Asian bottleneck is stronger than that in Europe, the East Asian and European simulated ORs  
469 are virtually identical (e.g., 2.3 vs 2.4 in average for DIND with  $X = 100$ , Figure 2B). These  
470 observations suggest that the ABC estimations of  $X$  will be moderately affected by potential  
471 demographic misspecifications.

472

473 **Accuracy of the ABC estimations**

474 To assess the accuracy of the estimations in the two sweep scenarios explored, we treated the  
475 simulations shown in Figures 2 (incomplete sweeps) and S3 (complete and incomplete  
476 sweeps) as pseudo-empirical data for which the parameter values are known. We found  
477 virtually unbiased estimations of  $X$  (Figure 3). An in-depth analysis (Figures S4-S12) showed  
478 a similar accuracy for the estimations of  $X_1$ ,  $X_2$  and  $X_3$  under the two assumed scenarios  
479 (Figures S4 and S8). When simulating complete sweeps, the neutrality statistics with low  
480 power to detect sweeps near to fixation, iHS or DIND, remained informative as they captured  
481 the enrichment of candidate SNPs around beneficial mutations at fixation (Figure S3). This  
482 illustrates the importance of incorporating in PSVs the SNPs in the vicinity of selection

483 targets to leverage information provided by the hitchhiking effects of selected variants. Our  
484 computer simulations showed that incorporating nearby SNPs in PSVs, which encompass  
485 70% of the simulated genome, did not dilute the true adaptive signals. Using such proportions  
486 of simulated PSVs allow to define large 1000G PSVs regions by incorporating all the  
487 functional mutations and neighboring variants to capture a large fraction of the existing  
488 selective sweeps, including selected regulatory mutations located up/down-stream genes (see  
489 the next section). However, because of largely overlapping posterior distributions, our method  
490 is not able to distinguish between estimations obtained when the real number of sweeps is low  
491 (real  $X < 10$ ) and those obtained under neutrality (real  $X = 0$ ). Finally, we noticed a  
492 diminished accuracy of  $S$  and  $T$ , e.g., underestimations and overestimations of  $S$  and  $T$   
493 respectively when the selection is strong and recent in average (Figures S5 and S9), together  
494 with large CIs that sometimes exceed the range of simulated priors (Figures S7 and S11).

495 The estimations of  $X$  based on the combined six ORs performed reasonably well. They  
496 are equal -- in expectation -- to the numbers of occurring sweeps (Figure 3), with CIs that  
497 consistently predict the range of uncertainty within which the true parameters are expected  
498 (Figures S6 and S10), as shown by the 90% *COV* indicated in Figures S4 and S8. Thus, our  
499 estimations of  $X$  appeared to be neither systematically biased upward by neutral genomic  
500 regions with spurious signatures of selection nor systematically biased downward by sweep  
501 regions with few candidate SNPs (i.e., false negatives). As expected, false positives do not  
502 contribute to ORs, owing to the use of ENVs as an internal neutral control. In the absence of  
503 selection ( $X = 0$  in simulated pseudo-empirical data), the estimations of  $X$  are close to 0 as  
504 ORs are close to one. Under selection ( $X > 0$  in simulated pseudo-empirical data), the  
505 expected proportion of false negatives is accounted by the model thanks to the simulation of  
506 parameters that normally drive the statistical power to detect sweeps, such as the age of  
507 selection, the intensity of selection or the demography. Altogether, our method accurately

508 estimates  $X$  from OR values; yet, the estimations were obtained under a known demography  
509 whereas model-based methods are expected to be affected by partially unknown demography.  
510 We thus evaluated the effect of incorrect demographic assumptions on the estimations of  $X$  in  
511 1000G populations.

512

### 513 **Low numbers of selective sweeps in the 1000G populations**

514 We applied our method to each African, European and East Asian 1000G population analyzed  
515 separately. To capture the selection signals as exhaustively as possible, the PSVs defined  
516 using the VEP annotation of the 1000G data (Materials and Methods, Figure S13) correspond  
517 to large genomic regions as variants near functional mutations must be included in PSVs to  
518 account for the clustering of candidate SNPs around putative selection targets. We used the  
519 ABC simulations with the same proportion of simulated PSVs (Figure 3) and 1000G ORs  
520 corrected for genomic variation in coverage, mutation and recombination rates (Kudaravalli et  
521 al. 2009) (Materials and Methods, Figure S14), which were found to be compatible with the  
522 simulated ORs used to perform the estimations (Figure S15). While the incomplete sweep  
523 scenario is not applicable to real populations since it ignores complete sweeps, the results  
524 obtained under this scenario are reported for comparison (Tables 1 and 2, Figures 4 and 5).  
525 We noticed a stronger influence of this assumption in Eurasia (Table 1), mainly in East Asia,  
526 likely because sweeps can reach fixation faster due to lower effective size (large proportions  
527 of the fixed differences between human populations are due to fixation events outside of  
528 Africa, mainly in East Asia) (Coop et al. 2009).

529 In all cases, we found lower numbers of sweeps than previously detected (Pybus et al.  
530 2015; Schrider and Kern 2017), i.e., 116 sweeps in average across all populations (Table 1)  
531 with 198 [150-240] sweeps at most in the CHS population (Table S2A) (see also Figure 4,  
532 Table S2B and Figures S16 and S17). This result was found to be robust to the inclusion of

533 background selection in simulated PSVs and to the mode of computation of the 1000G ORs  
534 (Materials and Methods), whether they are corrected for various covariates (Table 1) or not  
535 (Table S2C). In Africa, we found that an average of ~70 sweeps occurred over ~100,000 years  
536 account for the genome-wide selection signal, or probably a few less since older sweeps are  
537 ignored. Such sweeps are not expected to have inflated much the 1000G ORs computed from  
538 neutrality statistics known to have low power to detect old selection events in humans  
539 (Grossman et al. 2013; Voight et al. 2006).

540       Beside the estimations of  $X$ , we obtained average strengths of selection close to 0.01 and  
541 average ages of selection compatible with expected ages of selection, e.g., recent ages when  
542 considering incomplete sweeps only (Table 2, Tables S2A and S2B). However, the CIs  
543 obtained for  $S$  and  $T$  covered the whole range of the simulated priors (Figures S7 and S11),  
544 preventing precise inferences for these parameters. In addition, we mapped the genomic  
545 regions significantly enriched in candidate SNPs by means of a classic genome-wide scan of  
546 selection performed with the same neutrality statistics used to estimate  $X$  (Materials and  
547 Methods). Reassuringly, the regions that most contribute to the ORs capture many examples  
548 of previously reported selected regions (Table S3), including iconic cases such as *TLR5* in  
549 Africa, *LCT* in Northern Europe and *EDAR* in East Asia (Fan et al. 2016; Jeong and Di  
550 Rienzo 2014; Vitti et al. 2013) (Figures S18 and S19, File S1 for a short description of the  
551 genomic regions identified).

552       We found a low overlap of selection signals between continents (Materials and Methods),  
553 confirming previous observations (Voight et al. 2006), while the most enriched regions tend  
554 to be highly shared across the 1000G populations from the same continent (Figure S20,  
555 Tables S3A and S3B), in agreement with recent ages of selection estimated under the  
556 incomplete sweep scenario. Indeed, with the neutrality statistics used herein, incomplete  
557 sweeps, younger by definition and thus more continent-specific than complete sweeps that



558 may have occurred before the split of Eurasian populations, result in the highest enrichment in  
559 candidate SNPs and consequently the lowest empirical P-values in our selection scan. We also  
560 tested the sensitivity of our ABC estimations to the use of neutral reference populations for  
561 the XP-EHH statistics. We thus performed another round of ABC estimations after removing  
562 the XP-EHH ORs, which can be impacted by shared selective sweeps. We found that  
563 estimates of  $X$  were unchanged (Figure S21), with similar average ages of selection, e.g., 36.0  
564 [23.5-45.8] kya under the incomplete sweep scenario, suggesting that our estimations are not  
565 sensitive to the use of neutral outgroups.

566

### 567 **Consistent detection of low numbers of selective sweeps in humans**

568 To test the discrepancies observed between our results and the high numbers of sweeps  
569 previously reported (Schridder and Kern 2017), we used our highest estimations of  $X$ , i.e.,  
570 those obtained with the 1000G ORs averaged across chromosomes (Table 1). First,  
571 estimations performed using different combinations of neutrality statistics, or modifying  
572 several technical steps, were consistently found low (Figure S21, Table S2C). For example,  
573 excluding the most correlated neutrality statistics and their corresponding ORs does not affect  
574 our ABC estimations (Figures S21 and S22). We next checked the influence of the  
575 demographic model assumed (Figures S23 and S24). We adopted a strategy based on stress  
576 tests to evaluate the differences in numbers of estimated sweeps when modifying the  
577 demographic assumptions, e.g., replacing an expansion with a bottleneck or increasing the  
578 bottleneck intensity. We performed this by swapping empirical and simulated ORs from  
579 differing continental regions, e.g., estimations performed using the 1000G African ORs with a  
580 “wrong” East Asian simulated demography and inversely. In doing so, we obtained similar  
581 numbers of sweeps (Tables S2D and S2E).

582 The graphical explanations given in Figure S23 illustrate how the new estimations of  $X$   
583 may differ. For example, we re-analyzed the European populations under an East Asian  
584 demographic model, and found higher numbers of sweeps than initially estimated, i.e., 141  
585 [95-187] vs 115 [68-160] sweeps on average (Table S2E vs Table 1). These higher estimations  
586 of  $X$  are due to lower simulated ORs in East Asia compared to Europe (Figure 2B, Figure  
587 S3B). The estimations are increased because the East Asian ORs simulated with higher  $X$   
588 values provided a better fit with the European 1000G ORs (see the  $X_{Europe}^{(ASIdemo)}$  in the Figure  
589 S23). The low discrepancies obtained were confirmed using pseudo-empirical datasets  
590 simulated with bottlenecks four times stronger (or weaker) than in the ABC simulations used  
591 for the estimations (Figure S24). Overall, we did not obtain point estimates that largely  
592 exceeded the estimations shown in Tables 1 and S2. We acknowledge that new ABC  
593 estimations based on different demographic models may differ. Nevertheless, our results  
594 suggest that they should be of the same order of magnitude.

595 Lastly, we performed estimations of  $X$  in the studied 1000G populations using various  
596 distributions of  $s$ . In all cases, we found low numbers of selective sweeps, although we  
597 noticed a moderate influence of the shape of the distributions of  $s$  on the estimations of  $X$   
598 (Table S2F). As expected, the estimations of  $X$  are inversely correlated with the intensity of  
599 selection (e.g., two sweeps of high intensity should result in ORs similar to those obtained  
600 with more sweeps of lower intensity). For example, we estimated 70 [35-106] and 126 [90-  
601 162] sweeps in average assuming a flat (~20% of  $s < 0.01$  with an average equal to ~0.025)  
602 or a L-shape distribution of  $s$  (~90% of  $s < 0.01$  with an average equal to ~0.007). Under  
603 these two extreme scenarios, the posterior distributions obtained are largely non-overlapping  
604 but the point estimates are of the same order of magnitude (Table S2F). Note that we did not  
605 consider this flat distribution to set the composite distribution of  $s$  used to perform our main  
606 estimations (Tables 1 and 2, Figures 4 and 5), because assuming similar proportions of

607 sweeps with high and low intensity is an unrealistic assumption. To exclude that our  
608 estimations are due to an underestimation of  $X$  when  $X$  is large, we verified the ability of our  
609 method to infer high numbers of selection events when they are present in the data using  
610 pseudo-empirical WGS simulated with high numbers of sweeps. We used priors extended up  
611 to 1,000 sweeps per population ( $X \sim U(0, 1000)$ ) and found unbiased estimations of  $X$  (Figure  
612 S25). We then re-estimated  $X$  in each 1000G population, and found similar estimations (Table  
613 S2G), supporting further our findings of low numbers of sweeps occurred.

614

### 615 **A trend towards higher numbers of sweeps occurred in non-African populations**

616 We found more sweep signals in Eurasian populations (Table 1), an observation that has been  
617 attributed to a greater extent of genetic adaptation outside Africa due to drastic changes of  
618 environmental pressures (Coop et al. 2009; Granka et al. 2012; Pybus et al. 2015; Schrider  
619 and Kern 2017). For example, we detected a recent European-specific selection signal likely  
620 due to the late Pleistocene warming at the end of the last ice age in Northern hemisphere  
621 (Cooper et al. 2015), a cold climatic period which could favor both light skin pigmentation  
622 and increased sensitivity to UV-induced melanoma (Key et al. 2016; Lopez et al. 2014).  
623 Regulatory mutations, protective against melanoma as they downregulate the *PLEK2* gene in  
624 skin cells exposed to the sun (GTEx database) (Ardlie et al. 2015), exhibit genome-wide  
625 significant signals of positive selection in Europe only (Table S3, Figure S26, and File S2).  
626 However, we cannot exclude that the higher  $X$  estimated in Eurasia, particularly in East Asia  
627 (see the non-overlapping CIs in Table 1), are due to the demographic model used (our stress  
628 tests and simulations also indicated that the estimations of  $X$  in East Asia could be lower  
629 under models with weaker bottlenecks, Table S2E, Figure S24). A more formal comparison  
630 would require a perfectly known demography or the inclusion of alternative demographic  
631 models.

632 With respect to the numbers of sweeps as a function of the initial frequency of selected  
633 alleles,  $X_2$  and  $X_3$  ( $0.01 \leq p_{start} < 0.2$ ) were found in excess compared to  $X_1$  ( $X_1 < X_2 + X_3$ ,  
634 Figure 5), particularly in Eurasian populations. This observation may reflect the loss of  
635 selected alleles at very low frequency ( $X_1$ ) during the first generations of selection, which is  
636 exacerbated in bottlenecked populations relative to expanding populations. However, in all  
637 cases our results confirm sweeps from standing variation as main drivers of adaptation  
638 (Schridder and Kern 2017), since sweeps from *de novo* mutations are contained in the  $X_1$   
639 category (numbers of SDN ranged from 0 to  $X_1$ ). In practice, the neutrality statistics used to  
640 identify sweeps detect large frequency changes (Ferrer-Admetlla et al. 2014; Schridder and  
641 Kern 2017) whereas they are not sensitive to moderate changes in allele frequencies due to  
642 polygenic adaptation (Field et al. 2016; Pritchard and Di Rienzo 2010; Pritchard et al. 2010;  
643 Stephan 2016); some major loci controlling complex traits with extreme phenotypes recently  
644 selected (Perry et al. 2014) may resemble incomplete sweeps.

645

646

## DISCUSSION

647 Our study shows that ABC approaches can be helpful to formally assess the extent and nature  
648 of recent positive selection in humans, and provide valuable information about the proportion  
649 of true sweeps identified by selection scans. A result emerging from this study is the low  
650 contribution of sweeps to recent human adaptation, including sweeps from standing variation,  
651 which supports the rarity of recent selective sweeps in humans (Coop et al. 2009; Harris et al.  
652 2018; Hernandez et al. 2011; Pritchard et al. 2010). In our analyses, ~70% of the genome was  
653 considered as potentially influenced by selection and that all significant signals of positive  
654 selection found in intergenic regions were accounted for suggests that we capture a large  
655 fraction of the existing selective sweeps. Yet, our estimations are far lower than numbers of  
656 recent sweeps individually detected using classic selection scans; they correspond to ~35% of

657 the numbers of putatively selected regions identified (Table S3C). They also correspond to  
658 ~30% and ~10% of sweeps identified using machine learning algorithms based on simulations  
659 of the same demographic model used herein (Pybus et al. 2015) and on similar models of  
660 selection from standing variation (Schrider and Kern 2017), respectively. Such discrepancies  
661 are in agreement with studies in *Drosophila* that have shown that the numbers of sweeps  
662 identified using scans of selection are roughly an order of magnitude greater than would be  
663 predicted, likely because of high false positive rates (Jensen 2009; Teshima et al. 2006).

664 To evaluate the validity of our estimations, we compared estimations of adaptation rates  
665 obtained in human and *Drosophila*. We first compared with  $\alpha$ , the fraction of nonsynonymous  
666 mutations driven to fixation by positive selection (Messer and Petrov 2013), estimated in  
667 1000G African populations using ABC and 29,925 nonsynonymous fixed differences with  
668 chimps (Uricchio et al. 2019). With the number of sweeps estimated in Africa (Table 1)  
669 translated in number of sweeps per generation, we predicted  $\alpha$  in the human lineage assuming  
670 divergence times between human and chimp of ~7.9Mya (Moorjani et al. 2016a) (Tables 3  
671 and S4A). These predictions assume constant numbers of beneficial mutations per generation  
672 whereas these numbers did vary over time due to changes in environmental conditions or in  
673 effective sizes (Hawks et al. 2007). Our predicted values of  $\alpha$ , which can be higher given that  
674 we also considered selection targeting regulatory sites, are quite similar to those previously  
675 estimated by Uricchio et al. (2019) and Zhen et al. (2021) (Table 3). This indicates that our  
676 results are compatible with previous estimates obtained at broader evolutionary scales and  
677 using different methods and data.

678 We next compared our results with the numbers of beneficial mutations per generation per  
679 nucleotide site ( $\lambda$ ) estimated in *Drosophila* under a recurrent hitchhiking model (Kaplan et al.  
680 1989; Przeworski 2002; Wiehe and Stephan 1993) with sweeps occurring at random locations  
681 in the genome (Andolfatto 2007; Jensen et al. 2008; Li and Stephan 2006; Macpherson et al.

682 2007). The estimations of  $\lambda$  previously obtained in *Drosophila* largely vary across different  
683 studies (Table 4). Using the per generation numbers of sweeps obtained with our approach,  
684 which can predict  $\lambda$  in humans (Tables 4 and S4B), we found  $\lambda$  values similar or even higher  
685 than estimated in *Drosophila*. A recent analysis suggests that the species with a smaller  
686 population size and greater complexity (i.e., humans) may have stronger and/or more  
687 abundant new beneficial mutations than other species with much larger population sizes (i.e.,  
688 mice and *D. melanogaster*) (Zhen et al. 2021), suggesting that the complexity of the  
689 organisms affects adaptation rates. Otherwise, we found lower  $\lambda$  values when the intensity of  
690 selection estimated in *Drosophila* is much lower than that estimated in this study (Table 4).  
691 For example, the  $\lambda$  estimated by Andolfatto (2007) is approximately two orders of magnitude  
692 higher than ours, which roughly corresponds to the expected differences in  $N$  values between  
693 *Drosophila* and humans ( $N \sim 10^6$  vs  $N \sim 10^4$ ). Empirical evidence in great apes also suggests  
694 that the numbers of experienced sweeps per unit of time may increase with effective  
695 population sizes (Nam et al. 2017). Beside our predicted  $\lambda$ , our results suggest a lower impact  
696 of beneficial mutations on the human genome-wide diversity compared to *Drosophila*, as  
697 indicated by lower  $2N\lambda S$  values (Table S4B) (the reduction of heterozygosity below the  
698 neutral level as a function of the recombination distance to selected sites is essentially  
699 determined by this compound parameter) (Wiehe and Stephan 1993). Applying our  
700 methodology to *Drosophila* data may help to understand if lower  $2N\lambda S$  are due to lower  $\lambda$  or  
701 lower  $2NS$  driving lower rates of detectable beneficial mutations.

702 From a methodological standpoint, our ABC framework appears to provide unbiased  
703 estimations of  $X$  when the distribution on  $s$  is known. Although here we leveraged widely  
704 used neutrality statistics, the OR can be computed for any other kind of statistics, keeping in  
705 mind that background selection should be simulated (Johri et al. 2020) when using statistics  
706 sensitive to this selection regime. Our study presents however some limitations. We noticed a

707 dependency on the  $s$  distributions used, which suggests that a careful evaluation of their  
708 influence on the estimation of  $X$  is needed. We also observed a lack of accuracy in the  
709 estimation of  $S$  or  $T$ , indicating that other summary statistics that better describe the genomic  
710 extent of selection signals should be incorporated in selection analysis. Another important  
711 point that should be considered is the calibration of the period of selection studied, which  
712 depends on the sensitivity of the neutrality statistics to the age of selection simulated. For  
713 example, Li and Stephan (2006) used statistics with low power to detect selection events older  
714 than 60,000 years in *Drosophila*. Here, we used neutrality statistics with a low power to detect  
715 sweeps that are older than 100,000 years in humans (Grossman et al. 2013; Voight et al.  
716 2006). Such sweeps, which are ignored by the model, marginally contribute to empirical data  
717 (i.e., the 1000G ORs) and should have thus a low influence on the final estimations (Li and  
718 Stephan 2006). Because simulating whole chromosomes with selection is barely tractable, we  
719 concatenated shorter regions together and avoided the computation of neutrality statistics  
720 across the junctions to ensure that they were computed over genomic regions simulated  
721 according to human recombination rates. Also, our ABC estimates were obtained using the  
722 common assumption that the demography is known, which can affect the validity of model-  
723 based estimations. The minor differences in point estimates observed within continents are  
724 likely due to specific local demographic histories not accounted by the model, e.g., some  
725 lower  $X$  in the Finnish population are potentially due to strong bottlenecks (Bulik-Sullivan et  
726 al. 2015) not accounted by the CEU demography. However, our stress tests indicate that the  
727 low numbers of sweeps estimated herein are likely not due to demographic misspecifications,  
728 although more efforts are needed to refine the estimations. Future work simulating, in  
729 reasonable computation times, whole chromosomes with beneficial mutations in genic regions  
730 and incorporating recent population admixture as revealed by ancient genomes (Skoglund and  
731 Mathieson 2018) into ABC should provide more refined estimations of the true numbers of

732 selective sweeps occurred in humans. Furthermore, while our results suggest that our method  
733 is not sensitive to the use of neutral outgroups, the impact of this assumption needs to be  
734 formally evaluated using more complex scenarios of shared selection with adaptive gene flow  
735 (Laso-Jadart et al. 2017; Patin et al. 2017; Refoyo-Martinez et al. 2019).

736 To conclude, our study provides novel evidence in support of the paucity of selective  
737 sweeps in humans, with posteriors distributions of  $X$  that consistently excluded a number of  
738 sweeps higher than  $\sim 250$  in all investigated scenarios. The observed enrichments of candidate  
739 SNPs in and near genes were found to be lower than those obtained with computer  
740 simulations assuming high numbers of sweeps (Figure S25A-C). Furthermore, the  $\alpha$   
741 adaptation rates predicted from our estimations were found compatible with those estimated  
742 using different sources of data and methods. Our estimated numbers of recent sweeps in  
743 humans are helpful to better understand the outcomes of classic scans of selection, keeping in  
744 mind that genomic regions evolving under complex selection regimes can only be detected  
745 with powerful methods such as those recently implemented (Schridder and Kern 2017). The  
746 methodology proposed in this study also provides an alternative to methods implemented in  
747 *Drosophila*, and can also be applied to great apes (Nam et al. 2017; Prado-Martinez et al.  
748 2013; Schmidt et al. 2019) or other mammals (Freedman et al. 2016; Ihle et al. 2006; Librado  
749 et al. 2015; Roux et al. 2015; Stella et al. 2010). Adapting these approaches to study more  
750 complex selection regimes, such as polygenic adaptation, should provide a broader, and more  
751 precise, picture of the recent adaptive history of humans and other species.

752

### 753 **Acknowledgments**

754 We thank the IT infrastructure of Institut Pasteur (Paris) for the management of computational  
755 resources. We also thank two anonymous reviewers for their fruitful comments and  
756 suggestions which greatly improved the quality of our manuscript. This work was supported



757 by the Agence Nationale de la Recherche (ANR) grant “DEMOCHIPS” ANR-12-  
758 BSV7-0012. The Human Evolutionary Genetics laboratory is supported by the Institut  
759 Pasteur, the Collège de France, the CNRS, the Fondation Allianz-Institut de France, and the  
760 French Government’s Investissement d’Avenir program, Laboratoires d’Excellence  
761 ‘Integrative Biology of Emerging Infectious Diseases’ (ANR-10-LABX-62-IBEID) and  
762 ‘Milieu Intérieur’ (ANR-10-LABX-69-01).

763

#### 764 **Literature Cited**

765 Akey JM. 2009. Constructing genomic maps of positive selection in humans: Where do we go  
766 from here? *Genome Res.* 19(5):711-722.

767 Andolfatto P. 2007. Hitchhiking effects of recurrent beneficial amino acid substitutions in the  
768 *drosophila melanogaster* genome. *Genome Res.* 17(12):1755-1762.

769 Ardlie KG, DeLuca DS, Segre AV, Sullivan TJ, Young TR, Gelfand ET, Trowbridge CA,  
770 Maller JB, Tukiainen T, Lek M et al. 2015. The genotype-tissue expression (gtex)  
771 pilot analysis: Multitissue gene regulation in humans. *Science.* 348(6235):648-660.

772 Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO, Marchini JL,  
773 McCarthy S, McVean GA, Abecasis GR. 2015. A global reference for human genetic  
774 variation. *Nature.* 526(7571):68-74.

775 Barker JSF. 1962. The estimation of generation interval in experimental populations of  
776 *drosophila*. *Genet Res.* 3(3):388-404.

777 Barreiro LB, Ben-Ali M, Quach H, Laval G, Patin E, Pickrell JK, Bouchier C, Tichit M,  
778 Neyrolles O, Gicquel B et al. 2009. Evolutionary dynamics of human toll-like  
779 receptors and their different contributions to host defense. *PLoS Genet.*  
780 5(7):e1000562.

781 Barreiro LB, Laval G, Quach H, Patin E, Quintana-Murci L. 2008. Natural selection has  
782 driven population differentiation in modern humans. *Nat Genet.* 40(3):340-345.

783 Beaumont MA, Zhang W, Balding DJ. 2002. Approximate bayesian computation in  
784 population genetics. *Genetics.* 162(4):2025-2035.

785 Boyko AR, Williamson SH, Indap AR, Degenhardt JD, Hernandez RD, Lohmueller KE,  
786 Adams MD, Schmidt S, Sninsky JJ, Sunyaev SR et al. 2008. Assessing the  
787 evolutionary impact of amino acid mutations in the human genome. *PLoS Genet.*  
788 4(5):e1000083.

789 Bulik-Sullivan BK, Loh PR, Finucane HK, Ripke S, Yang J, Patterson N, Daly MJ, Price AL,  
790 Neale BM. 2015. Ld score regression distinguishes confounding from polygenicity in  
791 genome-wide association studies. *Nat Genet.* 47(3):291-295.

792 Chen H, Hey J, Slatkin M. 2015. A hidden markov model for investigating recent positive  
793 selection through haplotype structure. *Theor Popul Biol.* 99:18-30.

794 Coop G, Pickrell JK, Novembre J, Kudaravalli S, Li J, Absher D, Myers RM, Cavalli-Sforza  
795 LL, Feldman MW, Pritchard JK. 2009. The role of geography in human adaptation.  
796 *PLoS Genet.* 5(6):e1000500.

797 Coop G, Wen XQ, Ober C, Pritchard JK, Przeworski M. 2008. High-resolution mapping of  
798 crossovers reveals extensive variation in fine-scale recombination patterns among  
799 humans. *Science.* 319(5868):1395-1398.

800 Cooper A, Turney C, Hughen KA, Brook BW, McDonald HG, Bradshaw CJ. 2015.  
801 Paleoecology. Abrupt warming events drove late pleistocene holarctic megafaunal  
802 turnover. *Science.* 349(6248):602-606.

803 Delaneau O, Marchini J, Zagury JF. 2012. A linear complexity phasing method for thousands  
804 of genomes. *Nat Methods.* 9(2):179-181.

805 Deschamps M, Laval G, Fagny M, Itan Y, Abel L, Casanova JL, Patin E, Quintana-Murci L.  
806 2016. Genomic signatures of selective pressures and introgression from archaic  
807 hominins at human innate immunity genes. *Am J Hum Genet.* 98(1):5-21.

808 Eyre-Walker A. 2006. The genomic rate of adaptive evolution. *Trends Ecol Evol.* 21(10):569-  
809 575.

810 Fagny M, Patin E, Enard D, Barreiro LB, Quintana-Murci L, Laval G. 2014. Exploring the  
811 occurrence of classic selective sweeps in humans using whole-genome sequencing  
812 data sets. *Mol Biol Evol.* 31(7):1850-1868.

813 Fan S, Hansen ME, Lo Y, Tishkoff SA. 2016. Going global by adapting local: A review of  
814 recent human adaptation. *Science.* 354(6308):54-59.

815 Fay JC, Wu CI. 2000. Hitchhiking under positive darwinian selection. *Genetics.* 155(3):1405-  
816 1413.

817 Fenner JN. 2005. Cross-cultural estimation of the human generation interval for use in  
818 genetics-based population divergence studies. *Am J Phys Anthropol.* 128(2):415-423.

819 Ferrer-Admetlla A, Liang M, Korneliussen T, Nielsen R. 2014. On detecting incomplete soft  
820 or hard selective sweeps using haplotype structure. *Mol Biol Evol.* 31(5):1275-1291.

821 Field Y, Boyle EA, Telis N, Gao ZY, Gaulton KJ, Golan D, Yengo L, Rocheleau G, Froguel  
822 P, McCarthy MI et al. 2016. Detection of human adaptation during the past 2000  
823 years. *Science.* 354(6313):760-764.

824 Frazer KA, Ballinger DG, Cox DR, Hinds DA, Stuve LL, Gibbs RA, Belmont JW, Boudreau  
825 A, Hardenbol P, Leal SM et al. 2007. A second generation human haplotype map of  
826 over 3.1 million snps. *Nature.* 449(7164):851-861.

827 Freedman AH, Schweizer RM, Ortega-Del Vecchyo D, Han E, Davis BW, Gronau I, Silva  
828 PM, Galaverni M, Fan Z, Marx P et al. 2016. Demographically-based evaluation of  
829 genomic regions under selection in domestic dogs. *PLoS Genet.* 12(3):e1005851.

830 Granka JM, Henn BM, Gignoux CR, Kidd JM, Bustamante CD, Feldman MW. 2012. Limited  
831 evidence for classic selective sweeps in african populations. *Genetics*. 192(3):1049-  
832 1064.

833 Grossman SR, Andersen KG, Shlyakhter I, Tabrizi S, Winnicki S, Yen A, Park DJ, Griesemer  
834 D, Karlsson EK, Wong SH et al. 2013. Identifying recent adaptations in large-scale  
835 genomic data. *Cell*. 152(4):703-713.

836 Grossman SR, Shlyakhter I, Karlsson EK, Byrne EH, Morales S, Frieden G, Hostetter E,  
837 Angelino E, Garber M, Zuk O et al. 2010. A composite of multiple signals  
838 distinguishes causal variants in regions of positive selection. *Science*. 327(5967):883-  
839 886.

840 Gunther T, Schmid KJ. 2011. Improved haplotype-based detection of ongoing selective  
841 sweeps towards an application in arabidopsis thaliana. *BMC Res Notes*. 4:232.

842 Haller BC, Messer PW. 2017. Slim 2: Flexible, interactive forward genetic simulations. *Mol*  
843 *Biol Evol*. 34(1):230-240.

844 Harris RB, Sackman A, Jensen JD. 2018. On the unfounded enthusiasm for soft selective  
845 sweeps ii: Examining recent evidence from humans, flies, and viruses. *PLoS Genet*.  
846 14(12):e1007859.

847 Hawks J, Wang ET, Cochran GM, Harpending HC, Moyzis RK. 2007. Recent acceleration of  
848 human adaptive evolution. *Proc Natl Acad Sci USA*. 104(52):20753-20758.

849 Hermisson J, Pennings PS. 2005. Soft sweeps: Molecular population genetics of adaptation  
850 from standing genetic variation. *Genetics*. 169(4):2335-2352.

851 Hernandez RD, Kelley JL, Elyashiv E, Melton SC, Auton A, McVean G, Sella G, Przeworski  
852 M. 2011. Classic selective sweeps were rare in recent human evolution. *Science*.  
853 331(6019):920-924.

854 Hoggart CJ, Chadeau-Hyam M, Clark TG, Lampariello R, Whittaker JC, De Iorio M, Balding  
855 DJ. 2007. Sequence-level population simulations over large genomic regions.  
856 *Genetics*. 177(3):1725-1731.

857 Huff CD, Harpending HC, Rogers AR. 2010. Detecting positive selection from genome scans  
858 of linkage disequilibrium. *BMC Genomics*. 11:8.

859 Ihle S, Ravaoarimanana I, Thomas M, Tautz D. 2006. An analysis of signatures of selective  
860 sweeps in natural populations of the house mouse. *Mol Biol Evol*. 23(4):790-797.

861 Innan H, Kim Y. 2004. Pattern of polymorphism after strong artificial selection in a  
862 domestication event. *Proc Natl Acad Sci USA*. 101(29):10667-10672.

863 Jensen JD. 2009. On reconciling single and recurrent hitchhiking models. *Genome Biol Evol*.  
864 1:320-324.

865 Jensen JD. 2014. On the unfounded enthusiasm for soft selective sweeps. *Nat Commun*.  
866 5:5281.

867 Jensen JD, Thornton KR, Andolfatto P. 2008. An approximate bayesian estimator suggests  
868 strong, recurrent selective sweeps in drosophila. *PLoS Genet*. 4(9):e1000198.

869 Jeong C, Di Rienzo A. 2014. Adaptations to local environments in modern human  
870 populations. *Curr Opin Genet Dev*. 29:1-8.

871 Jin W, Xu S, Wang H, Yu Y, Shen Y, Wu B, Jin L. 2012. Genome-wide detection of natural  
872 selection in african americans pre- and post-admixture. *Genome Res*. 22(3):519-527.

873 Johri P, Charlesworth B, Jensen JD. 2020. Towards an evolutionarily appropriate null model:  
874 Jointly inferring demography and purifying selection. *Genetics*. 215(1):173-192.

875 Kaplan NL, Hudson RR, Langley CH. 1989. The "hitchhiking effect" revisited. *Genetics*.  
876 123(4):887-899.

877 Key FM, Fu Q, Romagne F, Lachmann M, Andres AM. 2016. Human adaptation and  
878 population differentiation in the light of ancient genomes. *Nat Commun*. 7:10775.

879 Kimura M. 1977. Preponderance of synonymous changes as evidence for the neutral theory of  
880 molecular evolution. *Nature*. 267(5608):275-276.

881 Kong A, Thorleifsson G, Gudbjartsson DF, Masson G, Sigurdsson A, Jonasdottir A, Walters  
882 GB, Jonasdottir A, Gylfason A, Kristinsson KT et al. 2010. Fine-scale recombination  
883 rate differences between sexes, populations and individuals. *Nature*. 467(7319):1099-  
884 1103.

885 Kudaravalli S, Veyrieras JB, Stranger BE, Dermitzakis ET, Pritchard JK. 2009. Gene  
886 expression levels are a target of recent natural selection in the human genome. *Mol*  
887 *Biol Evol*. 26(3):649-658.

888 Laso-Jadart R, Harmant C, Quach H, Zidane N, Tyler-Smith C, Mehdi Q, Ayub Q, Quintana-  
889 Murci L, Patin E. 2017. The genetic legacy of the indian ocean slave trade: Recent  
890 admixture and post-admixture selection in the makranis of pakistan. *Am J Hum Genet*.  
891 101(6):977-984.

892 Li H, Stephan W. 2006. Inferring the demographic history and rate of adaptive substitution in  
893 *drosophila*. *PLoS Genet*. 2(10):e166.

894 Librado P, Sarkissian CD, Ermini L, Schubert M, Jonsson H, Albrechtsen A, Fumagalli M,  
895 Yang MA, Gambo C, Seguin-Orlando A et al. 2015. Tracking the origins of yakutian  
896 horses and the genetic basis for their fast adaptation to subarctic environments. *Proc*  
897 *Natl Acad Sci USA*. 112(50):E6889-E6897.

898 Lopez S, Garcia O, Yurrebaso I, Flores C, Acosta-Herrera M, Chen H, Gardeazabal J,  
899 Careaga JM, Boyano MD, Sanchez A et al. 2014. The interplay between natural  
900 selection and susceptibility to melanoma on allele 374f of *slc45a2* gene in a south  
901 european population. *PLoS One*. 9(8):e104367.

902 Macpherson JM, Sella G, Davis JC, Petrov DA. 2007. Genomewide spatial correspondence  
903 between nonsynonymous divergence and neutral polymorphism reveals extensive  
904 adaptation in drosophila. *Genetics*. 177(4):2083-2099.

905 Maynard Smith J, Haigh J. 1974. The hitch-hiking effect of a favourable gene. *Genet Res*.  
906 23(1):23-35.

907 Mcdonald JH, Kreitman M. 1991. Adaptive protein evolution at the adh locus in drosophila.  
908 *Nature*. 351(6328):652-654.

909 Messer PW, Petrov DA. 2013. Frequent adaptation and the mcdonald-kreitman test. *Proc Natl*  
910 *Acad Sci USA*. 110(21):8615-8620.

911 Moorjani P, Amorim CE, Arndt PF, Przeworski M. 2016a. Variation in the molecular clock of  
912 primates. *Proc Natl Acad Sci USA*. 113(38):10607-10612.

913 Moorjani P, Sankararaman S, Fu QM, Przeworski M, Patterson N, Reich D. 2016b. A genetic  
914 method for dating ancient genomes provides a direct estimate of human generation  
915 interval in the last 45,000 years. *Proc Natl Acad Sci USA*. 113(20):5652-5657.

916 Myers S, Bottolo L, Freeman C, McVean G, Donnelly P. 2005. A fine-scale map of  
917 recombination rates and hotspots across the human genome. *Science*. 310(5746):321-  
918 324.

919 Nam K, Munch K, Mailund T, Nater A, Greminger MP, Krutzen M, Marques-Bonet T,  
920 Schierup MH. 2017. Evidence that the rate of strong selective sweeps increases with  
921 population size in the great apes. *Proc Natl Acad Sci USA*. 114(7):1613-1618.

922 Orr HA, Betancourt AJ. 2001. Haldane's sieve and adaptation from the standing genetic  
923 variation. *Genetics*. 157(2):875-884.

924 Patin E, Lopez M, Grollemund R, Verdu P, Harmant C, Quach H, Laval G, Perry GH,  
925 Barreiro LB, Froment A et al. 2017. Dispersals and genetic adaptation of bantu-  
926 speaking populations in africa and north america. *Science*. 356(6337):543-546.

927 Pavlidis P, Alachiotis N. 2017. A survey of methods and tools to detect recent and strong  
928 positive selection. *J Biol Res-Thessalon*. 24:7.

929 Perry GH, Foll M, Grenier JC, Patin E, Nedelec Y, Pacis A, Barakatt M, Gravel S, Zhou X,  
930 Nsohya SL et al. 2014. Adaptive, convergent origins of the pygmy phenotype in  
931 african rainforest hunter-gatherers. *Proc Natl Acad Sci USA*. 111(35):E3596-3603.

932 Peter BM, Huerta-Sanchez E, Nielsen R. 2012. Distinguishing between selective sweeps from  
933 standing variation and from a de novo mutation. *PLoS Genet*. 8(10): e1003011.

934 Pickrell JK, Coop G, Novembre J, Kudaravalli S, Li JZ, Absher D, Srinivasan BS, Barsh GS,  
935 Myers RM, Feldman MW et al. 2009. Signals of recent positive selection in a  
936 worldwide sample of human populations. *Genome Res*. 19(5):826-837.

937 Pool JE. 2015. The mosaic ancestry of the drosophila genetic reference panel and the d.  
938 *Melanogaster* reference genome reveals a network of epistatic fitness interactions. *Mol*  
939 *Biol Evol*. 32(12):3236-3251.

940 Prado-Martinez J, Sudmant PH, Kidd JM, Li H, Kelley JL, Lorente-Galdos B, Veeramah KR,  
941 Woerner AE, O'Connor TD, Santpere G et al. 2013. Great ape genetic diversity and  
942 population history. *Nature*. 499(7459):471-475.

943 Pritchard JK, Di Rienzo A. 2010. Adaptation - not by sweeps alone. *Nat Rev Genet*.  
944 11(10):665-667.

945 Pritchard JK, Pickrell JK, Coop G. 2010. The genetics of human adaptation: Hard sweeps,  
946 soft sweeps, and polygenic adaptation. *Curr Biol*. 20(4):R208-215.

947 Przeworski M. 2002. The signature of positive selection at randomly chosen loci. *Genetics*.  
948 160(3):1179-1189.

949 Przeworski M, Coop G, Wall JD. 2005. The signature of positive selection on standing  
950 genetic variation. *Evolution*. 59(11):2312-2323.



951 Pybus M, Luisi P, Dall'Olio GM, Uzkudun M, Laayouni H, Bertranpetit J, Engelken J. 2015.  
952 Hierarchical boosting: A machine-learning framework to detect and classify hard  
953 selective sweeps in human populations. *Bioinformatics*. 31(24):3946-3952.

954 Refoyo-Martinez A, da Fonseca RR, Halldorsdottir K, Arnason E, Mailund T, Racimo F.  
955 2019. Identifying loci under positive selection in complex population histories.  
956 *Genome Res*. 29(9):1506-1520.

957 Roux PF, Boitard S, Blum Y, Parks B, Montagner A, Mouisel E, Djari A, Esquerre D, Desert  
958 C, Boutin M et al. 2015. Combined qtl and selective sweep mappings with coding snp  
959 annotation and cis-eqtl analysis revealed park2 and jag2 as new candidate genes for  
960 adiposity regulation. *G3-Genes Genomes Genetics*. 5(4):517-529.

961 Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, Cotsapas C, Xie X, Byrne EH,  
962 McCarroll SA, Gaudet R et al. 2007. Genome-wide detection and characterization of  
963 positive selection in human populations. *Nature*. 449(7164):913-918.

964 Sawyer SA, Hartl DL. 1992. Population genetics of polymorphism and divergence. *Genetics*.  
965 132(4):1161-1176.

966 Sawyer SA, Parsch J, Zhang Z, Hartl DL. 2007. Prevalence of positive selection among  
967 nearly neutral amino acid replacements in drosophila. *Proc Natl Acad Sci USA*.  
968 104(16):6504-6510.

969 Schaffner SF, Foo C, Gabriel S, Reich D, Daly MJ, Altshuler D. 2005. Calibrating a  
970 coalescent simulation of human genome sequence variation. *Genome Res*.  
971 15(11):1576-1583.

972 Schmidt JM, de Manuel M, Marques-Bonet T, Castellano S, Andres AM. 2019. The impact of  
973 genetic adaptation on chimpanzee subspecies differentiation. *PLoS Genet*.  
974 15(11):e1008485.

975 Schrider DR, Kern AD. 2017. Soft sweeps are the dominant mode of adaptation in the human  
976 genome. *Mol Biol Evol.* 34(8):1863-1877.

977 Skoglund P, Mathieson I. 2018. Ancient genomics of modern humans: The first decade. *Annu*  
978 *Rev Genomics Hum Genet.* 19:381-404.

979 Smith NG, Eyre-Walker A. 2002. Adaptive protein evolution in drosophila. *Nature.*  
980 415(6875):1022-1024.

981 Stella A, Ajmone-Marsan P, Lazzari B, Boettcher P. 2010. Identification of selection  
982 signatures in cattle breeds selected for dairy production. *Genetics.* 185(4):1451-1461.

983 Stephan W. 2016. Signatures of positive selection: From selective sweeps at individual loci to  
984 subtle allele frequency changes in polygenic adaptation. *Mol Ecol.* 25(1):79-88.

985 Stephan W, Wiehe THE, Lenz MW. 1992. The effect of strongly selected substitutions on  
986 neutral polymorphism: Analytical results based on diffusion theory. *Theor Popul Biol.*  
987 41:237–254.

988 Teshima KM, Coop G, Przeworski M. 2006. How reliable are empirical genomic scans for  
989 selective sweeps? *Genome Res.* 16(6):702-712.

990 Tishkoff SA, Reed FA, Ranciaro A, Voight BF, Babbitt CC, Silverman JS, Powell K,  
991 Mortensen HM, Hirbo JB, Osman M et al. 2007. Convergent adaptation of human  
992 lactase persistence in africa and europe. *Nat Genet.* 39(1):31-40.

993 Uricchio LH, Petrov DA, Enard D. 2019. Exploiting selection at linked sites to infer the rate  
994 and strength of adaptation. *Nat Ecol Evol.* 3(6):977-984.

995 Vitti JJ, Grossman SR, Sabeti PC. 2013. Detecting natural selection in genomic data. *Annu*  
996 *Rev Genet.* 47:97-120.

997 Voight BF, Kudaravalli S, Wen X, Pritchard JK. 2006. A map of recent positive selection in  
998 the human genome. *PLoS Biol.* 4(3):e72.

999 Wiehe TH, Stephan W. 1993. Analysis of a genetic hitchhiking model, and its application to  
1000 DNA polymorphism data from drosophila melanogaster. *Mol Biol Evol.* 10(4):842-  
1001 854.

1002 Zeng K, Fu YX, Shi S, Wu CI. 2006. Statistical tests for detecting positive selection by  
1003 utilizing high-frequency variants. *Genetics.* 174(3):1431-1439.

1004 Zhen Y, Huber CD, Davies RW, Lohmueller KE. 2021. Greater strength of selection and  
1005 higher proportion of beneficial amino acid changing mutations in humans compared  
1006 with mice and drosophila melanogaster. *Genome Res.* 31(1):110-120.

1007

1008

## 1009 **Legends**

### 1010 **Figure 1. The clustering of candidate SNPs of selection in PSVs**

1011 (A) A sweep in a PSV region (yellow) with the target of selection indicated in red and neutral  
1012 nearby candidate SNPs indicated in black for SNPs unaffected by selection and in blue for  
1013 SNPs driven by selection (ENV regions are indicated in gray). (B) Enrichment of candidate  
1014 SNPs around the simulated targets of selection in 5Mb genomic regions simulated using the  
1015 African demographic model and the recombination and selection parameters used in this  
1016 study. The blue line in bold shows the proportions of candidate SNPs observed in non-  
1017 overlapping 100Kb regions averaged across 1,000 simulations. The blue surface shows the  
1018 minimum and maximum values obtained. The same statistics obtained using 5Mb regions  
1019 simulated in absence of selection are indicated in black bold and dashed lines respectively.  
1020 (C) An example of  $X = 5$  selective sweeps (each targeting the SNPs indicated in red) in a  
1021 whole-genome sequence (WGS) dataset, illustrating the expected enrichment of candidate  
1022 SNPs in PSVs.

1023

### 1024 **Figure 2. Relationships between the simulated ORs and $X$ assuming incomplete sweeps** 1025 **only**

1026 ORs obtained with fixed numbers of incomplete sweeps per simulated WGS data in African,  
1027 European and East Asian populations (1,000 simulated WGSs in each case). The frequency at  
1028 the onset of selection varies from  $1/2N$  to 0.2 (SDN and SSV merged together). (A)  
1029 Distributions of  $s$  (selection coefficient per sweep, left hand side) and  $t$  (age of selection per  
1030 sweep, right hand side) obtained after the exclusion of complete sweeps in Africa (yellow), in  
1031 Europe (blue) and in East Asia (green). (B) Simulated ORs for Fay & Wu's  $H$  (F&W-H),  $iHS$ ,  
1032  $DIND$ ,  $\Delta iHH$  and the two pairwise XP-EHHs.

1033

1034 **Figure 3. Accuracy of the ABC estimations of  $X$  in African, European and East Asian**  
1035 **simulated populations**

1036 Accuracy of the estimations of  $X$  under the incomplete sweep scenario (**A**) and under the  
1037 sweep scenario (**B**), assessed by means of simulated WGSs used as pseudo-empirical data,  
1038 each containing 0, 50, 100 or 150 selective sweeps (200 pseudo-empirical data in each  
1039 situation). Each boxplot display the 200 point estimates obtained. The horizontal red lines  
1040 indicate the true values of  $X$ . (**A**) For both pseudo-empirical data and ABC simulations, the  
1041 selective sweeps are incomplete only. The pseudo-empirical data are the simulated WGSs  
1042 shown in Figure 2B. The distributions of  $s$  and  $t$  used to simulate the WGS data are shown in  
1043 Figure 2A. (**B**) For both pseudo-empirical data and ABC simulations, the selective sweeps  
1044 can be either complete or incomplete. The pseudo-empirical data are the simulated WGSs  
1045 shown in Figure S3B. The distribution of  $t$  used to simulate the WGS data is a uniform  
1046 distribution whereas we used the distribution of  $s$  enriched in low values shown in Figure  
1047 S3A.

1048

1049 **Figure 4. Posterior distributions of  $X$  in African, European and East Asian 1000G**  
1050 **populations.**

1051 Posterior distributions of  $X$  obtained for each African (**A**), European (**B**) and East Asian (**C**)  
1052 1000G population analyzed separately (the distributions are shifted for visibility, the  
1053 population names are ranked in order of appearance in the plots, the YRI, CEU and CHB  
1054 populations are indicated in bold). These distributions correspond to the estimations shown in  
1055 Table 1 obtained using the 1000G ORs averaged across chromosomes. The left hand side  
1056 panels show the posterior distributions obtained under the incomplete sweep scenario. The  
1057 ABC simulations used to perform these estimations are those used in Figure 3A (ABC  
1058 simulations with incomplete selective sweeps only). The right hand side panels show the

1059 posterior distributions obtained under the sweep scenario. The ABC simulations used to  
1060 perform these estimations are those used in Figure 3B (ABC simulations with selective  
1061 sweeps that can be either complete or incomplete).

1062

1063 **Figure 5. Numbers of selective sweeps as a function of allele frequencies at the onset of**  
1064 **selection**

1065 Point estimates of  $X_1$ ,  $X_2$  and  $X_3$ , i.e., the numbers of selective sweeps with very low ( $1/2N \leq$   
1066  $p_{start} < 0.01$ ), low ( $0.01 \leq p_{start} < 0.1$ ) and intermediate ( $0.1 \leq p_{start} < 0.2$ ) initial  
1067 frequencies of the selected alleles. The estimations were obtained using the 1000G ORs  
1068 averaged across chromosomes and the ABC simulations used in Figure 4. **(A)** Point estimates  
1069 obtained under the incomplete sweep scenario. **(B)** Point estimates obtained under the sweep  
1070 scenario. The vertical bars show the 95% CIs edges averaged in a given continent (the  
1071 estimations obtained in each 1000G population can be found in Table S2A).

1072 **Table 1. Estimations of the number of selective sweeps.**

Continent <sup>a</sup>	Incomplete sweeps		Sweeps	
	X	X (BGS) <sup>b</sup>	X	X (BGS) <sup>b</sup>
<b>Africa<sup>c</sup></b>	62 [36-91]	71 [43-102]	68 [46-91]	63 [41-86]
<b>Europe<sup>c</sup></b>	71 [35-111]	74 [36-111]	115 [68-160]	137 [89-180]
<b>East Asia<sup>c</sup></b>	88 [50-127]	104 [66-143]	165 [119-211]	158 [109-206]
<b>All<sup>c</sup></b>	74 [41-110]	83 [48-119]	116 [78-154]	119 [80-157]
<b>Africa<sup>d</sup></b>	46 [21-73]	59 [34-87]	58 [37-81]	53 [32-75]
<b>Europe<sup>d</sup></b>	56 [26-90]	46 [17-79]	84 [40-130]	103 [57-148]
<b>East Asia<sup>d</sup></b>	60 [29-93]	73 [41-108]	128 [83-173]	124 [80-170]
<b>All<sup>d</sup></b>	54 [25-85]	59 [31-91]	90 [53-128]	93 [56-131]

1073 <sup>a</sup>Point estimates and 95% CIs edges averaged across populations of the same continental  
1074 origin (95% CIs are indicated in squared brackets). <sup>b</sup>BGS stands for background selection  
1075 simulated in PSVs. <sup>c</sup>Estimations obtained using the 1000G ORs averaged across  
1076 chromosomes. <sup>d</sup>Estimations obtained using the 1000G ORs computed merging chromosomes.

1077

1078 **Table 2. Estimations of the average strength and average age of selection.**

	Incomplete sweeps		Sweeps	
Continent <sup>a</sup>	S	T (Kya)	S	T (Kya)
<b>Africa<sup>b</sup></b>	0.010 [0.005-0.018]	43.9 [31.4-52.4]	0.013 [0.007-0.022]	53.6 [48.3-58.0]
<b>Europe<sup>b</sup></b>	0.017 [0.010-0.029]	28.2 [15.4-38.9]	0.010 [0.007-0.013]	52.0 [49.2-54.7]
<b>East Asia<sup>b</sup></b>	0.012 [0.006-0.022]	34.9 [23.1-44.5]	0.011 [0.007-0.017]	54.7 [51.1-57.9]
<b>All<sup>b</sup></b>	0.013 [0.007-0.023]	35.7 [23.3-45.2]	0.011 [0.007-0.017]	53.4 [49.5-56.9]
<b>Africa<sup>c</sup></b>	0.014 [0.008-0.023]	38.3 [24.2-48.5]	0.014 [0.007-0.024]	52.0 [46.1-56.6]
<b>Europe<sup>c</sup></b>	0.018 [0.010-0.032]	26.4 [12.6-38.5]	0.010 [0.007-0.013]	51.5 [48.4-54.1]
<b>East Asia<sup>c</sup></b>	0.013 [0.006-0.025]	30.0 [17.1-41.2]	0.012 [0.008-0.018]	53.6 [49.8-57.1]
<b>All<sup>c</sup></b>	0.015 [0.008-0.027]	31.6 [17.9-42.8]	0.012 [0.007-0.018]	52.3 [48.1-55.9]

1079 <sup>a</sup>Point estimates and 95% CIs edges averaged across populations of the same continental  
1080 origin (95% CIs are indicated in squared brackets). <sup>b</sup>Estimations obtained using the 1000G  
1081 ORs averaged across chromosomes. <sup>c</sup>Estimations obtained using the 1000G ORs computed  
1082 merging chromosomes).

1083



1084 **Table 3. Comparing with previous estimations of  $\alpha$  in the human lineage.**

<b>Human lineage</b>	<b><math>X^a</math></b>	<b>[95% CI]</b>	<b><math>X_g^b</math></b>	<b><math>\alpha^c</math></b>	<b>[95% CI]</b>
<b>Africa</b>	68	[46-91]	0.019	0.177	[0.120-0.236]
<b>Africa<sup>d</sup></b>	58	[37-81]	0.017	0.151	[0.096-0.211]
<b>Uricchio et al 2019</b>				0.135	[0.096-0.170]
<b>Zhen et al 2021<sup>e</sup></b>				0.060	
<b>Zhen et al 2021<sup>f</sup></b>				0.160	
<b>Human-chimp</b>					
<b>Zhen et al 2021<sup>e</sup></b>				0.110	
<b>Zhen et al 2021<sup>f</sup></b>				0.250	

1085 <sup>a</sup>Point estimates and 95% CIs edges averaged across populations of the same continental  
1086 origin (95% CIs are indicated in squared brackets). <sup>b</sup> $X_g$  stands for  $X$  per generation, i.e.,  $X$   
1087 divided by the 3,500 simulated generations. <sup>c</sup>Predictions using  $X$  per generation, 29,925  
1088 nonsynonymous and 7,9Ky of human-chimp divergence (95% CIs are indicated in squared  
1089 brackets). <sup>d</sup>Estimations obtained using the 1000G ORs computed merging chromosomes  
1090 (Table 1). <sup>e,f</sup>Estimations under the simple and complex models presented in Zhen et al.  
1091 (2021).

1092

1093 **Table 4. Comparing with *Drosophila*.**

<b>Human</b>	$X$	$X_g^a$	$S$	$Ns^b$	$\lambda$	$\lambda S$	$2N\lambda^b$
<b>Africa</b>	68	0.019	0.013	130	6.5E-12 <sup>c</sup>	8.4E-14	1.3E-07
<b>Africa<sup>d</sup></b>	58	0.017	0.014	140	5.5E-12 <sup>c</sup>	7.5E-14	1.1E-07
<b><i>Drosophila</i><sup>e</sup></b>							
<b>Jensen et al 2008</b>	-	-	0.011	27,500	7.9E-12	8.7E-14	3.9E-05
<b>Macpherson et al 2007</b>	-	-	0.010	15,000	3.6E-12	3.6E-14	1.1E-05
<b>Jensen et al 2008</b>	-	-	0.002	4,800	4.2E-11	8.4E-14	2.0E-04
<b>Li and Stephan 2006</b>	-	-	0.002	17,200	1.1E-11	2.2E-14	1.9E-04
<b>Andolfatto 2007</b>	-	-	1.2E-05	23	6.9E-10	8.3E-15	2.6E-03

1094 <sup>a</sup> $X_g$  stands for  $X$  per generation, i.e.,  $X$  divided by the 3,500 simulated generations. <sup>b</sup>Mean  $Ns$

1095 and  $2N\lambda$  computed from the estimations of  $S$  and  $\lambda$  using  $N = 10,000$  as reference for

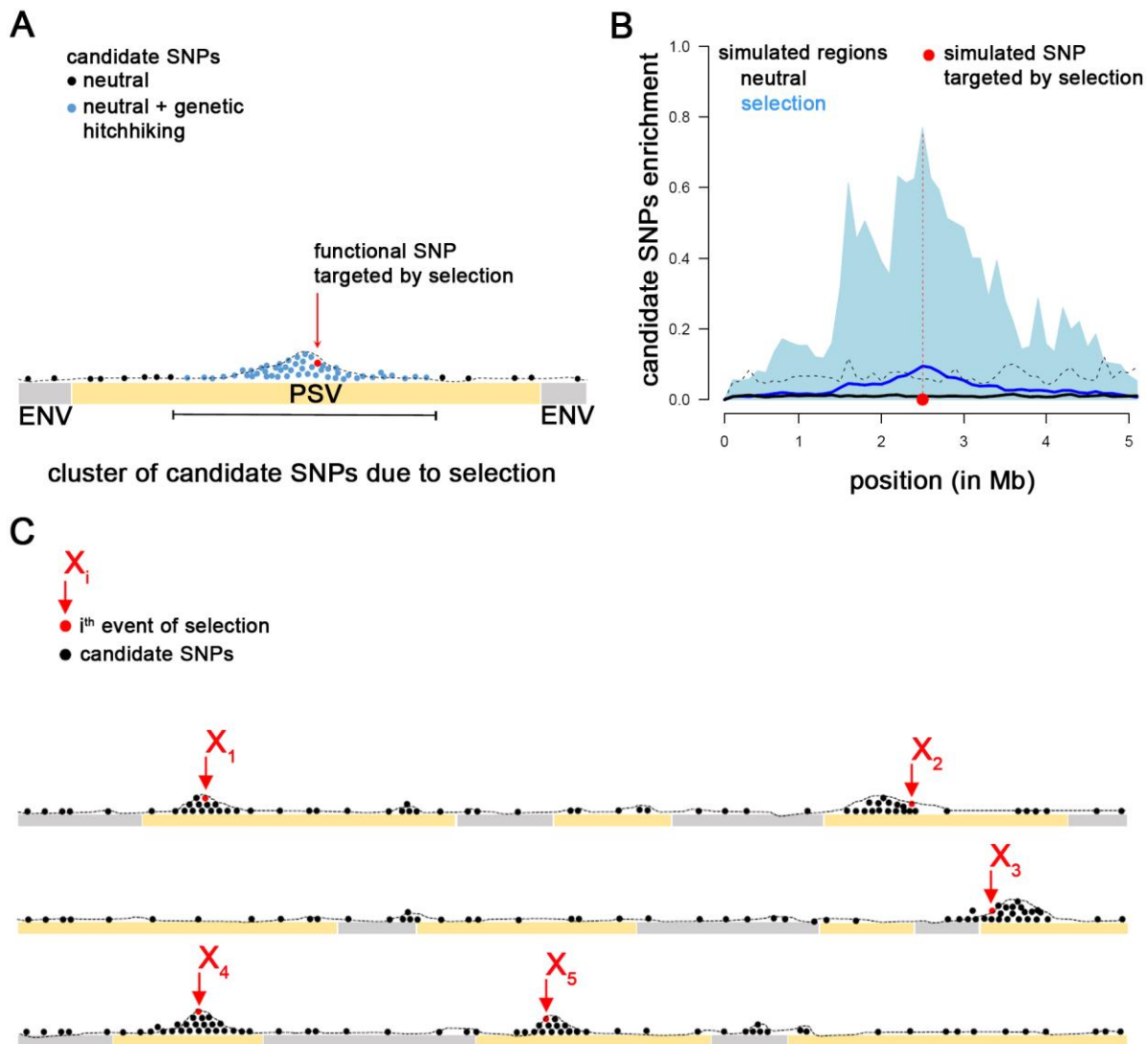
1096 humans effective sizes or, in *Drosophila*, the effective size values drawn from Jensen et al

1097 (2008). <sup>c</sup> $X$  per generation divided by the number of base pairs analyzed. <sup>d</sup>Estimations obtained

1098 using the 1000G ORs computed merging chromosomes (Tables 1 and 2). <sup>e</sup>Values drawn from

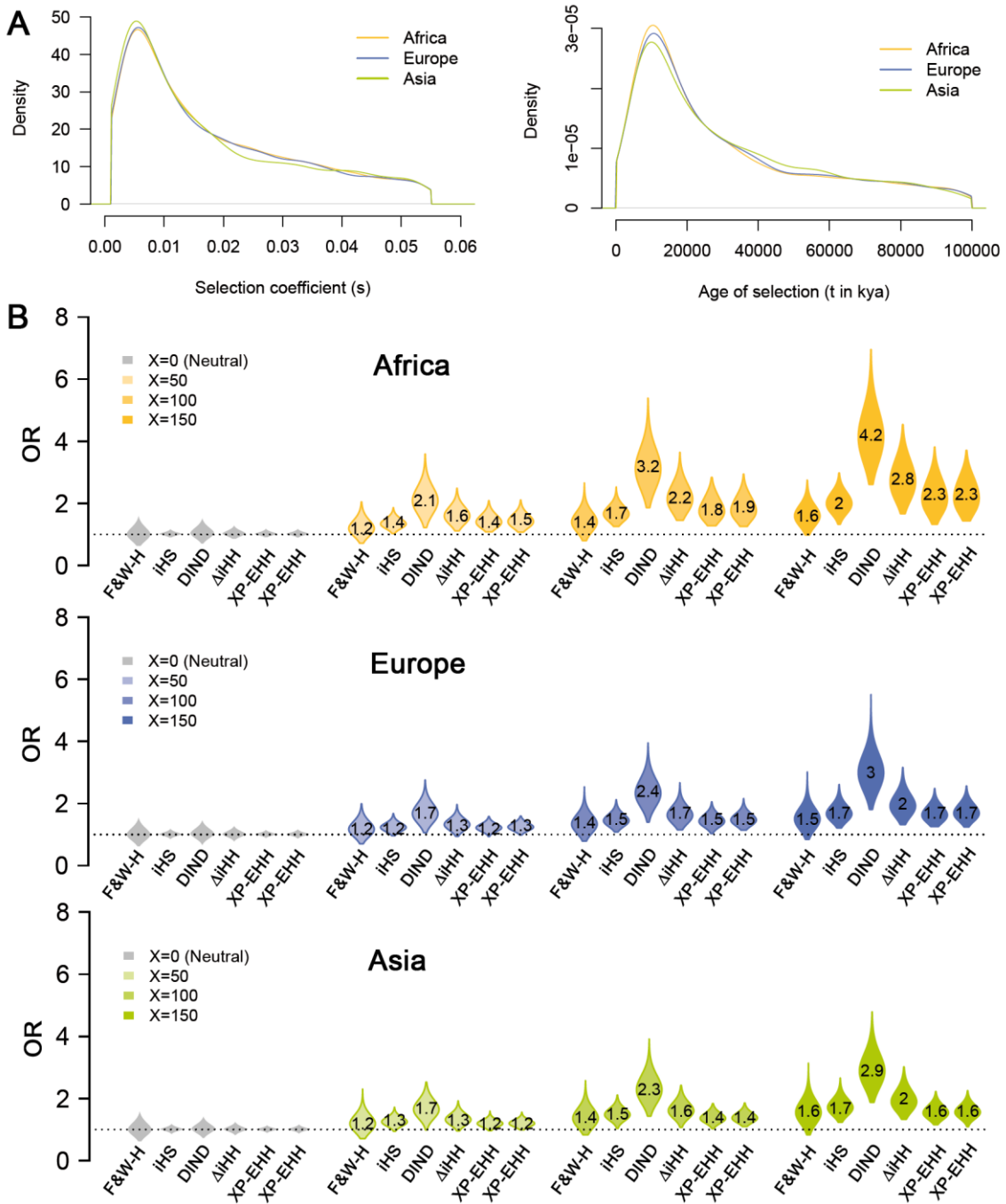
1099 Jensen et al (2008).

1100



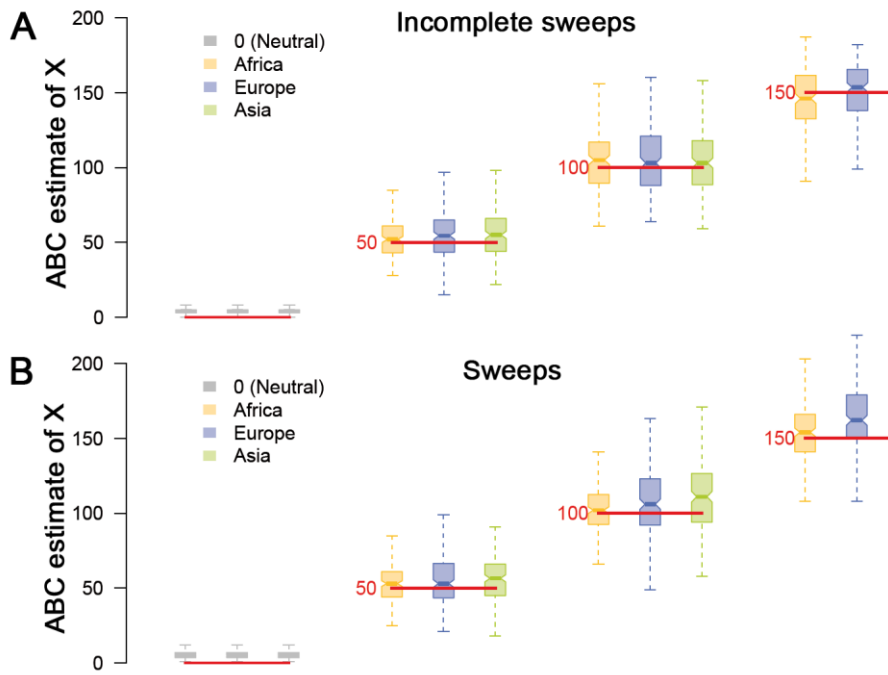
1102

1103



1105

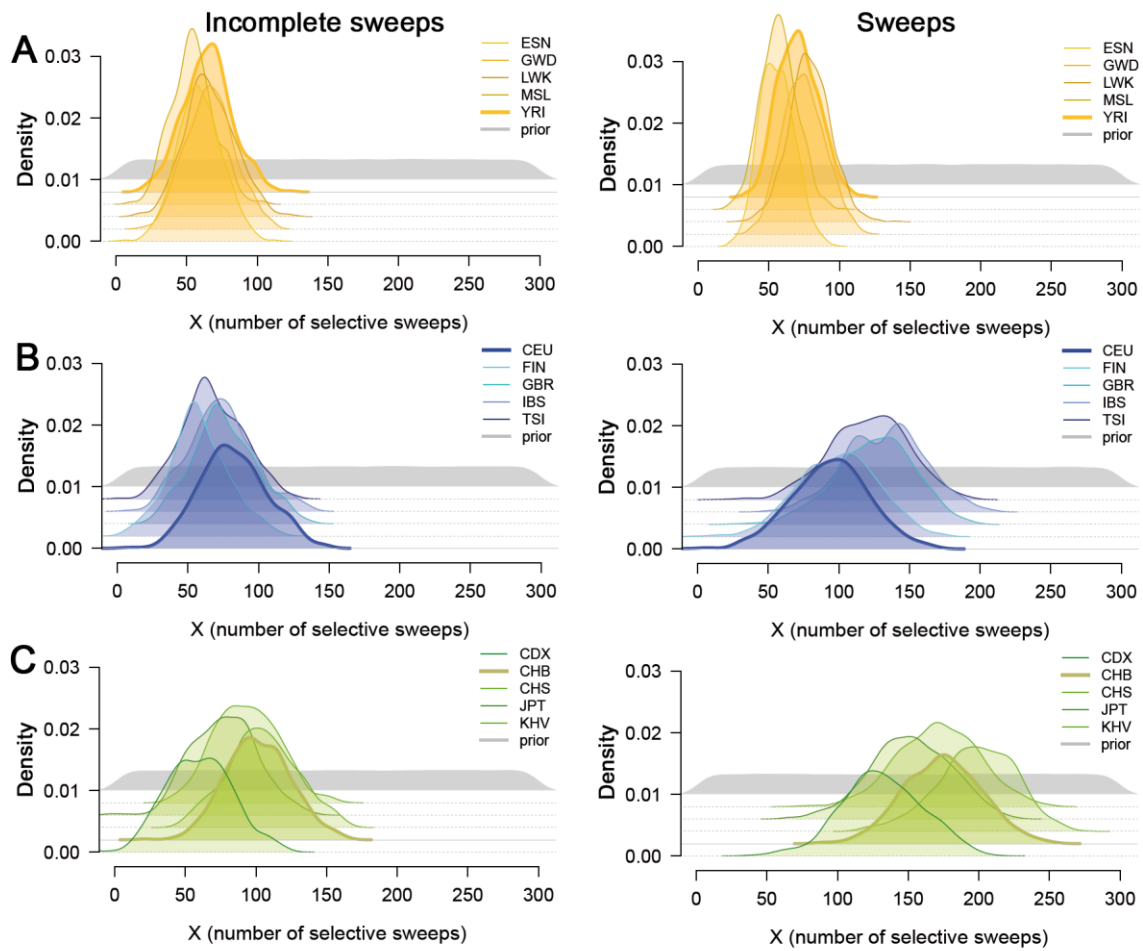
1106



1108

1109

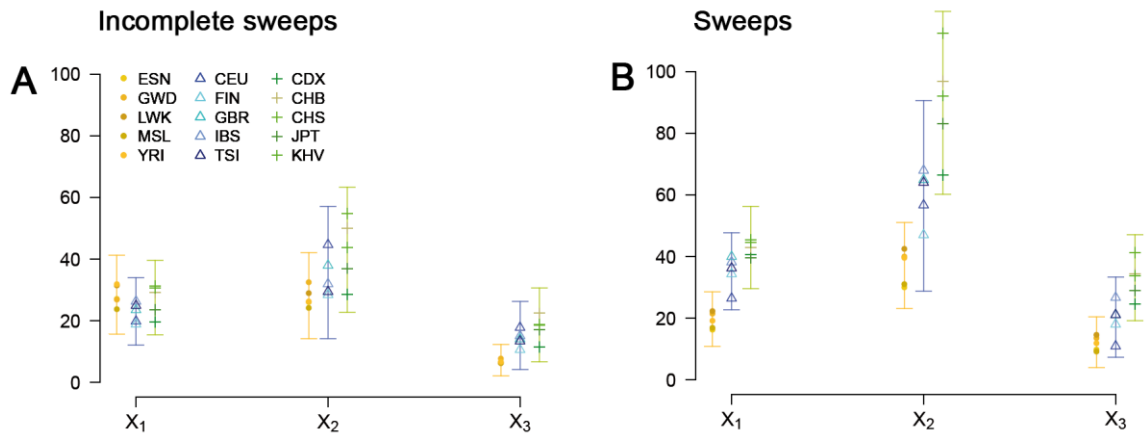
1110 **Figure 4**



1111

1112

1113 **Figure 5**



1114