

# Accuracy of demographic inferences from Site Frequency Spectrum: The case of the Yoruba population

Marguerite Lapierre<sup>\*,†</sup>, Amaury Lambert<sup>†,‡</sup>, and Guillaume Achaz<sup>\*,†</sup>

<sup>\*</sup>Atelier de Bioinformatique, UMR 7205 ISyEB, MNHN-UPMC-CNRS-EPHE, Muséum  
National d'Histoire Naturelle, 75005 Paris, France

<sup>†</sup>SMILE (Stochastic Models for the Inference of Life Evolution), UMR 7241 CIRB, Collège  
de France, CNRS, INSERM, PSL Research University, 75005 Paris, France

<sup>‡</sup>Laboratoire de Probabilités et Modèles Aléatoires (LPMA), UMR 7599, UPMC-CNRS,  
75005 Paris, France

Running title: Accuracy of demographic inferences

Key Words: human demography, model identifiability, coalescent theory, site frequency spectrum

Corresponding Author:

Marguerite Lapierre

Atelier de Bioinformatique

Muséum National d'Histoire Naturelle

Boîte Courrier 50, Bâtiment 139

45 rue Buffon

75005 PARIS

(33)(0)140 79 45 09

`marguerite.lapierre@mnhn.fr`

## Abstract

1  
2 Some methods for demographic inference based on the observed genetic diversity of current  
3 populations rely on the use of summary statistics such as the Site Frequency Spectrum (SFS).  
4 Demographic models can be either model-constrained with numerous parameters such as  
5 growth rates, timing of demographic events and migration rates, or model-flexible, with an  
6 unbounded collection of piecewise constant sizes. It is still debated whether demographic  
7 histories can be accurately inferred based on the SFS. Here we illustrate this theoretical issue  
8 on an example of demographic inference for an African population. The SFS of the Yoruba  
9 population (data from the 1000 Genomes Project) is fit to a simple model of population  
10 growth described with a single parameter (*e.g.*, founding time). We infer a time to the most  
11 recent common ancestor of 1.7 million years for this population. However, we show that the  
12 Yoruba SFS is not informative enough to discriminate between several different models of  
13 growth. We also show that for such simple demographies, the fit of one-parameter models  
14 outperforms the model-flexible method recently developed by Liu and Fu. The use of this  
15 method on simulated data suggests that it is biased by the noise intrinsically present in the  
16 data.

## INTRODUCTION

17

18 Inference of human population history based on demographic models for genomic data can  
19 complement archaeological knowledge, owing to the large amount of polymorphism data now  
20 available in human populations. Polymorphism data can be viewed as an imprint left by  
21 past demographic events on the current genetic diversity of a population (see, *e.g.*, review  
22 by POOL *et al.* 2010).

23 There are several means of analyzing this observed genetic diversity for demographic  
24 inference. The polymorphism data can be used to reconstruct a coalescence tree of the sam-  
25 pled individuals. The demography of the sampled population can be inferred by comparing  
26 this reconstructed tree with theoretical predictions under a constant size model (PYBUS  
27 *et al.* 2000). For example, in an expanding population, the reconstructed coalescent tree will  
28 have relatively longer terminal branches than the reference coalescent tree in a population  
29 of constant size. However, methods based on a single reconstructed tree are flawed because  
30 of recombination (LAPIERRE *et al.* 2016), since the genealogy of a recombining genome is  
31 described by as many trees as there are recombining loci.

32 The genome-wide distribution of allele frequencies is a function of the average genealogies,  
33 and can thus be used as a summary statistic for demographic inference. This distribution,  
34 called the Site Frequency Spectrum (SFS), reports the number of mutated sites at any  
35 given frequency. The demographic history of a population affects the shape of its SFS  
36 (ADAMS and HUDSON 2004; MARTH *et al.* 2004). For example, an expanding population  
37 carries an excess of low-frequency variants, compared with the expectation under a constant  
38 size model. The shape of the SFS is also altered by selection, which results in an excess  
39 of low- and high-frequency variants (FAY and WU 2000). However, selection acts mainly  
40 on the coding parts of the genome and the non-coding segments linked to them, while  
41 demography impacts the whole genome. Furthermore, unlike reconstructed trees, the SFS is  
42 not biased by recombination (WALL 1999). Quite on the contrary, by averaging the SFS over  
43 many correlated marginal genealogies, recombination lowers the variance of the SFS while

44 its expectation remains unchanged. Therefore, the SFS of a sample is a summary of the  
45 genetic diversity, averaged over all the genome due to recombination, that can be analyzed  
46 in terms of demography.

47 Several types of methods exist to infer the demography of a population based on its SFS. A  
48 specific demographic model can be tested by computing a pseudo-likelihood function for this  
49 model, based on the comparison of the observed SFS and the SFS estimated by Monte Carlo  
50 coalescent tree simulations (NIELSEN 2000; COVENTRY *et al.* 2010; NELSON *et al.* 2012).  
51 This method can be extended to infer demographic scenarios of several populations, using  
52 their joint SFS (EXCOFFIER *et al.* 2013). Methods based on Monte Carlo tree simulations  
53 are typically very costly in computation time. Other approaches rely on diffusion processes:  
54 they use the solution to the partial differential equation of the density of segregating sites  
55 as a function of time (GUTENKUNST *et al.* 2009; LUKIĆ *et al.* 2011).

56 Whereas all these methods are model-constrained, *i.e.*, they use the SFS to test the like-  
57 lihood of a given demographic model, more flexible methods have been developed. Recently,  
58 BHASKAR *et al.* (2015) derived exact expressions of the expected SFS for piecewise-constant  
59 and piecewise-exponential demographic models. LIU and FU (2015) developed a model-  
60 flexible method based on the SFS: the stairway plot. This method infers the piecewise-  
61 constant demography which maximizes the composite likelihood of the SFS, without any  
62 previous knowledge on the demography. This optimization is based on the estimation of a  
63 time-dependent population mutation rate  $\theta$ . Although they show that their method infers  
64 efficiently some theoretical demographies, they do not test the goodness of fit of the ex-  
65 pected SFS, reconstructed under the demography they infer, with the input SFS on which  
66 they apply their method.

67 All these methods are widely used for the inference of demography in humans and other  
68 species, but doubts remain on the identifiability of a population demography based on its  
69 SFS. It has been shown theoretically that certain population size functions are unidentifiable  
70 from the population SFS (MYERS *et al.* 2008; TERHORST and SONG 2015). MYERS *et al.*

71 (2008) showed that for any given population size function  $N(t)$ , there exists an infinite  
72 number of smooth functions  $F(t)$  such that  $\xi^N = \xi^{N+F}$  where  $\xi^N$  is the SFS of a population  
73 of size function  $N(t)$ . However, other theoretical works have recently shown that for many  
74 types of population size functions commonly used in demography studies, such as piecewise  
75 constant or piecewise exponential functions, demography can be inferred based on the SFS,  
76 provided the sample is large enough (BHASKAR and SONG 2014). These studies argued  
77 that the unidentifiability proven by MYERS *et al.* (2008) relied on biologically unrealistic  
78 population size functions involving high frequency oscillations near the present. Lately, two  
79 studies (KIM *et al.* 2015; TERHORST and SONG 2015) have provided bounds on the amount  
80 of demographic information contained in the SFS or in coalescent times.

81 In this study, we use the SFS of an African population (the Yoruba population, data from  
82 THE 1000 GENOMES PROJECT CONSORTIUM 2015) as an example of a somewhat simple  
83 demography, to illustrate the risks of over-confidence in demographic scenarios inferred.  
84 Namely, we highlight two issues potentially arising even in the case of simple demogra-  
85 phies: unidentifiability of models and poor goodness of fit of inferences. We first infer the  
86 Yoruba demography with a model-constrained method, using diverse one-parameter models  
87 of growth, and then with a model-flexible method, the stairway plot (LIU and FU 2015).  
88 For the model-constrained method, we test four different growth models derived from the  
89 standard neutral framework used in the vast majority of population genetics studies, also  
90 compared with a more uncommon type of model based on a branching process. Individual-  
91 based models such as the branching process are widely used in population ecology (LAMBERT  
92 2010): the population is modeled as individuals which die and give birth at given rates in-  
93 dependently. These models are not commonly used in population genetics although they  
94 provide interesting features of fluctuating population sizes for example, and benefit from a  
95 strong mathematical framework.

## MATERIALS AND METHODS

96

97 **1 000 Genomes Project data:** Variant calls from the 1 000 Genomes Project phase 3  
98 were downloaded from the project ftp site (THE 1000 GENOMES PROJECT CONSORTIUM  
99 2015). The sample size for the Yoruba population is  $n = 108$  individuals (polymorphism  
100 data available for both genome copies of each individual, *i.e.*,  $2n = 216$  sequences). We  
101 kept all single nucleotide bi-allelic variants to plot the sample SFS. The number of bi-allelic  
102 sites is  $S = 20\,417\,698$ . The average distance between two sites is 136 bp (median 81 bp).  
103 The number of sites for which the ancestral allele is known is  $S' = 19\,441\,528$ . To avoid  
104 possible bias due to sequencing errors, we ignored singletons (mutations appearing in only  
105 one chromosome of one individual in the sample) for the rest of the study. The implications  
106 of ignoring singletons are examined in the discussion.

107 **Site Frequency Spectrum definition and graphical representation:** The Site Fre-  
108 quency Spectrum (SFS) of a sample of  $n$  diploid individuals is described as the vector  
109  $\xi = (\xi_1, \xi_2, \dots, \xi_{2n-1})$  where for  $i \in [1, 2n - 1]$ ,  $\xi_i$  is the number of dimorphic (*i.e.*, with ex-  
110 actly two alleles) sites with derived form at frequency  $i/2n$ . To avoid potential orientation  
111 errors, we assumed that the ancestral form is unknown for all sites: we worked with a folded  
112 spectrum, where we consider the frequency of the less frequent (or minor) allele. In this  
113 case, the folded SFS is described as the vector  $\eta = (\eta_1, \eta_2, \dots, \eta_n)$  where  $\eta_i = \xi_i + \xi_{2n-i}$  for  
114  $i \in [1, n - 1]$  and  $\eta_n = \xi_n$ . The folded SFS of the Yoruba sample is plotted in Figure S1. For  
115 a better graphical representation, all SFS were transformed as follows: we plot  $\phi_i$  normalized  
116 by its sum, where

- 117 • for unfolded SFS,  $\phi_i = i \xi_i$  for  $i \in [1, 2n - 1]$
- 118 • for folded SFS,  $\phi_i = \eta_i \frac{i(2n-i)}{2n}$  for  $i \in [1, n - 1]$  and  $\phi_n = n \eta_n$

119 The transformed SFS has a flat expectation (*i.e.*, constant over all values of  $i$ ) under the  
120 standard neutral model (NAWA and TAJIMA 2008; ACHAZ 2009).

121 **Demographic models used for the model-constrained methods:** We inferred the  
122 demography of the Yoruba population using five growth models (Figure 1), compared with  
123 the predictions of the standard model with constant population size. Time is measured in  
124 coalescent units of  $2N$  generations, where the scaling parameter  $N$  has the same dimension  
125 as the current population size, which we will not estimate. Time starts at 0 (present time)  
126 and increases backward in time. Four models are based on the standard Kingman coalescent  
127 (KINGMAN 1982), amended with demography. Three of them are described with an explicit  
128 demography: either *Linear* growth since time  $\tau$ , *Exponential* growth at rate  $1/\tau$  or *Sudden*  
129 growth from a single ancestor to the entire population at time  $\tau$ . We also use another model  
130 based on the Kingman coalescent, with an implicit demography: the *Conditioned* model.  
131 This model is based on a standard constant size model, but the Time to the Most Recent  
132 Common Ancestor ( $T_{MRC A}$ ) is conditioned on being reached before time  $\tau$ . The fifth model,  
133 *Birth-Death*, is not based on the standard Kingman coalescent, but on a critical branching  
134 process measured in units of  $2N$  generations. Forward in time, the process starts with a  
135 founding event of one individual. Individuals give birth and die at equal rate 1. The process  
136 is conditioned on not becoming extinct before a period of time  $\tau$ , and on reaching on average  
137  $2N$  individuals.

138 **Stairway plot inference on the Yoruba SFS:** We applied the model-flexible stairway  
139 plot method developed by LIU and FU (2015) on the unfolded Yoruba SFS. Inferences are  
140 made on 200 SFS as suggested by their method. We use the script they provide to create  
141 199 bootstrap samples of the Yoruba SFS. We also ignore the singletons for this method,  
142 and use the default parameter values suggested in their paper for the optimization.

143 **SFS simulation with demography:** We used two different method to simulate SFS under  
144 the four demographic models derived from the Kingman coalescent (*Linear*, *Exponential*,  
145 *Sudden* and *Conditioned*) or under a piecewise-constant demography reconstructed by the



146 stairway plot method.

147 *Method 1:* Simulate  $l$  independent topologies under the Kingman coalescent on which mu-  
148 tations are placed at rate  $\theta$  (population mutation rate) (HUDSON *et al.* 1990). This allows  
149 to simulate the SFS of  $l$  independent loci.

150 *Method 2:* Another way to simulate SFS is using the following formula:

$$\mathbb{E}[\xi_i] = \frac{\theta}{2} \sum_{k=2}^{2n-i+1} k \mathbb{E}[t_k] \mathbb{P}(k, i) \quad (1)$$

151 where  $\theta$  is the population mutation rate,  $t_k$  is the time during which there are  $k$  lines in the  
152 tree (hereafter named state  $k$ ) and  $\mathbb{P}(k, i)$  is the probability that a randomly chosen line at  
153 state  $k$  gives  $i$  descendants in the sample of size  $2n$  (*i.e.*, at state  $2n$ ) (FU 1995). For all  
154 models, the neutrality assumption ensures that

$$\mathbb{P}(k, i) = \frac{\binom{2n-i-1}{k-2}}{\binom{2n-1}{k-1}}$$

155 for  $i \in [1, 2n - 1]$  and  $k \in [2, 2n - i + 1]$ . Using this probability allows to average over the  
156 space of topologies. This reduces considerably computation time since the space of topologies  
157 is very large, and produces smooth SFS for which only the  $t_k$  need to be simulated to obtain  
158 the expectations  $\mathbb{E}[t_k]$ .

159 The expectations  $\mathbb{E}[t_k]$  are obtained as follow: for  $k \in [2, 2n]$ , times in the standard  
160 coalescent  $t_k^*$  are drawn from an exponential distribution of parameter  $\binom{k}{2}$ . For the *Linear*  
161 and *Exponential* models, and for the piecewise-constant demographies reconstructed by the  
162 stairway plot method, these times are then rescaled to take into account the given explicit  
163 demography (see, *e.g.*, HEIN *et al.* 2004, chap.4). For the *Sudden* model, we assume the  
164 coalescence of all lineages at time  $\tau$  if the common ancestor has not been reached yet.  
165 For the *Conditioned* model, we keep only simulations for which  $\sum_{k=2}^{2n} t_k^* \leq \tau$  where  $\tau$  is the  
166 model parameter. The expectations  $\mathbb{E}[t_k]$  are obtained by averaging over  $10^7$  simulations.

167 Alternatively, the expectations  $\mathbb{E}[t_k]$  could also be obtained with analytic formulae provided  
168 by POLANSKI and KIMMEL (2003).

169 For the *Birth-Death* model, we use the explicit formula for the SFS given in DELAPORTE  
170 *et al.* (2016).

171 We normalize the SFS computed under all these models so that their sum equals 1. This  
172 normalization removes the dependence on the mutation rate parameter  $\theta$ . Consequently, the  
173 standard model has no parameters while all others have exactly one ( $\tau$ ).

174 **Optimization of the parameter  $\tau$ :** For each demographic model, we optimize the pa-  
175 rameter  $\tau$  by minimizing the weighted square distance  $d^2$  between the observed SFS of the  
176 Yoruba population and the predicted SFS under the model (simulated with *Method 2*). Both  
177 SFS are normalized for comparison. The distance is computed for all  $\tau$  values in a given  
178 interval (no specific optimization method was used to find the minimum). With  $\tilde{\eta}^{model}$  and  
179  $\tilde{\eta}^{obs}$  the folded and normalized SFS in the tested model and in the data respectively,

$$d^2(\tilde{\eta}^{model}, \tilde{\eta}^{obs}) = \sum_{i=2}^n \frac{(\tilde{\eta}_i^{model} - \tilde{\eta}_i^{obs})^2}{\tilde{\eta}_i^{model}}$$

180 The sum starts at  $i = 2$  because we ignore  $\tilde{\eta}_1^{obs}$ , corresponding to singletons. To calculate  
181 the distance  $d^{2'}$  between the SFS predicted by two models A and B, we weight the terms by  
182 the mean of the two models:

$$d^{2'}(\tilde{\eta}^A, \tilde{\eta}^B) = \sum_{i=2}^n \frac{(\tilde{\eta}_i^A - \tilde{\eta}_i^B)^2}{(\tilde{\eta}_i^A + \tilde{\eta}_i^B)/2}$$

183 **Inference of the Yoruba demography with  $\partial a \partial i$ :** We inferred the demography of the  
184 Yoruba population with the software  $\partial a \partial i$  v1.7 (GUTENKUNST *et al.* 2009), testing the three  
185 models of explicit demography (*Linear*, *Exponential* and *Sudden*). The demographic models  
186 were specified so that the only parameter to optimize is  $\tau$  like for the distance-based inference  
187 method. Singletons were masked and the method was applied on the folded Yoruba SFS.  
188 Details on the demographic functions and parameter values used for the optimization in  $\partial a \partial i$

189 are provided in the Supplementary Methods. We ran the method 100 times for each model  
190 and kept the parameter value with the best maximum log composite likelihood over the 100  
191 runs. In Figure S4, we plot the best log composite likelihood of the 100 runs.

192 **Scaling of the coalescent time:** Optimized values of the parameter  $\hat{\tau}$  for each model are  
193 expressed in coalescent time units, *i.e.*, scaled in  $2N(0)$  generations. As the model population  
194 size at time zero,  $2N(0)$ , is unknown, to scale these coalescent time units in numbers of  
195 generations and consequently in years, we used the expected number of mutations per site  
196  $M$ . From the dataset, we have  $M^{obs} = S/L$  where  $S$  is the number of single nucleotide  
197 mutations (a  $k$ -allelic SNP accounts for  $k - 1$  mutations) and  $L$  is the length of the accessible  
198 sequenced genome in the 1 000 genomes project (90% of the total genome length, THE 1000  
199 GENOMES PROJECT CONSORTIUM). For the theoretical value, we get that  $M^{theo} = \mu \hat{T}_{tot} C$ ,  
200 where we know the mutation rate  $\mu$  from the literature and the total tree length expressed  
201 in coalescent time units  $\hat{T}_{tot}$  from the SFS simulations. Here  $C$  is the coalescent factor, that  
202 is the number of generations per coalescent time unit, also corresponding to  $2N_e(0)$  where  
203  $N_e(0)$  is the effective population size at present time. The total number of generations in  
204 the tree is  $\hat{T}_{tot} C$  from which we derive the total number of mutations per site  $M^{theo}$ . Thus,  
205 using the observed value  $M^{obs}$ , we can estimate  $C$  by  $S/(\mu L \hat{T}_{tot})$ . We assumed a mutation  
206 rate of  $1.2 \times 10^{-8}$  per base pair per generation (CONRAD *et al.* 2011; CAMPBELL *et al.* 2012;  
207 KONG *et al.* 2012). With the coalescent factor  $C$ , we can then convert a coalescent time  
208 unit into a number of generations, or into a number of years assuming 24 years as generation  
209 time (SCALLY and DURBIN 2012).

210 **Graphical representation of the inferred demographies:** To represent the inferred  
211 explicit demographies (models *Linear*, *Exponential* and *Sudden*), we plot the shape of the  
212 demography with the optimized value  $\hat{\tau}$  for each model. For the implicit demographies  
213 (models *Conditioned* and *Birth-Death*), as there is no explicit demographic shape, we plot

214 the mean trajectory of fixation of a new allele in the population: forward in time, these  
215 fixation trajectories illustrate the expansion of the descendance of the sample's ancestor in  
216 the population (see the Supplementary Methods for details).

217 **Comparing the model-constrained and model-flexible methods to infer *Linear***  
218 **growth:** We applied both methods (the one-parameter inference method and the stairway  
219 plot method) on SFS simulated under *Linear* growth. To test the stairway plot method  
220 on a *Linear* growth demography, we simulate 200 independent SFS using *Method 1*, with  
221 sample size  $2n = 216$ ,  $\theta = 100$  (arbitrary value removed by normalization) and a founding  
222 time  $\tau = 2.48$  (estimated for the Yoruba population, see Results). The SFS are simulated  
223 with either  $10^3$ ,  $10^4$  or  $10^5$  independent loci. We scaled the simulated SFS to obtain a total  
224 number of  $S = 20\,417\,698$  variants, so that the total number of variants in the simulated SFS  
225 is the same as in the Yoruba SFS. We ran the stairway plot method on these 200 independent  
226 SFS with the default parameter values suggested in the method, and with the same mutation  
227 rate ( $1.2 \times 10^{-8}$  per base pair per generation) and generation time (24 years) as in our study.  
228 We report the median demography of these 200 independent inferences.

229 To test the one-parameter inference method on these SFS simulated under the *Linear*  
230 model, we run the parameter optimization on a SFS simulated with either  $10^3$ ,  $10^4$ ,  $10^5$  or  
231  $10^6$  loci. The search of the parameter value that minimizes the distance  $d^2$  was optimized  
232 with a Newton-Raphson algorithm. Derivatives were calculated at  $\tau \pm 0.05$  where  $\tau$  is the  
233 parameter value being optimized. The optimization stopped when the optimization step of  
234 the parameter value was smaller than  $10^{-3}$ .

235 **Data and software availability** The 1000 genomes project data used in this study is  
236 publicly available at <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/>.  
237 The code in Python and C written for the study is available at [https://github.com/](https://github.com/lapierreM/Yoruba_demography)  
238 [lapierreM/Yoruba\\_demography](https://github.com/lapierreM/Yoruba_demography). The code in C used for the *Method 1* of SFS simulation is

239 available upon request to G. ACHAZ.

## 240 RESULTS

241 We inferred the demography of the Yoruba population (Africa), from the whole-genome  
242 polymorphism data of 108 individuals (data from the 1000 Genomes Project, THE 1000  
243 GENOMES PROJECT CONSORTIUM), with SFS-based methods, either model-constrained or  
244 model-flexible.

245 It has been shown that human populations have been growing since their emergence  
246 in Africa, and that African populations were supposedly not affected by the Out-of-Africa  
247 bottleneck described for Eurasian populations (MARTH *et al.* 2004; ATKINSON *et al.* 2008;  
248 GUTENKUNST *et al.* 2009; GRONAU *et al.* 2011; TENNESSEN *et al.* 2012). Analyses using  
249 the PSMC method (LI and DURBIN 2011) have shown a reduction of the African popula-  
250 tion size after the divergence with non-African populations. However, MAZET *et al.* (2016)  
251 have recently shown that these analyses could be biased by population structure. Based on  
252 this previous knowledge, for the model-constrained method, we chose to infer the Yoruba  
253 demography with simple models of growth, *i.e.*, with only one phase of growth characterized  
254 by a single parameter. These five models are: *Linear*, *Exponential* or *Sudden* growth, a  
255 *Conditioned* model where the  $T_{MRCA}$  is conditioned on being smaller than the given param-  
256 eter, and a critical *Birth-Death* model based on a branching process (Figure 1). To infer the  
257 Yoruba demography, we fit the SFS predicted under each model with the observed Yoruba  
258 SFS (all SFS are folded). The SFS were normalized to remove the population mutation  
259 rate parameter  $\theta$ , so that each model is characterized by one single parameter  $\tau$  which has  
260 the dimension of a time duration. We fit this parameter by least-square distance between  
261 the observed SFS and the predicted SFS, and by maximum likelihood using the *∂a∂i* soft-  
262 ware (GUTENKUNST *et al.* 2009). For the model-flexible inference, we used the stairway  
263 plot method developed recently by LIU and FU (2015), which infers a piecewise-constant  
264 demography based on the SFS. For this method, the number of parameters to be estimated

265 is determined by a likelihood-ratio test. It can range from 1 to  $2n - 1$  where  $2n$  is the number  
266 of sequences in the sample.

267 The Yoruba SFS was constructed by taking into account the entire genome. Removing  
268 the coding parts of the genome to avoid potential bias due to selection does not affect the  
269 shape of the SFS substantially (Figure S2), since the coding parts represent a very small  
270 fraction of the human genome. The first bin of the observed SFS, accounting for mutations  
271 found in one chromosome of one individual in the sample (black dot in the observed SFS  
272 in Figure 3B), seemed to lie outside the rest of the distribution. This could be due to  
273 sequencing errors being considered as singletons (ACHAZ 2008), and thus we chose to ignore  
274 this value for the model optimization. We have also made sure that the SFS shape was not  
275 affected greatly by the sample size. We compared the SFS of a subsample of half the Yoruba  
276 individuals ( $2n = 108$ ) with the full sample SFS ( $2n = 216$ ) (Figure S3). This shows that  
277 the only bin of the SFS which is significantly affected by this subsampling is the first one,  
278 containing the singletons. As we ignore it in our study, it does not influence our results.

279 The analysis of the Yoruba SFS with the stairway plot method results in a complex  
280 demography with several bottlenecks in the last 160 000 years (Figure 2). The current  
281 effective population size  $N_e(0)$  is 28 500 (time 0 does not correspond to present time as we  
282 ignored singletons, see discussion). The demographic history earlier than 160 000 years ago  
283 shows spurious patterns that should not be interpreted, according to LIU and FU (2015).

284 The inference of the Yoruba demography with one-parameter models was done by min-  
285 imizing the distance between observed and predicted SFS. This gave an optimized value  $\hat{\tau}$   
286 of the parameter  $\tau$  (Figure 3A and Table 1) (with  $\hat{\tau}$  in coalescent units, *Linear*:  $\hat{\tau} = 2.48$ ,  
287 *Exponential*:  $\hat{\tau} = 1.82$ , *Sudden*:  $\hat{\tau} = 1.36$ , *Conditioned*:  $\hat{\tau} = 1.89$ , *Birth-Death*:  $\hat{\tau} = 2.28$ ).  
288 Plotting the predicted SFS with the optimized parameter value  $\hat{\tau}$  confirmed their goodness  
289 of fit with the observed Yoruba SFS (Figure 3B). Compared to the standard model with-  
290 out demography, the addition of just one parameter allows for a surprisingly good fit of  
291 the observed Yoruba SFS. The Yoruba demography thus seems to be compatible with a

292 simple scenario of growth. On the other hand, the demography inferred by the stairway  
293 plot predicts a SFS which does not fit well the observed Yoruba SFS: the distance between  
294 the observed Yoruba SFS and the expected SFS under the stairway plot demography is ten  
295 times the distance between any of the one-parameter model SFS and the data (Figure 3B  
296 and Table 1).

297 The best fitting SFS under each of the five demographic models all have a square dis-  
298 tance  $d^2$  of the order of  $10^{-4}$  with the observed Yoruba SFS (Figure 3A and Table 1) and  
299 have highly similar shapes (Figure 3B). This suggests that the five demographic models used  
300 to infer the demography of the Yoruba are hard to distinguish based only on the observed  
301 SFS. To validate the use of a least square distance to find the best fitting SFS, we also  
302 inferred the Yoruba demography using the  $\partial a \partial i$  software. This model-constrained method  
303 based on the SFS uses a diffusion approximation to simulate SFS and a likelihood framework  
304 for the parameter optimization. We tested the three models of explicit demography (*Linear*,  
305 *Exponential* and *Sudden* growth) parametrized in the same way as in our method. The best  
306 parameter values found by  $\partial a \partial i$  by maximum log composite likelihood are the same as by  
307 our method (with  $\hat{\tau}$  in coalescent units, *Linear*:  $\hat{\tau} = 2.48$ , *Exponential*:  $\hat{\tau} = 1.82$ , *Sudden*:  
308  $\hat{\tau} = 1.36$ ). Moreover, the log composite likelihoods of the best fitting SFS under each model  
309 are on the same scale (the likelihoods are directly comparable because the number of param-  
310 eters is the same for each model) : *Linear*:  $\ln(L) = -3107$ , *Exponential*:  $\ln(L) = -3953$ ,  
311 *Sudden*:  $\ln(L) = -3393$  (Figure S4). They rank the explicit demography models in the same  
312 order as the least square distance  $d^2$  would rank them: the best model is *Linear* growth,  
313 then *Sudden* and finally *Exponential* growth.

314 We computed the expected  $T_{MRCA}$  based on the predicted SFS using (1): as the SFS  
315 predicted under each model are very similar, it means that they have roughly the same  
316 estimated time durations  $t_k$  while there are  $k$  branches in the coalescent tree of the Yoruba  
317 sample. From these expected  $t_k$  we can compute  $T_{MRCA} = \sum_{i=2}^{2n} t_k$ . This is the  $T_{MRCA}$  of the  
318 sample, but we can assume that it is the same as the  $T_{MRCA}$  of the population, because with

319 such a large sample size, the probability that the  $T_{MRCA}$  of the population is different from  
320 the  $T_{MRCA}$  of the sample becomes very small. Under each of four models (excluding the  
321 *Birth-Death* model for which there is no obvious common time scaling), the expected  $T_{MRCA}$   
322 for the Yoruba population is 1.3 in coalescent units. By using the number of mutations  
323 per site in the data and the total tree length inferred from the simulations, we scaled back  
324 this  $T_{MRCA}$  in number of generations and in years, assuming a mutation rate of  $1.2 \times 10^{-8}$   
325 per base pair per generation (CONRAD *et al.* 2011; CAMPBELL *et al.* 2012; KONG *et al.*  
326 2012) and a generation time of 24 years (SCALLY and DURBIN 2012) (see Methods). The  
327  $T_{MRCA}$  of the Yoruba population inferred under the four demographic models is of 87 100  
328 generations corresponding to 1.7 million years. The inferred demographic models, with  
329 scaling in coalescent units, number of generations and number of years, are shown in Figure  
330 4. The coalescent unit of 67 000 estimated to scale the inferred coalescent times in number  
331 of years corresponds to a present effective population size  $N_e(0)$  of 33 500.

332 The demography inferred by the stairway plot method for the Yoruba population is a  
333 piecewise-constant demography showing much more complex patterns of growth and bottle-  
334 necks than the one-parameter models (Figure 2). Moreover, the expected SFS under this  
335 inferred demography does not fit well the observed Yoruba SFS (Figure 3B). To understand  
336 what could produce such a complex demography, we simulated SFS under a *Linear* growth  
337 with the founding time  $\hat{\tau} = 2.48$  inferred for the Yoruba population. We simulated three sets  
338 of 200 SFS, with respectively  $10^3$ ,  $10^4$ , and  $10^5$  loci, to obtain SFS with more or less noise  
339 (solid lines on Figure 5A). We applied the two inference methods to these SFS. The median  
340 demographies inferred by the stairway plot method are strongly affected by the noise of the  
341 SFS, as shown on Figure 5B. When the number of simulated loci is very large (median of 200  
342 independent demographies inferred with  $10^6$  loci), the stairway plot gives a good approxima-  
343 tion of the true demography, and the expected SFS under the inferred demography fits the  
344 input SFS. However, for smaller numbers of loci (median of 200 independent demographies  
345 inferred with  $10^5$  loci or less), the stairway plot shows complex patterns of growth and bot-



346 tlenecks incompatible with the true demography, and the expected SFS under the inferred  
347 demographies do not fit the input SFS. On the contrary, the one-parameter method infers a  
348 *Linear* demography with a founding time close to the true value for SFS simulated with  $10^4$   
349 loci or more (Table 2).

## 350 DISCUSSION

351 In this study, we fit the SFS of the Yoruba population with five simple demographic models  
352 of growth described by one parameter. Surprisingly, even though these five models are  
353 quite distinct in the way they model population growth, fitting them on the Yoruba data  
354 results in strongly similar SFS, which all show an excellent goodness of fit with the observed  
355 Yoruba SFS. Fitting the same SFS with the stairway plot method (LIU and FU 2015), a  
356 model-flexible method which infers a piecewise-constant demography, resulted in a complex  
357 demography with several bottlenecks in the last 160 000 years. The poor goodness of fit of  
358 the expected SFS under this inferred demography with the Yoruba SFS indicates that this  
359 complex demography is not to be trusted and suggests that the way the method estimates  
360 the number of change points is too flexible.

361 The results obtained by the model-constrained and model-flexible methods showed some  
362 similarities: the current population size  $N_e(0)$  of about 30 000 inferred with the stairway  
363 plot corresponds roughly to the coalescent unit of 67 000 generations (equivalent to  $2N_e(0)$   
364 in the coalescent theory) found with the one-parameter models. Similarly, the  $T_{MRCA}$  of  
365  $\sim 1.7$  million years inferred with the one-parameter models seems to match with the last  
366 time point of the stairway plot, at about 1.9 million years.

367 We hypothesize that the complexity of the demography inferred by the stairway plot  
368 method is caused by the irregularities of the observed Yoruba SFS. Two concurrent non-  
369 exclusive explanations can be put forward for these irregularities. First, they can be due  
370 to the sampling and thus be considered as noise that should not be interpreted as evidence  
371 for demography. Second, these irregularities could be biologically relevant and result from

372 a very complex demographic history. To assess the impact of noise on the stairway plot  
373 method, we tested it on simulated SFS under a *Linear* growth. These SFS were simulated  
374 with different numbers of independent loci: the more loci, the less noise in the simulated  
375 SFS. The stairway plot inference on these SFS shows that the method is strongly affected  
376 by the noise in the SFS simulated data: whereas the demography inferred for a smooth SFS  
377 (corresponding to a high number of independent loci) corresponds to the true demography  
378 approximated as piecewise constant, the demographies inferred for smaller numbers of loci  
379 show complex patterns of bottlenecks and deviate strongly from the true demography. It  
380 could be that this method captures the signal contained in these irregularities and infers  
381 a demography taking them into account, whereas the one-parameter models fit the global  
382 trend of the SFS shape and can thus infer the true demography for much smaller numbers  
383 of loci. One solution could be to constrain the number of parameters allowed for model-  
384 flexible methods: it seems that determining it by likelihood-ratio test, as it is done in the  
385 stairway plot method, is not conservative enough, as it does not prevent from overfitting  
386 the noise. If the number of parameters was forced to be small, the method might capture  
387 the global trend of the demography and avoid this issue. The SFS reconstructed under the  
388 demographies inferred by the stairway plot, however, differ strongly from the input SFS.  
389 If the issue was the overfitting of noise, we would expect the reconstructed SFS to fit the  
390 data more closely. The method is clearly biased by noise on the SFS but it remains unclear  
391 why. It would require further investigation to analyze how the different characteristics of  
392 this particular method, such as the parametrization of population size history, respond to  
393 noise, and what is responsible for this bias.

394 The five one-parameter demographic models all predict virtually the same SFS for the  
395 Yoruba population. Therefore, they also predict the same  $T_{MRCA}$  for the Yoruba population.  
396 This  $T_{MRCA}$  of  $\sim 1.3$  in coalescent units corresponds, with our scaling of coalescent time based  
397 on the number of mutations per site, to  $\sim 1.7$  million years. This estimation is similar to  
398 results concerning the whole human population, obtained by BLUM and JAKOBSSON (2011)

399 or reviewed in GARRIGAN and HAMMER (2006). Although the commonly admitted date  
400 of emergence of the anatomically modern human is around 200 000 years ago, BLUM and  
401 JAKOBSSON showed that finding a much older  $T_{MRCA}$  was compatible with the single-origin  
402 hypothesis, assuming a certain ancestral effective population size. These ancient times to  
403 most recent common ancestor could also be explained by gene flow in a structured ancestral  
404 population (GARRIGAN and HAMMER 2006).

405 Although all five models predict the same  $T_{MRCA}$ , the inferred demographies differ sub-  
406 stantially between the models (Figure 3A). In the time range further beyond the  $T_{MRCA}$ , no  
407 information is carried by the sample. Thus, the inferred demographies differ in this time  
408 range (Figure 4), making the inferred founding time of the Yoruba population unreliable.

409 Our results with one-parameter models are reproducible with another model-constrained  
410 method,  $\partial a \partial i$ , which uses different approaches both for the theoretical SFS simulations (dif-  
411 fusion approximation) and the parameter optimization (composite likelihood). This shows  
412 that, for models having the same number of parameters, a distance-based approach finds  
413 the same ranking of models as a likelihood framework, while being computationally less in-  
414 tensive. Furthermore, the distance-based approach allows for intuitive evidence on the fact  
415 that these different models actually all perform very well to fit the Yoruba SFS: the small  
416 differences of distance between the best SFS predicted by each model and the observed SFS  
417 could be due only to the noise in the observed SFS and thus do not mean that one model is  
418 better than another.

419 Among the five tested demographic models, two pairs of models seem to predict partic-  
420 ularly similar SFS (pairs of models with the two smallest values of  $d^2$  in Table 1). First,  
421 the *Linear* (L) and *Exponential* (E) growth models predict almost identical SFS for the  
422 Yoruba population ( $d^2(\tilde{\eta}^L, \tilde{\eta}^E) = 2.2 \times 10^{-5}$ ). Figure 4 shows that, in the time range where  
423 information is conveyed by the mean coalescent tree of the population, *i.e.*, between present  
424 time and the  $T_{MRCA}$ , these two demographies are very similar. This explains why their SFS  
425 are almost indistinguishable, and shows that in this parameter range, it is impossible to dis-

426 tinguish linear from exponential growth. Second, the SFS predicted under the two models  
427 with implicit demography, *Conditioned* (C) and *Birth-Death* (BD), are so similar that they  
428 are undistinguishable in Figure 3B ( $d^2(\tilde{\eta}^C, \tilde{\eta}^{BD}) = 3.5 \times 10^{-6}$ ). This raises a question on  
429 how these two models, based on different processes — a Wright-Fisher model or a branching  
430 process — compare and in particular why their SFS are so similar.

431 As we compute the distance statistic to optimize the models on normalized SFS, the  
432 information of the magnitude of the SFS (often referred to as  $\theta$ , the population mutation  
433 rate) is lost. However, as the inferred SFS under the five demographic models all have the  
434 same shape, the constant  $\theta$  by which they should be multiplied to fit the real, not normalized,  
435 Yoruba SFS would be the same for all five models. Thus, this information would not allow  
436 to choose which model infers the most realistic value of  $\theta$ .

437 The outlying first bin of the Yoruba SFS, corresponding to singletons, was removed  
438 from our inference because it can be affected by sequencing errors. As the relatively low to  
439 moderate coverage of the 1000 Genomes project could also result in an underestimation of  
440 doubletons and tripletons, we optimized  $\tau$  masking also these values. It did not change the  
441 estimation of  $\hat{\tau}$  and thus had no effect on the inferred demographies. As the first bin of the  
442 SFS accounts for the mutations that occur in the terminal branches of the coalescent tree, a  
443 large part of the excess of singletons can be due to very recent and massive growth. Recent  
444 studies with deep sequencing coverage have shown that there is a large abundance of rare  
445 variants in human populations (COVENTRY *et al.* 2010; NELSON *et al.* 2012; GAZAVE *et al.*  
446 2014). As the dataset we used for this study had a limited sample size and low-coverage,  
447 we focused on the inference of demography in the more distant past. Thus, because of both  
448 sequencing errors and incompatibility with our one-parameter models, singletons were not  
449 taken into account. Our inferences concern the population before this recent and massive  
450 growth. It should also be noted that LIU and FU (2015) emphasize that the strength of their  
451 method is in capturing recent demographic history. Thus, ignoring singletons, although it is  
452 an existing feature of their software, might not be the most appropriate use of the stairway

453 plot.

454 For non-African human population, the SFS based on the 1 000 Genomes Project data  
455 are not monotonous: their shape is more complex than the SFS of the Yoruba population.  
456 Thus, one-parameter models cannot capture the complexity of the demographic histories  
457 underlying these types of observed SFS. Even for the Yoruba population, capturing the  
458 recent growth event, by taking into account the singletons, would have required adding  
459 another parameter. The stairway plot method shows more flexibility and could capture the  
460 signal for more complex demographic histories, provided that the number of independent  
461 loci is very large so that there is no bias due to noise.

462 Overall, this study shows that even in the case of a simple demography, the scenario  
463 inferred by the stairway plot, a model-flexible method, can show spuriously complex patterns  
464 of growth and decline and can predict SFS poorly fitting with the initial SFS data. This  
465 might be explained by overfitting of the method to the noise present in the observed SFS,  
466 which can be expected for a reasonable number of loci. We also show that simple models  
467 described by one parameter can have an excellent goodness of fit to the data and avoid the  
468 issue of noise overfitting. The results indicate that the demography of the Yoruba population  
469 is compatible with simple one-parameter models of growth, and that the expected  $T_{MRC A}$  of  
470 this population can be estimated at  $\sim 1.7$  million years. However, the SFS is not sufficient to  
471 determine which model better characterizes the Yoruba demographic growth, and estimations  
472 of the founding time of the population, that depend on the chosen model, are thus unreliable.  
473 More generally, this study illustrates the issue of non-identifiability of demographies based  
474 on the SFS of a finite sample.

475 Our comparison of a model-constrained method using one parameter models with a  
476 model-flexible method using a potentially large number of parameters highlights the im-  
477 portance of the model complexity. How many parameters should we use to “properly”  
478 characterize a demography? We argue that low complexity models should be tested first.  
479 For model-flexible methods, the number of parameters is usually unbounded and determined

480 by successive likelihood ratio tests. This statistical framework implies that a certain risk is  
481 taken at each successive step, and that with the repetition of steps, errors can potentially  
482 be made. For example, these errors can lead to spurious inferences in noisy data (*i.e.*, any  
483 real data). We recommend (visually) monitoring the improvement in goodness of fit when  
484 adding new parameters on statistical grounds. Examination of the intermediate steps of  
485 fitting would likely prevent an unnecessary increase in the model complexity.

## 486 ACKNOWLEDGMENTS

487 We thank Cécile Delaporte for preliminary work on this project and Simon Boitard, Michael  
488 Blum, Konrad Lohse and three anonymous reviewers for useful comments on the manuscript.  
489 G.A. and M.L. acknowledge support from the grant ANR-12-NSV7-0012 Demochips from the  
490 Agence Nationale de la Recherche (France). M.L. is funded by the PhD program ‘Interfaces  
491 pour le Vivant’ of UPMC Univ Paris 06. G.A., A.L. and M.L. thank the *Center for Inter-*  
492 *disciplinary Research in Biology* for funding.

## LITERATURE CITED

- ACHAZ, G., 2008 Testing for neutrality in samples with sequencing errors. *Genetics* *179*(3):  
1409–1424.
- ACHAZ, G., 2009 Frequency spectrum neutrality tests: one for all and all for one. *Genet-*  
*ics* *183*(1): 249–258.
- ADAMS, A. M. and R. R. HUDSON, 2004 Maximum-likelihood estimation of demographic  
parameters using the frequency spectrum of unlinked single-nucleotide polymorphisms.  
*Genetics* *168*(3): 1699–1712.
- ATKINSON, Q. D., R. D. GRAY, and A. J. DRUMMOND, 2008 mtDNA variation pre-  
dicts population size in humans and reveals a major Southern Asian chapter in human  
prehistory. *Molecular biology and evolution* *25*(2): 468–474.

- BHASKAR, A. and Y. S. SONG, 2014 Descartes' rule of signs and the identifiability of population demographic models from genomic variation data. *Ann. Statist.* *42*(6): 2469–2493.
- BHASKAR, A., Y. R. WANG, and Y. S. SONG, 2015 Efficient inference of population size histories and locus-specific mutation rates from large-sample genomic variation data. *Genome research* *25*(2): 268–279.
- BLUM, M. G. and M. JAKOBSSON, 2011 Deep divergences of human gene trees and models of human origins. *Molecular biology and evolution* *28*(2): 889–898.
- CAMPBELL, C. D., J. X. CHONG, M. MALIG, A. KO, B. L. DUMONT, L. HAN, L. VIVES, B. J. O'ROAK, P. H. SUDMANT, J. SHENDURE, M. ABNEY, C. OBER, and E. E. EICHLER, 2012 Estimating the human mutation rate using autozygosity in a founder population. *Nat Genet* *44*(11): 1277–1281.
- CONRAD, D. F., J. E. KEEBLER, M. A. DEPRISTO, S. J. LINDSAY, Y. ZHANG, F. CASALS, Y. IDAGHDOUR, C. L. HARTL, C. TORROJA, K. V. GARIMELLA, M. ZILVERSMIT, R. CARTWRIGHT, G. ROULEAU, M. DALY, E. A. STONE, M. E. HURLES, and P. AWADALLA, 2011 Variation in genome-wide mutation rates within and between human families. *Nature genetics* *43*(7): 712–714.
- COVENTRY, A., L. M. BULL-OTTERSON, X. LIU, A. G. CLARK, T. J. MAXWELL, J. CROSBY, J. E. HIXSON, T. J. REA, D. M. MUZNY, L. R. LEWIS, and OTHERS, 2010 Deep resequencing reveals excess rare recent variants consistent with explosive population growth. *Nature communications* *1*: 131.
- DELAPORTE, C., G. ACHAZ, and A. LAMBERT, 2016 Mutational pattern of a sample from a critical branching population. *Journal of mathematical biology*: 1–38.
- EXCOFFIER, L., I. DUPANLOUP, E. HUERTA-SANCHEZ, V. C. SOUSA, and M. FOLL, 2013 Robust Demographic Inference from Genomic and SNP Data. *PLoS Genet* *9*(10): 1–17.

- FAY, J. C. and C.-I. WU, 2000 Hitchhiking under positive Darwinian selection. *Genetics* **155**(3): 1405–1413.
- FU, Y.-X., 1995 Statistical properties of segregating sites. *Theoretical population biology* **48**(2): 172–197.
- GARRIGAN, D. and M. F. HAMMER, 2006 Reconstructing human origins in the genomic era. *Nature Reviews Genetics* **7**(9): 669–680.
- GAZAVE, E., L. MA, D. CHANG, A. COVENTRY, F. GAO, D. MUZY, E. BOERWINKLE, R. A. GIBBS, C. F. SING, A. G. CLARK, and OTHERS, 2014 Neutral genomic regions refine models of recent rapid human population growth. *Proceedings of the National Academy of Sciences* **111**(2): 757–762.
- GRONAU, I., M. J. HUBISZ, B. GULKO, C. G. DANKO, and A. SIEPEL, 2011 Bayesian inference of ancient human demography from individual genome sequences. *Nature genetics* **43**(10): 1031–1034.
- GUTENKUNST, R. N., R. D. HERNANDEZ, S. H. WILLIAMSON, and C. D. BUSTAMANTE, 2009 Inferring the Joint Demographic History of Multiple Populations from Multidimensional SNP Frequency Data. *PLoS Genet* **5**(10): 1–11.
- HEIN, J., M. SCHIERUP, and C. WIUF, 2004 *Gene genealogies, variation and evolution: a primer in coalescent theory*. Oxford University Press, USA.
- HUDSON, R. R. and OTHERS, 1990 Gene genealogies and the coalescent process. *Oxford surveys in evolutionary biology* **7**(1): 44.
- KIM, J., E. MOSSEL, M. Z. RÁ CZ, and N. ROSS, 2015 Can one hear the shape of a population history? *Theoretical population biology* **100**: 26–38.
- KINGMAN, J. F. C., 1982 The coalescent. *Stochastic processes and their applications* **13**(3): 235–248.



- KONG, A., M. L. FRIGGE, G. MASSON, S. BESENBACHER, P. SULEM, G. MAGNUS-  
SON, S. A. GUDJONSSON, A. SIGURDSSON, A. JONASDOTTIR, A. JONASDOTTIR,  
W. S. W. WONG, G. SIGURDSSON, G. B. WALTERS, S. STEINBERG, H. HELGA-  
SON, G. THORLEIFSSON, D. F. GUDBJARTSSON, A. HELGASON, O. T. MAGNUS-  
SON, U. THORSTEINSDOTTIR, and K. STEFANSSON, 2012 Rate of de novo mutations and the  
importance of father's age to disease risk. *Nature* 488(7412): 471–475.
- LAMBERT, A., 2010 Population genetics, ecology and the size of populations. *Journal of  
mathematical biology* 60(3): 469–472.
- LAPIERRE, M., C. BLIN, A. LAMBERT, G. ACHAZ, and E. P. ROCHA, 2016 The im-  
pact of selection, gene conversion, and biased sampling on the assessment of microbial  
demography. *Molecular biology and evolution*: msw048.
- LI, H. and R. DURBIN, 2011 Inference of human population history from individual whole-  
genome sequences. *Nature* 475(7357): 493–496.
- LIU, X. and Y.-X. FU, 2015 Exploring population size changes using SNP frequency  
spectra. *Nature genetics* 47(5): 555–559.
- LUKIĆ, S., J. HEY, and K. CHEN, 2011 Non-equilibrium allele frequency spectra via  
spectral methods. *Theoretical population biology* 79(4): 203–219.
- MARTH, G. T., E. CZABARKA, J. MURVAI, and S. T. SHERRY, 2004 The Allele Frequency  
Spectrum in Genome-Wide Human Variation Data Reveals Signals of Differential Demo-  
graphic History in Three Large World Populations. *Genetics* 166(1): 351–372.
- MAZET, O., W. RODRIGUEZ, S. GRUSEA, S. BOITARD, and L. CHIKHI, 2016 On  
the importance of being structured: instantaneous coalescence rates and human  
evolution[mdash]lessons for ancestral population size inference[quest]. *Heredity* 116(4):  
362–371.
- MYERS, S., C. FEFFERMAN, and N. PATTERSON, 2008 Can one learn history from the

- allelic spectrum? *Theor Popul Biol* 73(3): 342–8.
- NAWA, N. and F. TAJIMA, 2008 Simple method for analyzing the pattern of DNA polymorphism and its application to SNP data of human. *Genes & genetic systems* 83(4): 353–360.
- NELSON, M. R., D. WEGMANN, M. G. EHM, D. KESSNER, P. S. JEAN, C. VERZILLI, J. SHEN, Z. TANG, S.-A. BACANU, D. FRASER, and OTHERS, 2012 An abundance of rare functional variants in 202 drug target genes sequenced in 14,002 people. *Science* 337(6090): 100–104.
- NIELSEN, R., 2000 Estimation of population parameters and recombination rates from single nucleotide polymorphisms. *Genetics* 154(2): 931–942.
- POLANSKI, A. and M. KIMMEL, 2003 New explicit expressions for relative frequencies of single-nucleotide polymorphisms with application to statistical inference on population growth. *Genetics* 165(1): 427–436.
- POOL, J. E., I. HELLMANN, J. D. JENSEN, and R. NIELSEN, 2010 Population genetic inference from genomic sequence variation. *Genome research* 20(3): 291–300.
- PYBUS, O. G., A. RAMBAUT, and P. H. HARVEY, 2000 An integrated framework for the inference of viral population history from reconstructed genealogies. *Genetics* 155(3): 1429–1437.
- SCALLY, A. and R. DURBIN, 2012 Revising the human mutation rate: implications for understanding human evolution. *Nat Rev Genet* 13(10): 745–753.
- TENNESSEN, J. A., A. W. BIGHAM, T. D. O’CONNOR, W. FU, E. E. KENNY, S. GRAVEL, S. MCGEE, R. DO, X. LIU, G. JUN, and OTHERS, 2012 Evolution and functional impact of rare coding variation from deep sequencing of human exomes. *science* 337(6090): 64–69.
- TERHORST, J. and Y. S. SONG, 2015 Fundamental limits on the accuracy of demographic

inference based on the sample frequency spectrum. *Proceedings of the National Academy of Sciences* *112*(25): 7677–7682.

THE 1000 GENOMES PROJECT CONSORTIUM, 2015 A global reference for human genetic variation. *Nature* *526*(7571): 68–74.

WALL, J. D., 1999 Recombination and the power of statistical tests of neutrality. *Genetical Research* *74*: 65–79.

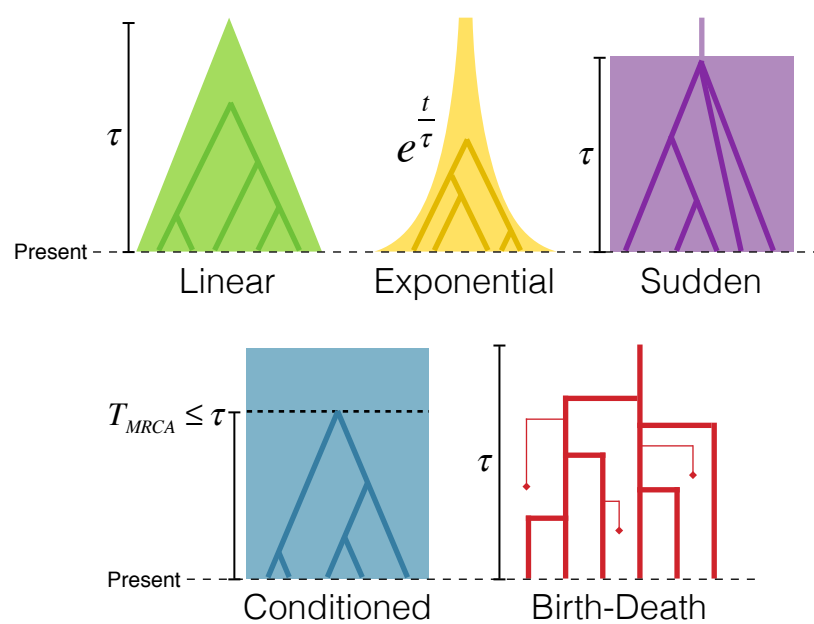


Figure 1: The five demographic models. Each model has one single time parameter  $\tau$ .

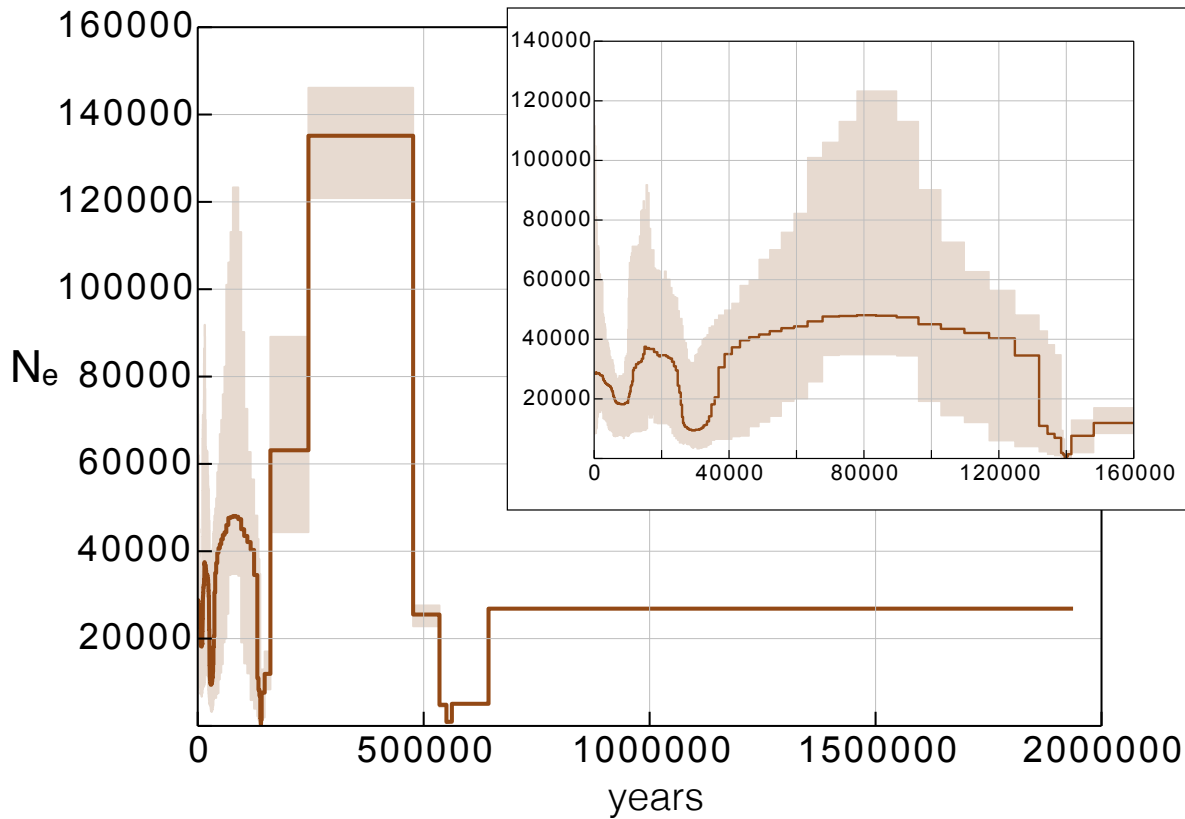


Figure 2: Stairway plot inference of the Yoruba demography. The inferred effective size  $N_e$  of the Yoruba population is plotted from present time (0) to the past. The inset is a zoom between 0 and 160 000 years. The thick brown line is the median  $N_e$ , the light brown area is the [2.5, 97.5] percentiles interval. The inference is based on 200 bootstrap samples of the unfolded Yoruba SFS. The singletons are not taken into account for the optimization of the stairway plot.

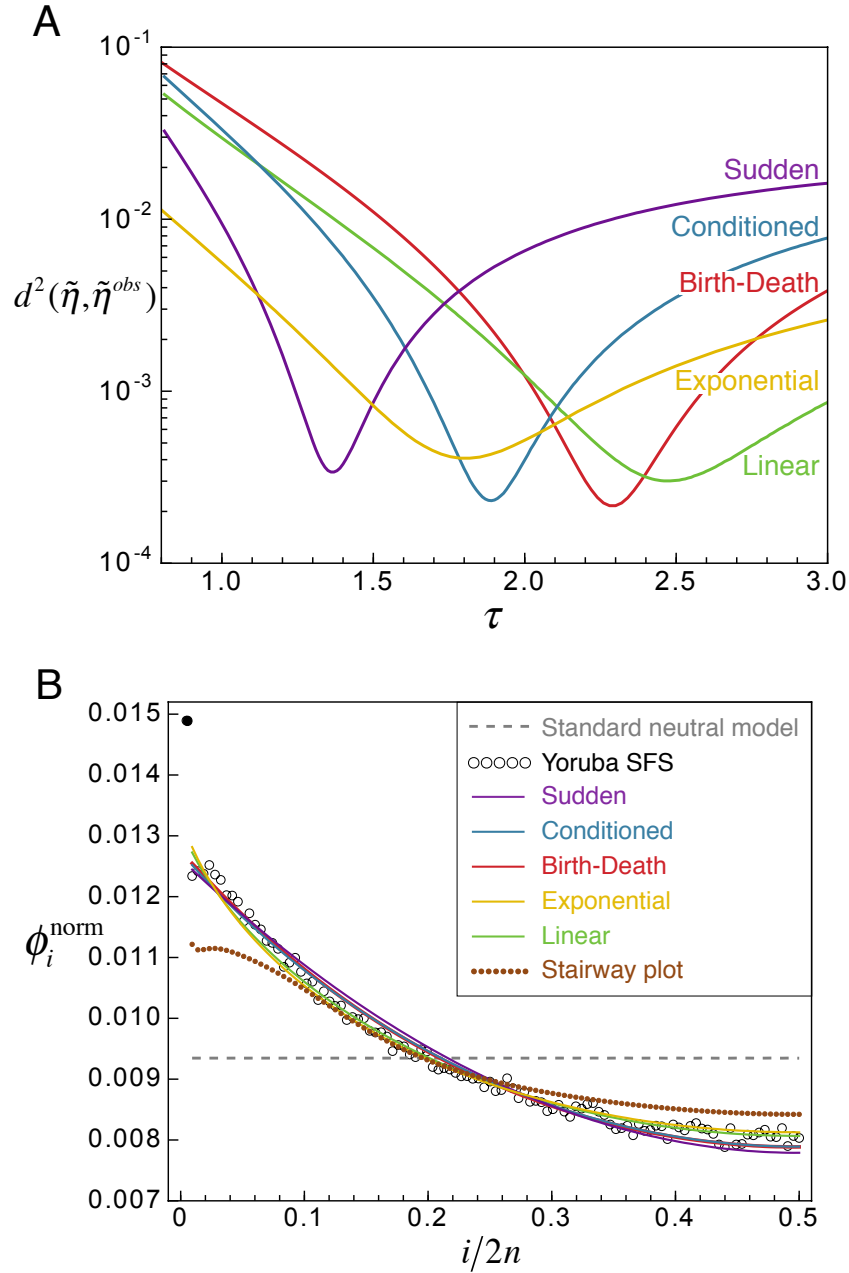


Figure 3: Inference of the Yoruba demography with one-parameter models. A) Weighted square distance  $d^2(\tilde{\eta}, \tilde{\eta}^{obs})$  between the normalized Yoruba SFS  $\tilde{\eta}^{obs}$  and the normalized predicted SFS  $\tilde{\eta}$  under each of the five models, depending on the value of the parameter  $\tau$  (Purple: *Sudden*, Blue: *Conditioned*, Red: *Birth-Death*, Yellow: *Exponential*, Green: *Linear*). B) Predicted SFS under each of the five models, with the optimized value  $\hat{\tau}$  of the parameter, and under the demography inferred by the stairway plot (brown dotted line). The Yoruba SFS is shown in empty circles. The first dot, colored in black, accounting for the singletons, was not taken into account for the optimization of  $\tau$  to avoid potential bias due to sequencing errors. The grey dashed line is the expected SFS under the standard neutral model without demography. Colors match the plot above (the predicted SFS under the models *Birth-Death* and *Conditioned* are indistinguishable). The SFS are folded, transformed and normalized (see Methods). 30

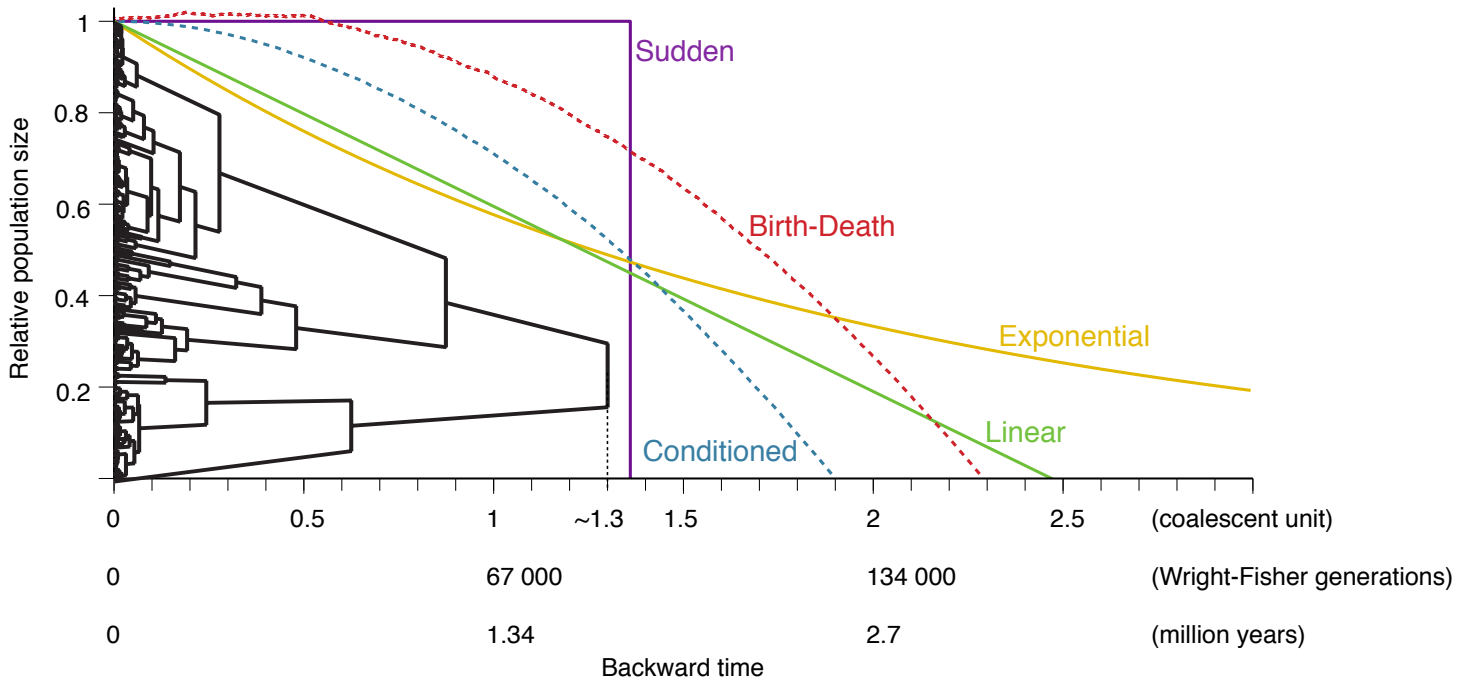


Figure 4: Demographic histories and reconstructed tree estimated from the Yoruba SFS. The tree shown has internode durations  $t_k$  during which there are  $k$  lineages consistent with the SFS (the topology was chosen uniformly among ranked binary trees with  $2n$  tips). Time is given in coalescent units, and scaled in number of generations and in millions of years. The demographic histories (solid lines: explicit models, dashed lines: implicit models) are plotted with their optimized  $\hat{\tau}$  values. See the supplementary methods for details on the demographic histories plotted for the models with implicit demographies (*Birth-Death* and *Conditioned*)

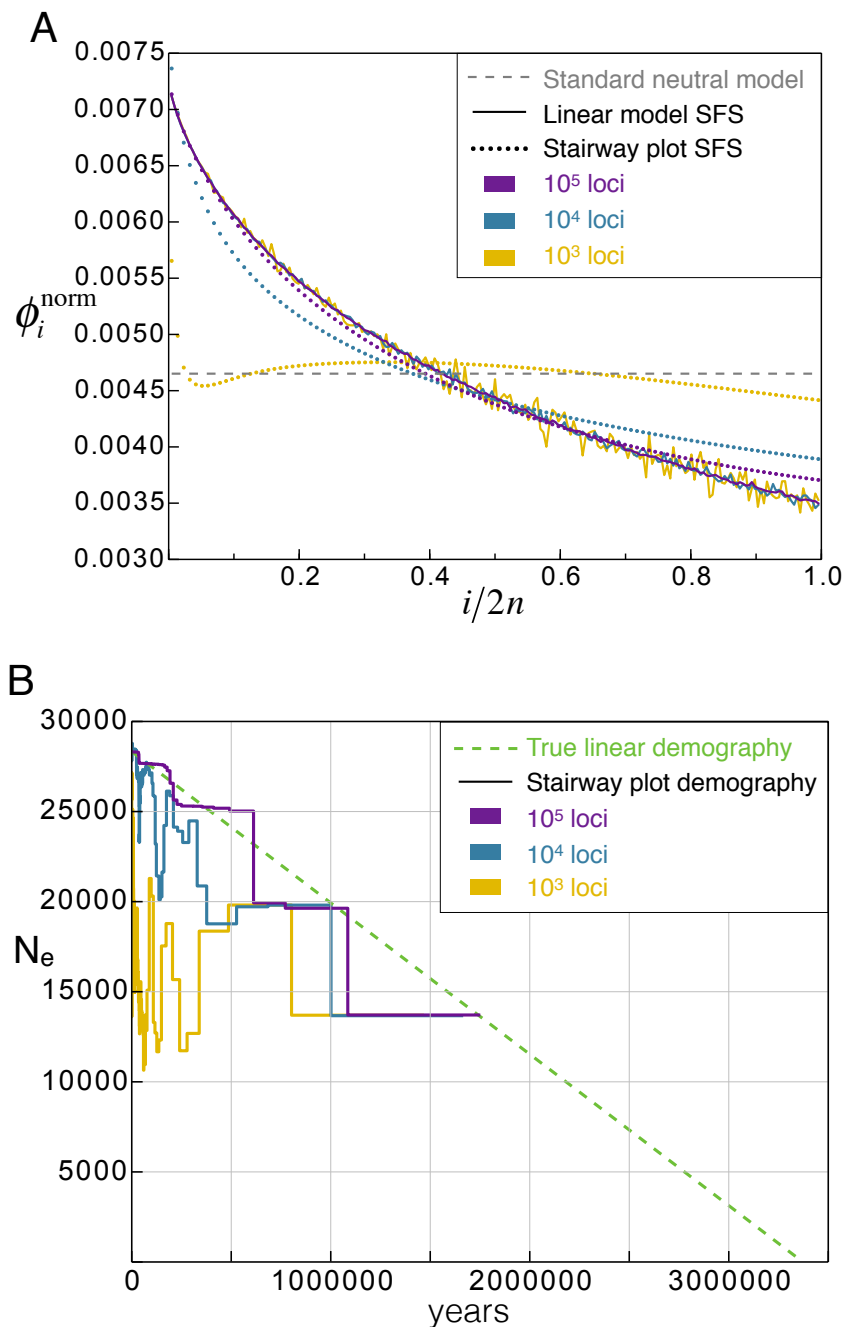


Figure 5: Stairway plot inference of a linear demography SFS with noise. A) Solid lines: mean of 200 SFS simulated independently under the *Linear* growth model, with either  $10^5$  loci (purple),  $10^4$  loci (blue) or  $10^3$  loci (yellow). Dotted lines: expected SFS under the demography reconstructed by the stairway plot method for different number of loci (same colors than solid lines). The grey dashed line is the expected SFS under the standard neutral model without demography. The SFS are transformed and normalized (see Methods). B) Stairway plot demographic inference: Median of 200 independent demographies inferred with 200 independently simulated SFS for each number of loci (colors match the plot above). The true demography is the green dashed line. The inferred effective size  $N_e$  is plotted from present time (0) to the past.



	Data	<i>Linear</i>	<i>Exponential</i>	<i>Sudden</i>	<i>Conditioned</i>	<i>Birth-Death</i>
<i>Linear</i>	$3.0 \times 10^{-4}$	0				
<i>Exponential</i>	$4.1 \times 10^{-4}$	$2.2 \times 10^{-5}$	0			
<i>Sudden</i>	$3.4 \times 10^{-4}$	$3.5 \times 10^{-4}$	$5.5 \times 10^{-4}$	0		
<i>Conditioned</i>	$2.3 \times 10^{-4}$	$1.6 \times 10^{-4}$	$5.5 \times 10^{-4}$	$3.7 \times 10^{-5}$	0	
<i>Birth-Death</i>	$2.2 \times 10^{-4}$	$1.7 \times 10^{-4}$	$3.1 \times 10^{-4}$	$4.1 \times 10^{-5}$	$3.5 \times 10^{-6}$	0
Stairway plot	$2.9 \times 10^{-3}$	$3.1 \times 10^{-3}$	$3.3 \times 10^{-3}$	$2.8 \times 10^{-3}$	$2.8 \times 10^{-3}$	$2.9 \times 10^{-3}$

Table 1: Least-square distance  $d^2$  between pairs of observed Yoruba SFS and optimized SFS under the five demographic models or the stairway plot method.

Number of loci	5% percentile	Mean $\hat{\tau}$	95% percentile
$10^3$	2.569	2.713	2.893
$10^4$	2.463	2.503	2.540
$10^5$	2.473	2.485	2.498
$10^6$	2.478	2.483	2.487

Table 2: Inference of the founding time  $\hat{\tau}$  under the *Linear* model on SFS with noise. Mean, 5% and 95% percentile of the founding time inferred with a *Linear* model. The SFS on which the inference is made are simulated with a founding time  $\tau$  of 2.48, with different number of loci, using the method with topology reconstruction.