

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

## INTRODUCTION

16

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

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

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

43 variance of the SFS while its expectation remains unchanged (WALL 1999). Therefore, the  
44 SFS of a sample is a summary of the genetic diversity, averaged over all the genome due to  
45 recombination, that can be analyzed in terms of demography.

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

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

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

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

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

## 93 MATERIALS AND METHODS

94 **1 000 Genomes Project data:** Variant calls from the 1000 Genomes Project phase 3  
95 were downloaded from the project ftp site (THE 1000 GENOMES PROJECT CONSORTIUM

2015). The sample size for the Yoruba population is  $n = 108$  individuals (polymorphism data available for both genome copies of each individual, *i.e.*,  $2n = 216$  sequences). We kept all single nucleotide bi-allelic variants to plot the sample SFS. To avoid possible bias due to sequencing errors, we ignored singletons (mutations appearing in only one chromosome of one individual in the sample) for the rest of the study.

**Site Frequency Spectrum definition and graphical representation:** The Site Frequency Spectrum (SFS) of a sample of  $n$  diploid individuals is described as the vector  $\xi = (\xi_1, \xi_2, \dots, \xi_{2n-1})$  where for  $i \in [1, 2n - 1]$ ,  $\xi_i$  is the number of dimorphic (*i.e.*, with exactly two alleles) sites with derived form at frequency  $i/2n$ . To avoid potential orientation errors, we assumed that the ancestral form is unknown for all sites: we worked with a folded spectrum, where we consider the frequency of the less frequent (or minor) allele. In this 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  $i \in [1, n - 1]$  and  $\eta_n = \xi_n$ . For a better graphical representation, all SFS were transformed as follows: we plot  $\phi_i$  normalized by its sum, where

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

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

**Demographic models used for the model-constrained method:** We inferred the demography of the Yoruba population using five growth models (Figure 1), compared with the predictions of the standard model with constant population size. Time is measured in coalescent units of  $2N$  generations, where the scaling parameter  $N$  has the same dimension as the current population size, which we will not estimate. Time starts at 0 (present time) and increases backward in time. Four models are based on the standard Kingman coalescent

120 (KINGMAN 1982), amended with demography. Three of them are described with an explicit  
121 demography: either *Linear* growth since time  $\tau$ , *Exponential* growth at rate  $1/\tau$  or *Sudden*  
122 growth from a single ancestor to the entire population at time  $\tau$ . We also use another model  
123 based on the Kingman coalescent, with an implicit demography: the *Conditioned* model.  
124 This model is based on a standard constant size model, but the Time to the Most Recent  
125 Common Ancestor ( $T_{MRC A}$ ) is conditioned on being reached before time  $\tau$ . The fifth model,  
126 *Birth-Death*, is not based on the standard Kingman coalescent, but on a critical branching  
127 process measured in units of  $2N$  generations. In forward time (from the past to the present),  
128 the process starts with a founding event of one individual. Individuals give birth and die at  
129 equal rate 1. The process is conditioned on not becoming extinct before a period of time  $\tau$ ,  
130 and on reaching on average  $2N$  individuals.

131 **Stairway plot inference on the Yoruba SFS:** We applied the model-flexible stairway  
132 plot method developed by LIU and FU (2015) on the unfolded Yoruba SFS. Inferences are  
133 made on 200 SFS as suggested by their method. We use the script they provide to create  
134 199 bootstrap samples of the Yoruba SFS. We also ignore the singletons for this method,  
135 and use the default parameter values suggested in their paper for the optimization.

136 **SFS simulation with demography:** We used two different method to simulate SFS under  
137 the four demographic models derived from the Kingman coalescent (*Linear*, *Exponential*,  
138 *Sudden* and *Conditioned*) or under a piecewise-constant demography reconstructed by the  
139 stairway plot method.

140 *Method 1:* A first method is to simulate  $l$  independent topologies under the Kingman coa-  
141 lescent on which mutations are placed at rate  $\theta$  (population mutation rate) (HUDSON *et al.*  
142 1990). This allows to simulate the SFS of  $l$  independent loci.



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

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

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

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

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

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

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

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

165 **Optimization of the parameter  $\tau$ :** For each demographic model, we optimize the pa-  
166 rameter  $\tau$  by minimizing the weighted square distance  $d$  between the observed SFS of the  
167 Yoruba population and the predicted SFS under the model. With  $\eta^{model}$  and  $\eta^{obs}$  the folded  
168 SFS in the tested model and in the data respectively,

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

169 The sum starts at  $i = 2$  because we ignore  $\eta_1^{obs}$ , corresponding to singletons, to avoid bias  
170 due to sequencing errors. To calculate the distance  $d'$  between the SFS predicted by two  
171 models A and B, we weight the terms by the mean of the two models:

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

172 **Scaling of the coalescent time:** Optimized values of the parameter  $\hat{\tau}$  for each model are  
173 expressed in coalescent time units, *i.e.*, scaled in  $2N$  generations. As the present population  
174 size  $N$  is unknown, to scale these coalescent time units in numbers of generations and conse-  
175 quently in years, we used the expected number of mutations per site  $M$ . From the dataset,  
176 we have  $M^{obs} = S/L$  where  $S$  is the number of single nucleotide mutations (a  $k$ -allelic SNP  
177 accounts for  $k - 1$  mutations) and  $L$  is the length of the accessible sequenced genome in the  
178 1 000 genomes project (90% of the total genome length, THE 1000 GENOMES PROJECT  
179 CONSORTIUM). We can state that  $M^{theo} = \mu \hat{T}_{tot} C$ , where we know the mutation rate  $\mu$   
180 from the literature and the total tree length expressed in coalescent time units  $\hat{T}_{tot}$  from  
181 the SFS simulations. Here  $C$  is the coalescent factor, that is the number of generations per  
182 coalescent time unit, also corresponding to  $2N_e(0)$  where  $N_e(0)$  is the effective population  
183 size at present time. The total number of generations in the tree is  $\hat{T}_{tot} C$  from which we  
184 derive the total number of mutations per site  $M^{theo}$ . Thus, equaling  $M^{obs}$  with  $M^{theo}$ , we  
185 can estimate  $C$  by  $S/(\mu L \hat{T}_{tot})$ . We assumed a mutation rate of  $1.2 \times 10^{-8}$  per base pair  
186 per generation (CONRAD *et al.* 2011; CAMPBELL *et al.* 2012; KONG *et al.* 2012). With the  
187 coalescent factor  $C$ , we can then convert a coalescent time unit into a number of generations,

188 or into a number of years assuming 24 years as generation time (SCALLY and DURBIN 2012).

189 **Graphical representation of the inferred demographies:** To represent the inferred  
190 explicit demographies (models *Linear*, *Exponential* and *Sudden*), we plot the shape of the  
191 demography with the optimized value  $\hat{\tau}$  for each model. For the implicit demographies  
192 (models *Conditioned* and *Birth-Death*), as there is no explicit demographic shape, we plot the  
193 mean trajectory of fixation of a new allele in the population: in forward time, these fixation  
194 trajectories illustrate the expansion in the population of the descendance of the sample's  
195 ancestor. For the *Conditioned* model, we use the Wright-Fisher diffusion conditioned upon  
196 fixation (LAMBERT 2008) to simulate trajectories of fixation:

$$dX_t = (1 - X_t)dt + \sqrt{X_t(1 - X_t)} dB_t$$

197 where  $X_t$  is the random variable accounting for the frequency of the allele at time  $t$  and  $B_t$  is  
198 Brownian motion. We simulate the trajectories starting at  $X_0 = 0.01$  with  $dt = 0.0001$  and  
199 we stop the trajectories when  $X_t$  reaches 1. To account for the specificity of the *Conditioned*  
200 model, we keep only trajectories that reach fixation in a time smaller than the optimized  
201 parameter value  $\hat{\tau}$ . Similarly, for the Birth-Death model, we use the critical Feller diffusion  
202 (LAMBERT 2008):

$$dX_t = \sqrt{2X_t}dB_t$$

203 and we run trajectories until time reaches the optimized parameter value  $\hat{\tau}$ . We keep tra-  
204 jectories for which  $X_{\hat{\tau}} \in (U_n, U_{n+1})$ , where  $U_k = \sum_{i=1}^k V_i$  and the  $V_i$ 's are independent expo-  
205 nential random variables with mean  $1/n$ , which amounts to conditioning upon sampling  $n$   
206 individuals at time  $\hat{\tau}$ . For both models, we average over 5 000 trajectories.

207 **Comparing the model-constrained and model-flexible methods for *Linear* de-**  
208 **mography inference:** We applied both methods on simulated SFS under the *Linear* model.  
209 To test the stairway plot method on a *Linear* model demography, we simulate 200 indepen-  
210 dent SFS using method 1, with  $\theta = 100$  (arbitrary value removed by normalization) and a

211 foundation time  $\tau$  (here we used  $\hat{\tau} = 2.48$  that we estimated for the Yoruba population, see  
212 Results). The SFS are simulated with either  $10^3$ ,  $10^4$  or  $10^5$  independent loci. We run the  
213 stairway plot method with the default parameter values suggested in the method, and with  
214 the same mutation rate ( $1.2 \times 10^{-8}$  per base pair per generation) and generation time (24  
215 years) than in our study.

216 To test the one-parameter inference method on these SFS simulated under the *Linear*  
217 model, we run the parameter optimization on a SFS simulated with either  $10^3$ ,  $10^4$ ,  $10^5$  or  
218  $10^6$  loci. The search of the parameter value that minimizes the distance  $d$  was optimized  
219 with a Newton-Raphson algorithm. Derivatives were calculated at  $t \pm 0.05$  where  $t$  is the  
220 parameter value being optimized. The optimization stopped when the optimization step of  
221 the parameter value was smaller than  $10^{-3}$ .

222 **Data and software availability** The 1000 genomes project data used in this study is  
223 publicly available at <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/>.  
224 The code in Python and C written for the study is available at [https://github.com/](https://github.com/lapierreM/Yoruba_demography)  
225 [lapierreM/Yoruba\\_demography](https://github.com/lapierreM/Yoruba_demography). The code in C used for the method 1 of SFS simulation is  
226 available upon request to G. ACHAZ.

## 227 RESULTS

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

232 It has been shown that human populations have been growing since their emergence in  
233 Africa, and that African populations were supposedly non-affected by the Out of Africa  
234 bottleneck described for Eurasian populations (MARTH *et al.* 2004; GUTENKUNST *et al.*  
235 2009). Based on this previous knowledge, for the model-constrained method, we chose to

236 infer the Yoruba demography with simple models of growth, *i.e.*, with only one phase of  
237 growth characterized by a single parameter. These five models are: *Linear*, *Exponential*  
238 or *Sudden* growth, a *Conditioned* model where the  $T_{MRCA}$  is conditioned on being smaller  
239 than the given parameter, and a critical *Birth-Death* model based on a branching process  
240 (Figure 1). To infer the Yoruba demography with this method, we fit the SFS predicted under  
241 each model with the observed Yoruba SFS (all SFS are folded). The SFS were normalized  
242 to remove the population mutation rate parameter  $\theta$ , so that each model is characterized by  
243 one single parameter  $\tau$  which has the dimension of a time duration. We fit this parameter  
244 by least-square distance between the observed SFS and the predicted SFS. For the model-  
245 flexible inference, we used the stairway plot method developed recently by LIU and FU  
246 (2015), which infers a piecewise-constant demography based on the SFS. For this method,  
247 the number of parameters to be estimated is determined by a likelihood-ratio test. It can  
248 range from 1 to  $2n - 1$  where  $2n$  is the number of sequences in the sample.

249 The Yoruba SFS was constructed by taking into account the entire genome. Removing  
250 the coding parts of the genome to avoid potential bias due to selection does not affect the  
251 shape of the SFS (Figure S1) as the coding parts represent a very small fraction of the  
252 human genome. The first bin of the observed SFS, accounting for mutations found in one  
253 chromosome of one individual in the sample (black dot in the observed SFS in Figure 3B)  
254 seemed to lie outside the rest of the distribution. Because this could be due to sequencing  
255 errors being considered as singletons (ACHAZ 2008), we chose to ignore this value for the  
256 model optimization.

257 The analysis of the Yoruba SFS with the stairway plot method results in a complex  
258 demography with several bottlenecks in the last 160 000 years (Figure 2). The current  
259 effective population size  $N_e(0)$  is 28 500. The demographic history earlier than 160 000 years  
260 shows spurious patterns that should not be interpreted, according to LIU and FU (2015).

261 The inference of the Yoruba demography with one-parameter models was done by min-  
262 imizing the distance between observed and predicted SFS. This gave an optimized value  $\hat{\tau}$

263 of the parameter  $\tau$  (Figure 3A and Table 1) (with  $\hat{\tau}$  in coalescent units, *Linear*:  $\hat{\tau} = 2.48$ ,  
264 *Exponential*:  $\hat{\tau} = 1.79$ , *Sudden*:  $\hat{\tau} = 1.36$ , *Conditioned*:  $\hat{\tau} = 1.89$ , *Birth-Death*:  $\hat{\tau} = 2.28$ ).  
265 Plotting the predicted SFS with the optimized parameter value  $\hat{\tau}$  confirmed their goodness  
266 of fit with the observed Yoruba SFS (Figure 3B). Compared to the standard model with-  
267 out demography, the addition of just one parameter allows for a surprisingly good fit of  
268 the observed Yoruba SFS. The Yoruba demography thus seems to be compatible with a  
269 simple scenario of growth. On the other hand, the demography inferred by the stairway  
270 plot predicts a SFS which does not fit well the observed Yoruba SFS: the distance between  
271 the observed Yoruba SFS and the expected SFS under the stairway plot demography is ten  
272 times the distance between any of the one-parameter model SFS and the data (Figure 3B  
273 and Table 1).

274 The best fitting SFS under each of the five demographic models all have a least square  
275 distance  $d$  of the order of  $10^{-4}$  with the observed Yoruba SFS (Figure 3A and Table 1) and  
276 have highly similar shapes (Figure 3B). This shows that the five demographic models used  
277 to infer the demography of the Yoruba are indistinguishable based on only the observed SFS.  
278 To back up this assertion, we computed the expected  $T_{MRCA}$  based on the predicted SFS  
279 using (1): as the SFS predicted under each model are very similar, it means that they have  
280 roughly the same estimated time durations  $t_k$  while there are  $k$  branches in the coalescent  
281 tree of the Yoruba sample. From these expected  $t_k$  we can compute  $T_{MRCA} = \sum_{i=2}^{2n} t_k$ . Under  
282 each of four models (excluding the *Birth-Death* model for which there is no obvious common  
283 time scaling), the inferred  $T_{MRCA}$  for the Yoruba population is 1.3 in coalescent units. By  
284 using the number of mutations per site in the data and the total tree length inferred from  
285 the simulations, we scaled back this  $T_{MRCA}$  in number of generations and in years, assuming  
286 a mutation rate of  $1.2 \times 10^{-8}$  per base pair per generation (CONRAD *et al.* 2011; CAMPBELL  
287 *et al.* 2012; KONG *et al.* 2012) and a generation time of 24 years (SCALLY and DURBIN 2012)  
288 (see Methods). The  $T_{MRCA}$  of the Yoruba population inferred under the four demographic  
289 models is of 87 100 generations corresponding to 1.7 million years. The inferred demographic

290 models, with scaling in coalescent units, number of generations and number of years, are  
291 shown on Figure 4. The coalescent unit of 67 000 estimated to scale the inferred coalescent  
292 times in number of years corresponds to a present effective population size  $N_e(0)$  of 33 500.

293 The demography inferred by the stairway plot method for the Yoruba population is a  
294 piecewise-constant demography showing much more complex patterns of growth and bottle-  
295 necks than our one-parameter models (Figure 2). Moreover, the expected SFS under this  
296 inferred demography does not fit well the observed Yoruba SFS (Figure 3B). To understand  
297 what could produce such a complex demography, we simulated SFS under the *Linear* model  
298 with the foundation time  $\hat{\tau} = 2.48$  inferred for the Yoruba population. The SFS were sim-  
299 ulated with different numbers of loci, to obtain SFS with more or less noise (solid lines on  
300 Figure 5A). We applied the two inference methods to these SFS. The demographies inferred  
301 by the stairway plot method are strongly affected by the noise of the SFS, as shown on  
302 Figure 5B. When the number of simulated loci is very large ( $200 \times 10^6$  loci), the stairway  
303 plot gives a good approximation of the true demography, and the expected SFS under the  
304 inferred demography fits the input SFS. However, for smaller numbers of loci ( $200 \times 10^5$  loci  
305 or less), the stairway plot shows complex patterns of growth and bottlenecks incompatible  
306 with the true demography, and the expected SFS under the inferred demographies do not fit  
307 the input SFS. On the contrary, the one-parameter method infers a *Linear* demography with  
308 a foundation time close to the true value for SFS simulated with  $10^4$  loci or more (Table 2).

## 309 DISCUSSION

310 In this study, we fit the SFS of the Yoruba population with five simple demographic models  
311 of growth described by one parameter. Surprisingly, even though these five models are  
312 quite distinct in the way they model population growth, their fitting on the Yoruba data  
313 results in strongly similar SFS, which all show an excellent goodness of fit with the observed  
314 Yoruba SFS. Fitting the same SFS with the stairway plot method (LIU and FU 2015), a  
315 model-flexible method which infers a piecewise-constant demography, resulted in a complex

316 demography with several bottlenecks in the last 160 000 years. The poor goodness of fit of  
317 the expected SFS under this inferred demography with the Yoruba SFS indicates that this  
318 complex demography is not to be trusted and suggests that the way the method estimates  
319 the number of change points is too flexible.

320 The results obtained by the model-constrained and model-flexible methods showed some  
321 similarities: the current population size  $N_e(0)$  of about 30 000 inferred with the stairway plot  
322 corresponds roughly to the coalescent unit of 67 000 years (equivalent to  $2N_e(0)$  in the coa-  
323 lescent theory) found with the one-parameter models. Similarly, the  $T_{MRC A}$  of  $\sim 1.7$  million  
324 years inferred with the one-parameter models seems to match with the last time point of the  
325 stairway plot, at about 1.9 million years.

326 We postulate that the complexity of the demography inferred by the stairway plot method  
327 is due to the fitting of irregularities of the observed Yoruba SFS. Two concurrent non-  
328 exclusive explanations can be put forward for these irregularities. First, they can be due  
329 to the sampling and thus be considered as noise that should not be interpreted as evidence  
330 for demography. Second, these irregularities could be biologically relevant and result from  
331 a very complex demographic history. To assess the impact of noise on the stairway plot  
332 method, we tested it on simulated SFS under the *Linear* model. These SFS were simulated  
333 with different numbers of independent loci: the more loci, the less noise in the simulated  
334 SFS. The stairway plot inference on these SFS shows that the method is strongly affected  
335 by the noise in the SFS simulated data: whereas the demography inferred for a smooth SFS  
336 (corresponding to a high number of independent loci) corresponds to the true demography  
337 approximated as piecewise constant, the demographies inferred for smaller numbers of loci  
338 show complex patterns of bottlenecks and deviate strongly from the true demography. This  
339 method captures the signal contained in these irregularities and infers a demography taking  
340 them into account, whereas the one-parameter models fit the global trend of the SFS shape  
341 and can thus infer the true demography for much smaller numbers of loci. One solution  
342 could be to constrain the number of parameters allowed for model-flexible methods: it seems



343 that determining it by likelihood-ratio test, as it is done in the stairway plot method, is  
344 not conservative enough, as it does not prevent from overfitting the noise. If the number  
345 of parameters was forced to be small, the method might capture the global trend of the  
346 demography and avoid this issue.

347 The five one-parameter demographic models all predict virtually the same SFS for the  
348 Yoruba population. This implies that they predict the same  $T_{MRCA}$  for the Yoruba popula-  
349 tion. This  $T_{MRCA}$  of  $\sim 1.3$  in coalescent units corresponds, with our scaling of coalescent time  
350 based on the number of mutations per site, to  $\sim 1.7$  million years. This estimation is simi-  
351 lar to results concerning the whole human population, obtained by BLUM and JAKOBSSON  
352 (2011) or reviewed in GARRIGAN and HAMMER (2006). Although the commonly admitted  
353 date of emergence of the anatomically modern human is around 200 000 years ago, BLUM and  
354 JAKOBSSON showed that finding a much older  $T_{MRCA}$  was compatible with the single-origin  
355 hypothesis, assuming a certain ancestral effective population size. These ancient times to  
356 most recent common ancestor could also be explained by gene flow in a structured ancestral  
357 population (GARRIGAN and HAMMER 2006).

358 Although all five models predict the same  $T_{MRCA}$ , the predictions of the population's  
359 foundation time differ largely between the models (Figure 2A). The comparison of the in-  
360 ferred demographies (Figure 3) suggests that in the time range further behind the  $T_{MRCA}$ ,  
361 little information is carried by the sample. Thus, the inferred demographies differ in this  
362 time range, making the inferred foundation time of the Yoruba population unreliable.

363 Among the five tested demographic models, two pairs of models seem to predict partic-  
364 ularly similar SFS (pairs of models with the two smallest values of  $d$  in Table 1). First, the  
365 *Linear* (L) and *Exponential* (E) growth models predict almost identical SFS for the Yoruba  
366 population ( $d(\eta^L, \eta^E) = 2.2 \times 10^{-5}$ ). Figure 3 shows that, in the time range where informa-  
367 tion is conveyed by the mean coalescent tree of the population, *i.e.*, between present time  
368 and the  $T_{MRCA}$ , these two demographies are very similar. This explains why their SFS are  
369 almost indistinguishable, and shows that in this parameter range, it is impossible to distin-

370 guish linear from exponential growth. Second, the two models with implicit demography,  
371 *Conditioned* (C) and *Birth-Death* (BD), predict so similar SFS that they are completely  
372 overlaid on Figure 2B ( $d(\eta^C, \eta^{BD}) = 3.5 \times 10^{-6}$ ). This raises a question on how these two  
373 models, based on different processes — a Wright-Fisher model or a branching process —  
374 compare and in particular why their SFS are so similar.

375 The outlying first bin of the Yoruba SFS, corresponding to singletons, was removed from  
376 our inference because it can be spoiled by sequencing errors. However, as this first bin  
377 accounts for the mutations that occur in the terminal branches of the coalescent tree, this  
378 excess of singletons could alternatively be due to very recent and massive growth, aspects  
379 that we cannot capture with our one-parameter demographic models.

380 For non-African human population, the SFS based on the 1 000 Genomes Project data  
381 are not monotonous: their shape is more complex than the SFS of the Yoruba population.  
382 Thus, one-parameter models cannot capture the complexity of the demographic histories  
383 underlying these types of observed SFS. The stairway plot method shows more flexibility  
384 and could capture the signal for more complex demographic histories, provided that the  
385 number of independent loci is very large so that there is no bias due to noise.

386 Overall, this study shows that even in the case of a simple demography, the scenario  
387 inferred by a model-flexible method like the stairway plot can show spuriously complex  
388 patterns of growth and decline and can predict SFS poorly fitting with the initial SFS data.  
389 This can be explained by overfitting of the method to the noise present in the observed  
390 SFS, which can be expected for a reasonable number of loci. We also show that simple  
391 models described by one parameter can have an excellent goodness of fit to the data and  
392 avoid the issue of noise overfitting. The results indicate that the demography of the Yoruba  
393 population is compatible with simple one-parameter models of growth, and that the expected  
394  $T_{MRC A}$  of this population can be estimated at  $\sim 1.7$  million years. However, the SFS does  
395 not allow to determine which model characterizes better the Yoruba demographic growth,  
396 and estimations of the foundation time of the population, that depend on the chosen model,

397 are thus unreliable. More generally, this study illustrates the issue of non-identifiability  
398 of demographies based on the SFS. It also highlights the need to constrain model-flexible  
399 methods to avoid interpreting noise as signal in demographic inferences.

## 400 ACKNOWLEDGMENTS

401 We thank Cécile Delaporte for preliminary work on this project and Simon Boitard, Michael  
402 Blum and Konrad Lohse for useful comments on the manuscript. G.A. and M.L acknowledge  
403 support from the grant ANR-12-NSV7-0012 Demochips from the Agence Nationale de la  
404 Recherche (France). M.L. is funded by the PhD program ‘Interfaces pour le Vivant’ of  
405 UPMC Univ Paris 06. G.A., A.L. and M.L. thank the *Center for Interdisciplinary Research*  
406 *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. *Genetics* **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.
- 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.

- 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.
- 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. HELGASON, 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., 2008 Population Dynamics and Random Genealogies. *Stochastic Models* 24(sup1): 45–163.
- 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 impact of selection, gene conversion, and biased sampling on the assessment of microbial demography. *Molecular biology and evolution*: msw048.
- 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 Demographic History in Three Large World Populations. *Genetics* 166(1): 351–372.
- 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.
- 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.
- 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.

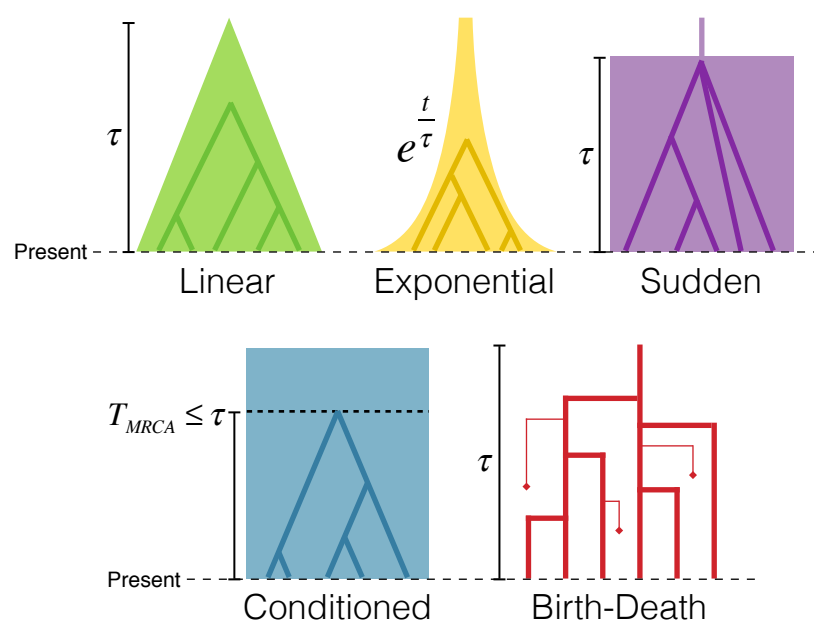


Figure 1: The five demographic models. Each model has one single time parameter  $\tau$ .



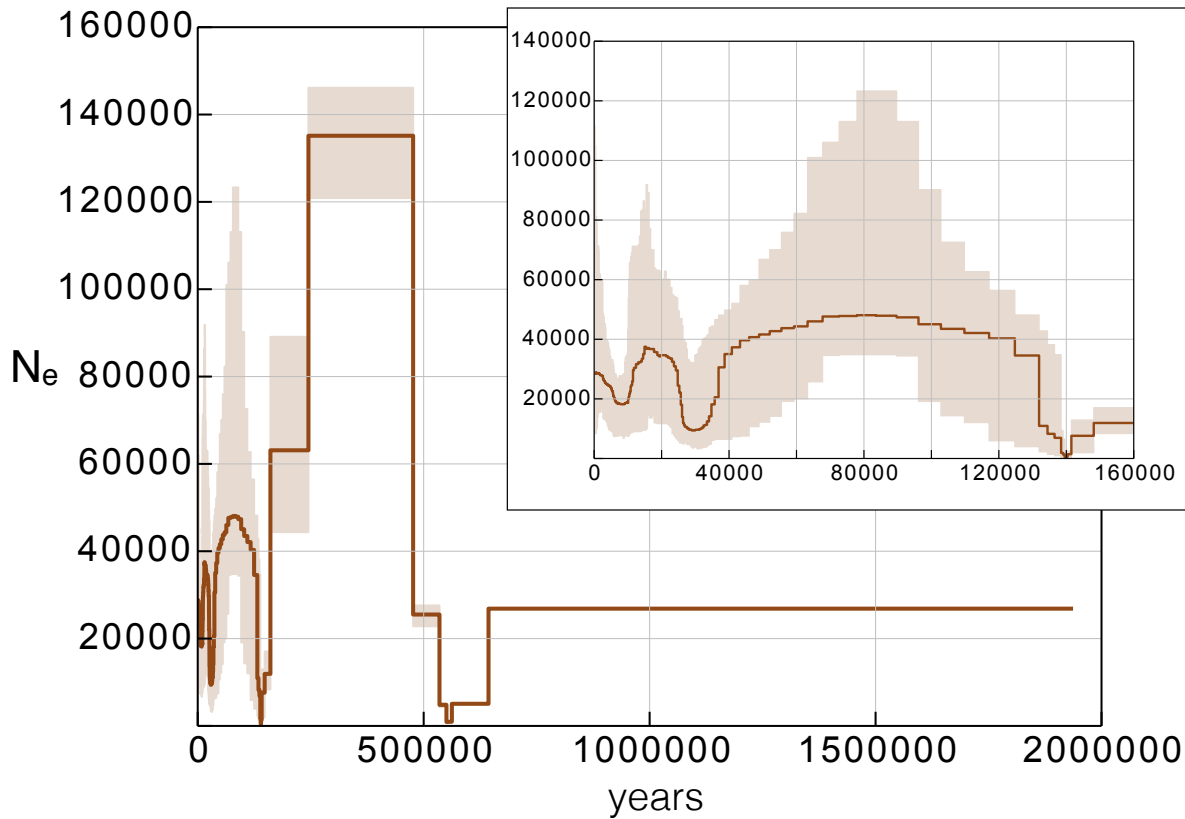


Figure 2: Stairway plot inference of the Yoruba demography. The inferred effective size  $N_e$  of the Yoruba population is plotted from present time (0) to the past. The inset is a zoom between 0 and 160 000 years. The thick brown line is the median  $N_e$ , the light brown area is the [2.5, 97.5] percentiles 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.

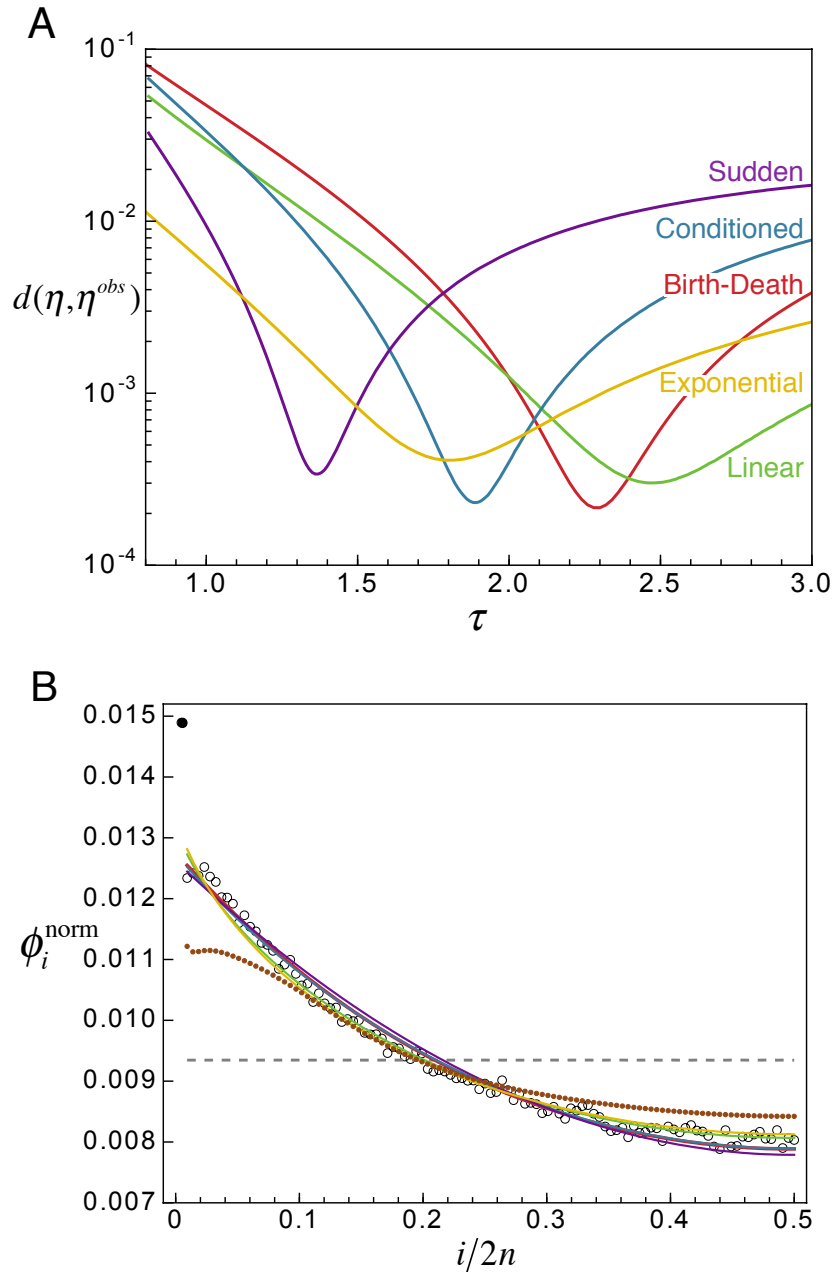


Figure 3: Inference of the Yoruba demography with one-parameter models. A) Weighted square distance  $d(\eta, \eta^{obs})$  between the Yoruba SFS  $\eta^{obs}$  and the predicted SFS  $\eta$  under each of the five models, depending on the value of the parameter  $\tau$  (Purple: *Sudden*, Blue: *Conditioned*, 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).

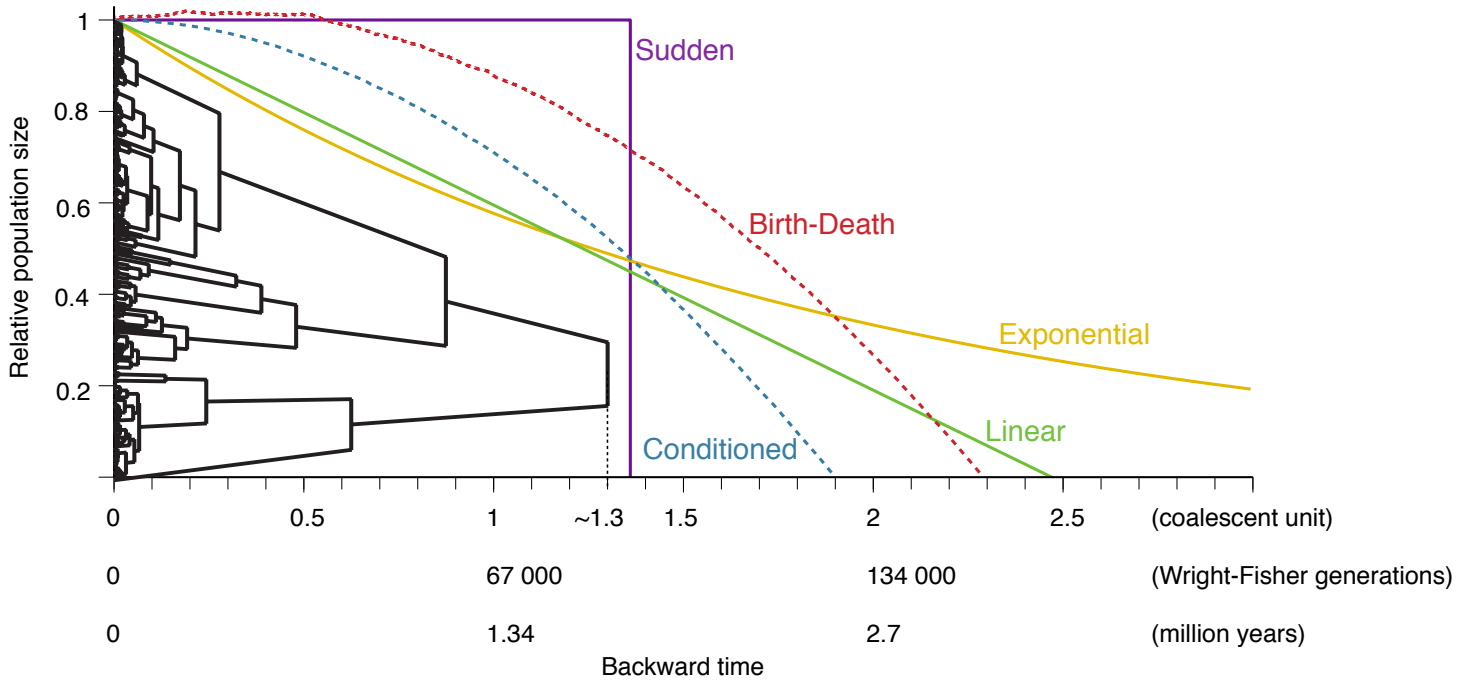


Figure 4: Demographic histories and reconstructed tree estimated from the Yoruba SFS. The tree shown has internode durations  $t_k$  during which there are  $k$  lineages consistent with the SFS (the topology was chosen uniformly among ranked binary trees with  $2n$  tips). Time is given in coalescent 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.

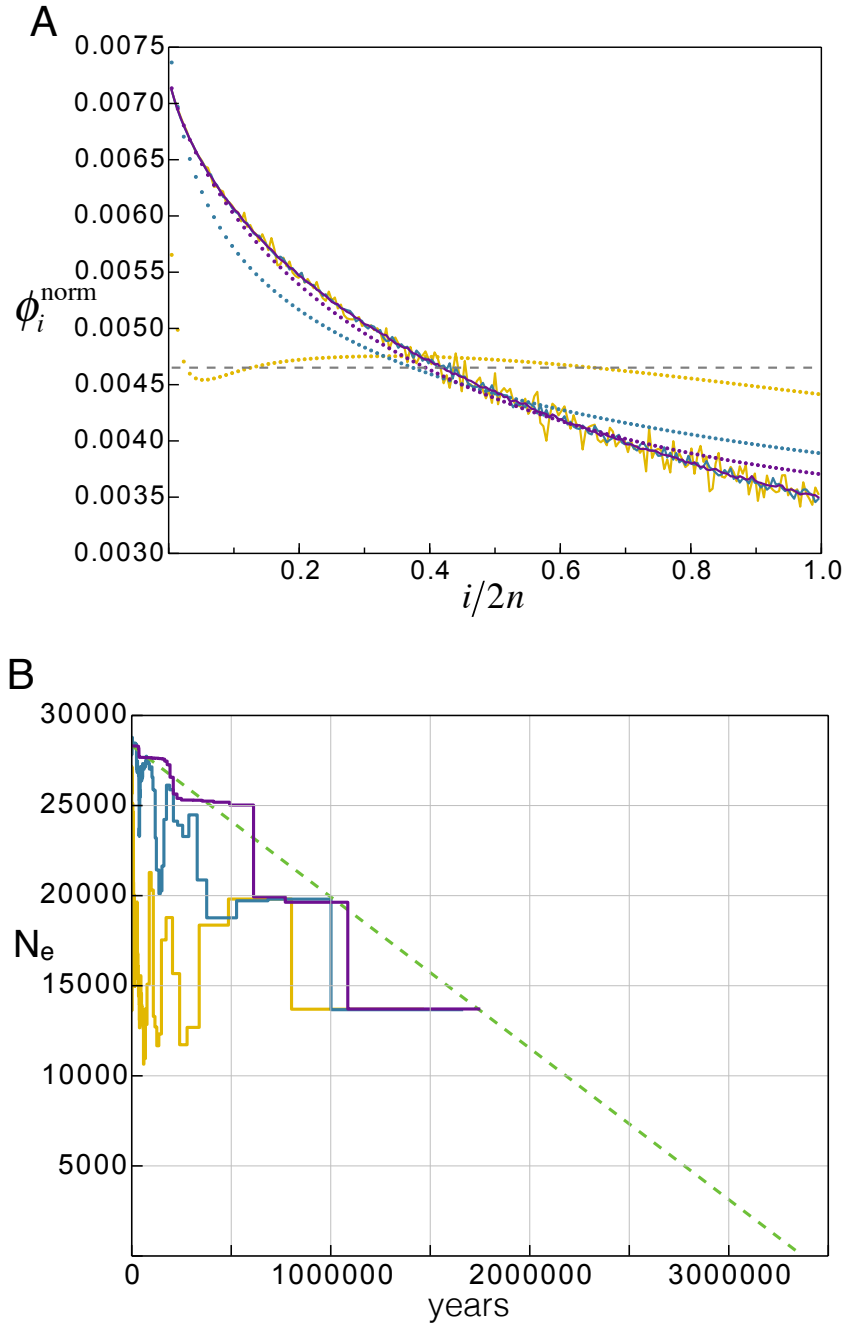


Figure 5: Stairway plot inference of a linear demography SFS with noise. A) Solid lines: mean of 200 SFS simulated independently under the *Linear* growth model, with either 10<sup>5</sup> loci (purple), 10<sup>4</sup> loci (blue) or 10<sup>3</sup> loci (yellow). Dotted lines: expected SFS under the 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 on the 200 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  $d$  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 foundation time  $\hat{\tau}$  under the *Linear* model on SFS with noise. Mean, 5% and 95% percentile of the foundation time inferred with a *Linear* model. The SFS on which the inference is made are simulated with a foundation time  $\tau$  of 2.48, with different number of loci, using the method with topology reconstruction.