

1 **Genomic inference of a human super bottleneck in the Early Stone Age**

2

3 Wangjie Hu^{1,†}, Ziqian Hao^{1,†}, Pengyuan Du¹, Yi-Hsuan Pan^{2,*}, Haipeng Li^{1,3,*}

4

5 ¹CAS Key Laboratory of Computational Biology, Shanghai Institute of Nutrition and
6 Health, University of Chinese Academy of Sciences, Chinese Academy of Sciences,
7 Shanghai 200031, China.

8 ²Key Laboratory of Brain Functional Genomics of Ministry of Education, School of
9 Life Science, East China Normal University, Shanghai 200062, China.

10 ³Center for Excellence in Animal Evolution and Genetics, Chinese Academy of
11 Sciences, Kunming 650223, China.

12

13 †These authors contributed equally.

14

15 *Corresponding Authors: yxpan@sat.ecnu.edu.cn; lihaipeng@picb.ac.cn

16

17 Short title: Human super bottleneck

18

19

20 **Abstract**

21 The demographic history has been a foundation of human evolutionary studies for
22 more than a century. In this study, we developed a novel method referred to as the fast
23 infinitesimal time coalescent (FitCoal) process. This method allows the accurate
24 calculation of the composite likelihood of a site frequency spectrum and provides the
25 precise inference of recent and ancient demographic history. Genomic sequences of
26 the 1000 Genomes Project and the Human Genome Diversity Project – Centre
27 d’Etude du Polymorphisme Humain panel were analyzed. Results showed that all ten
28 African populations had a population super bottleneck, a small effective size of
29 approximately 1,280 breeding individuals between 813 and 930 thousand years ago,
30 and a 20-fold rapid growth at the end of the bottleneck. The super bottleneck caused a
31 loss of 65.85% in current human genetic diversity, but it may have separated our
32 ancestors from other hominins. Further analysis confirmed the existence of the super
33 bottleneck in all 40 non-African populations. Our results provide new insights into
34 human evolution in the Early Stone Age.

35

36

37 Inferring demographic history from genomic information has played an important
38 role in population genetics. It uncovers prehistoric evolutionary events and deepens
39 our understanding about the evolution of human and other species³⁻⁷. Multiple
40 methods have been developed to infer demographic history with a predefined
41 demographic model⁸⁻¹². These methods require prior knowledge about the species
42 being investigated and estimation parameters by fitting in summary statistics such as
43 site frequency spectrum (SFS). In contrast, model-free methods do not need a
44 predefined model when inferring demography¹³⁻¹⁸. As SFS plays an essential role in
45 demographic inference, many efforts have been attempted to derive its analytical
46 formula under certain demographic models¹⁹⁻²¹.

47 To precisely infer recent and ancient demography, we developed the fast
48 infinitesimal time coalescent (FitCoal) process (Fig. 1) that calculates expected
49 branch length for each SFS type under arbitrary demographic models. It is effective
50 for a wide range of sample sizes in the calculation of the composite likelihood of a
51 given SFS^{8,9}. To infer the demographic history, FitCoal first maximizes the likelihood
52 with the constant size model and then increases the number of inference time intervals
53 and re-maximizes the likelihood until the best model is found. As inference time
54 intervals are variable to avoid long fixed ancient time intervals, the precision in the
55 inference of ancient demographic events can be improved. FitCoal does not need prior
56 information on demography, and its accuracy can be confirmed by simulation.

57 With African hominid fossils, the origin of anatomically modern humans has been
58 determined to be approximately 200 thousand years (kyr) ago²². Although the
59 demographic history of anatomically modern humans has been intensively
60 studied^{14,15,17,23-26}, it is conceivable that many new insights remain to be explored. In
61 this study, we used FitCoal to analyze genomic sequences of the 1000 Genomes
62 Project phase 3 (1000GP)¹ and the Human Genome Diversity Project–Centre d’Etude
63 du Polymorphisme Humain panel (HGPD-CEPH)². Results revealed a super
64 bottleneck in all 10 African populations between 813 and 930 kyr ago with an
65 effective population size of 1,280. According to the coalescent theory and simulations,

66 the super bottleneck cannot be directly inferred in the non-African populations
67 because fewer coalescent events remained in the non-African populations during the
68 bottleneck period. Instead, a hidden effect of the bottleneck was found in all 40
69 non-African populations. Our results suggest that our ancestors experienced a super
70 bottleneck and the effective size of our ancestors remained small for a very long
71 period of time (~117,000 years). The super bottleneck may have separated our
72 ancestors from other hominins^{27,28}.

73

74 **Fast infinitesimal time coalescent process**

75 As determination of expected branch length for each SFS type is essential for
76 theoretical population genetics and demographic inference^{8,9}, we developed the fast
77 infinitesimal time coalescent (FitCoal) process to accomplish the task (Fig. 1). For
78 FitCoal analysis, each of millions of time intervals Δt was set extremely small, and
79 the population size was assumed to be constant within each infinitesimal time interval.
80 The probabilities of all states were calculated backward in time. During each Δt , the
81 branches were categorized according to their state. For each state, the branch length
82 was multiplied by its probability and population size and then transformed to calculate
83 the expected branch length of each SFS type. Because the expected branch length of a
84 SFS type is equal to the sum of the expected branch length of this type during each
85 time interval, the latter can be rescaled and tabulated, making the calculation of the
86 expected branch lengths extremely fast under arbitrary demographic histories.
87 Hereafter, tabulated FitCoal is referred to as FitCoal for short, unless otherwise
88 indicated.

89

90 **FitCoal demographic inference**

91 After the expected branch lengths were obtained, the composite likelihood of the
92 SFS observed in a sample was calculated^{8,9,15,29}. As each single nucleotide
93 polymorphism (SNP) was treated independently, FitCoal did not need phased
94 haplotype data. When inferring demography, the likelihood was maximized in a wide
95 range of demographic scenarios. The FitCoal likelihood surface is smooth (Fig. S1),
96 so it is efficient to maximize the likelihood. FitCoal considered both instantaneous
97 populations size changes¹⁴⁻¹⁶ and long-term exponential changes of population in

98 order to generate various demographic scenarios.

99

100 **Demographic inference on simulated data**

101 The accuracy of FitCoal was validated by simulation and comparing its
102 demographic inferences with those of PSMC¹⁴ and stairway plot¹⁵ (Fig. 2). Six
103 different demographic models, examined in the study by Liu and Fu¹⁵, were
104 considered by simulating 200 independent data sets under each model. The medians
105 and 95% confidence intervals of demography were then determined by FitCoal with
106 the assumption that a generation time is 24 years^{15,30} and the mutation rate is
107 1.2×10^{-8} per site per generation for human populations^{15,31-33}. Our results (Fig. 2)
108 confirmed that SFS allows precise recovery of the demographic history³⁴.

109 FitCoal was found to precisely infer demographic histories (Fig. 2), and the
110 inference accuracy was improved by increasing sample size and length of sequence
111 (Fig. S2). In general, the confidence intervals of FitCoal inferred histories were
112 narrower than those of PSMC and stairway plot, except for those with insufficient
113 information on ancestral populations (Fig. 2c). The proportion of the most recent
114 change type inferred from the six different models mentioned above also showed that
115 FitCoal can distinguish instantaneous and exponential changes (Table S1).

116 Since a demographic event may affect every SFS type, demographic history can
117 be inferred using a subset of SFS. Results of simulation confirmed that FitCoal
118 accurately determined demographic history based on truncated SFSs (Fig. S3, 4), thus
119 reducing the impact of other factors, such as positive selection (Fig. S5) and
120 sequencing error, on FitCoal analysis.

121

122 **Demographic inference of African populations**

123 To infer the demographic histories of African populations, seven African
124 populations in the 1000GP¹ were analyzed by FitCoal. Only non-coding regions,
125 defined by GENCODE³⁵, were used in order to avoid the effect of purifying selection.
126 To avoid the potential effect of positive selection³⁶, high-frequency mutations were
127 excluded from the analysis.

128 Results showed that all seven African populations had a super bottleneck around
129 914 (854–1,003) kyr ago and that this bottleneck was relieved about 793 (772–815)
130 kyr ago (Fig. 3a-c, S6; Table S2). The average effective population size of African

131 populations during the bottleneck period was determined to be 1,270 (770–2,030).
132 Although traces of the bottleneck were observed in previous studies, the bottleneck
133 was ignored because its signatures were too weak to be noticed^{1,2,14,16,17}. After the
134 bottleneck was relieved, the population size was increased to 27,080 (25,300–29,180),
135 a 20-fold increase, around 800 kyr ago. This population size remained relatively
136 constant until the recent expansion.

137 To avoid the potential effects of low sequencing depth (~ 5x) of non-coding
138 regions in the 1000GP on the analysis, the autosomal non-coding genomic
139 polymorphism of HGDP-CEPH data set with high sequencing coverage (~35x) was
140 used. In total, populations with more than 15 individuals each were examined. Results
141 showed that the super bottleneck occurred between 859 (856–864) and 1,257 (1,042–
142 1,527) kyr ago in all three African populations in HGDP-CEPH (Fig. 3d-f, S7; Table
143 S3), and the average population size during the bottleneck period was 1,300 (908–
144 1,670). This number was very similar to that (1,270) estimated from the data of
145 1000GP.

146 After the bottleneck was relieved, the population sizes of the two HGDP-CEPH
147 agriculturalist populations were increased to 27,300 and 27,570 (Fig. 3e, S7; Table
148 S3), consistent with the 1000GP estimate of 27,280. However, the Biaka, a
149 hunter-gatherer population, had a larger population size of 35,330, suggesting a deep
150 divergence between this and other agriculturalist populations³⁷⁻³⁹. The Biaka
151 population was found to have a recent population decline (Fig. 3d, S7), as previously
152 observed². These results suggest that hunter-gatherer populations were widely spread
153 and decreased when agriculturalist populations were expanded.

154 To provide a precise inference of the super bottleneck, the results from the two
155 data sets were combined. After analyzing the inferred time of instantaneous change of
156 ten populations, the super bottleneck was inferred to last for about 117,000 years,
157 from 813 (772–864; s.e.m.: 11.02) to 930 (854–1,042; s.e.m.: 23.52) kyr ago. The
158 effective size during the bottleneck period was precisely determined to be 1,280 (767–
159 2,031; s.e.m.: 131). A loss of 65.85% in current genetic diversity of human
160 populations was estimated because of the bottleneck.

161

162 **Demographic inference of non-African populations**

163 No super bottleneck was directly observed in all 19 non-African populations in
164 1000GP (Fig. 3a-c, S6; Table S4). The ancestral population size of these populations
165 was determined to be 20,260 (18,850–22,220), similar to that determined in previous
166 studies^{2,14,16,17}. The population size started to decline around 368 (175–756) kyr ago in
167 1000GP non-African populations, suggesting that African and non-African
168 divergence occurred much earlier than the out-of-African migration^{1,2,14,16,17,24}.
169 European and South Asian populations were found to have a relatively weaker
170 bottleneck than East Asian populations, and the bottleneck severity was found to
171 correlate with their geographic distance to Africa, consistent with the observed
172 correlation between heterozygosity and geographic distance^{40,41}. A weak bottleneck
173 was observed in American populations, probably because of recent admixture¹. All
174 1000GP non-African populations were found to increase in size recently.

175 The super bottleneck was also not directly detected in all 21 HGDP-CEPH
176 non-African populations (Fig. 3d-f, S7; Table S5). The ancestral population size of
177 these populations was determined to be 20,030 (19,060–21,850), very similar to that
178 (20,260) estimated from 1000GP. These populations started to decline 367 (167–628)
179 kyr ago. A positive correlation was also observed between the severity of
180 out-of-African bottleneck and their geographic distance to Africa. The Middle East
181 populations had the weakest bottleneck, while the Maya, an American population, had
182 the strongest bottleneck. Similar to 1000GP non-African populations, most
183 HGDP-CEPH non-African populations were found to increase in size recently, except
184 an isolated Kalash population, consistent with previous studies^{2,42}.

185

186 **Super bottleneck in the Early Stone Age**

187 Although a super bottleneck was detected in all 10 African populations, such
188 bottleneck was not directly detected in all 40 non-African populations. To investigate
189 this phenomenon, simulations were performed with three 1000GP demographic
190 models, designated Bottleneck I, II, and III (Fig. 4). Bottleneck I simulated the
191 average inferred demographic history of African populations with the super
192 bottleneck, and Bottleneck II and III simulated the demography of non-African
193 populations without and with the super bottleneck. Both Bottleneck I and II were
194 inferred correctly as a corresponding bottleneck was found in all simulated data sets
195 (Table S6). However, no super bottleneck was detected in Bottleneck III simulations.

196 The super bottleneck was found to cause a population size gap between the true model
197 and inferred demographic history, suggesting a hidden effect of the super bottleneck
198 on non-African populations. Simulations were then extended to HGDP-CEHP
199 populations with Bottleneck models IV–VI, and similar results were obtained (Fig. S8;
200 Table S7). When simulations were performed on three artificial models (Bottleneck
201 VII–IX) with various demographic parameters, a population size gap was still
202 detected (Fig. S9; Table S8). These results suggest a hidden effect of the super
203 bottleneck on non-African populations.

204 The population size gap was found in both 1000GP and HGDP-CEPH data sets
205 (Fig. 3a,d). The average population sizes of non-African populations were determined
206 to be 20,260 and 20,030, respectively, while those of African agriculturalist
207 populations were 27,080 and 27,440, respectively in these two data sets. The observed
208 population size gap was 7,020, probably due to the hidden effect of the super
209 bottleneck on non-African populations.

210 The reasons were then investigated why the super bottleneck had different effects
211 on African and non-African populations. Results showed that non-African populations
212 had the out-of-African bottleneck, but African populations lacked such bottleneck.
213 Therefore, the standard coalescent time of non-African populations was larger than
214 that of African populations (Fig. 3c,f). As African populations had more coalescent
215 events occurred during the bottleneck period, the bottleneck was more readily
216 detected. The mathematical proof on this issue was described in the supplemental
217 material.

218

219 **Discussion**

220 In this study, we develop FitCoal, a novel model-flexible tool for demographic
221 inference. Key characteristic features of FitCoal are that the composite likelihood can
222 be rapidly calculated based on expected branch lengths and that inference time
223 intervals are variable during the demographic inference. Since coalescent events
224 become rare when tracing backward in time, the length of time interval is usually set
225 to increase progressively¹⁴⁻¹⁷. Although this strategy can capture recent demographic
226 events, it may miss ancient ones. As FitCoal uses variable time intervals, it can give
227 an accurate inference for both recent and ancient demographic events.

228 The most important discovery with FitCoal in this study is that human ancestors
229 had a super bottleneck. The ancient population size reduction around 930 kyr ago was
230 likely to be due to the formation of geographically isolated populations driven by
231 culture or ecological factors. In addition, we found that our ancestors had a very small
232 effective size of approximately 1,280 breeding individuals during the bottleneck
233 period. This number is comparable in the same magnitude in the effective size of
234 mammals threatened by extinction⁴³. We also detected an instantaneous recovery in
235 all ten African populations with a 20-fold population growth during a short time
236 period around 813 kyr ago (Fig. 3, S10). The earliest archaeological evidence for
237 human control of fire was found in Israel 790 kyr ago^{44,45}. As the control of fire
238 profoundly affected social evolution⁴⁶ and brain size⁴⁷, it may be associated with the
239 big bang in population size at the end of the super bottleneck.

240 The super bottleneck, which started about one million years ago, might be
241 associated with a speciation event in the early human evolution^{27,28}. The questions
242 about where the small ancient population dwelt, how they survived for such a long
243 time, and how they diverged from other hominin groups remain to be investigated. As
244 Neandertals and Denisovans diverged with the modern human between 270 and 440
245 kyr ago^{48,49}, it is conceivable that they also had the super bottleneck. In the future, a
246 more detailed picture of human evolution in the Early Stone Age may be revealed
247 when more genomic sequences of African populations and archaic hominins and more
248 advanced population genomics methods become available.

249

250 **Methods**

251 **Standard coalescent time and time in generations.** The population size is denoted
252 $N(\cdot)$, representing the demographic history. Time τ represents one-point scaled time
253 since the time in a generation is scaled by $2N(0)$. Time t is usually scaled by
254 $2N(t)$ generations^{20,34,50,51}. To distinguish it from the one-point scaled time τ , time t
255 is designated as the standard coalescent time.

256

257 **Fast infinitesimal time coalescent (FitCoal) process.** The FitCoal calculates the
258 expected branch length for each type of site frequency spectrum (SFS) under arbitrary
259 demographic history $N(\cdot)$. We assume that a sample is obtained by randomly taken n

260 sequences from the population. The sample is designated to be state l ($l = 2, \dots, n$)
 261 at time t if it has exactly l ancestral lineages at this time. The probability of state l
 262 at time t is denoted $p_l(t)$. In a coalescent tree, a branch is designated to be type i if
 263 it has exactly i descendants. We have

$$264 \quad \frac{d}{dt} p_l(t) = \begin{cases} \binom{l+1}{2} p_{l+1}(t) - \binom{l}{2} p_l(t) & l = 2, \dots, n-1 \\ -\binom{l}{2} p_l(t) & l = n \end{cases}.$$

265 When Δt is extremely small (Fig. 1), there is at most one coalescent event during t
 266 and $t + \Delta t$, leading to

$$267 \quad p_l(t + \Delta t) = \begin{cases} \binom{l+1}{2} \Delta t p_{l+1}(t) + (1 - \binom{l}{2} \Delta t) p_l(t) & l = 2, \dots, n-1 \\ (1 - \binom{l}{2} \Delta t) p_l(t) & l = n \end{cases}.$$

268 The branch length is in units of generations. The expected branch length of state
 269 l during t and $t + \Delta t$ is calculated as $\int_t^{t+\Delta t} 2N(t) p_l(t) l dt$. The probability that a
 270 branch of state l is of type i is $\frac{\binom{n-i-1}{l-2}}{\binom{n-1}{l-1}}$ (ref.²⁰). The expected branch length of type

271 i of state l during t and $t + \Delta t$ is $\int_t^{t+\Delta t} 2N(t) p_l(t) l \frac{\binom{n-i-1}{l-2}}{\binom{n-1}{l-1}} dt$. Therefore, the
 272 expected branch length $BL_i(N(\cdot))$ of type i is

$$273 \quad \sum_{l=2}^{n-i+1} \int_0^\infty 2N(t) p_l(t) l dt \frac{\binom{n-i-1}{l-2}}{\binom{n-1}{l-1}}.$$

274 A FitCoal time partition is denoted by $\{t_0, t_1, \dots, t_m\}$, where $0 = t_0 < t_1 <$
 275 $\dots < t_m$. We have $p_l(t_0) = \begin{cases} 1 & l = n \\ 0 & \text{else} \end{cases}$. For a large positive number m , if t_m is
 276 large and $(t_k - t_{k-1})$ is small for $k = 1, \dots, m$, then

$$277 \quad p_l(t_k) = \begin{cases} (1 - \binom{l}{2}(t_k - t_{k-1})) p_l(t_{k-1}) & l = n \\ (1 - \binom{l}{2}(t_k - t_{k-1})) p_l(t_{k-1}) + \binom{l+1}{2}(t_k - t_{k-1}) p_{l+1}(t_{k-1}) & \text{else} \end{cases},$$

278 where $k = 1, \dots, m$.

279 The expected branch length of type i is calculated as

$$280 \quad BL_i(N(\cdot)) = \sum_{l=2}^{n-i+1} l \frac{\binom{n-i-1}{l-2}}{\binom{n-1}{l-1}} (\sum_{k=1}^m 2N(t_{k-1}) p_l(t_{k-1}) (t_k - t_{k-1})).$$

281 To determine the time partition, we required that the coalescent probability was
 282 less than 10^{-4} during t_{k-1} and t_k ($k = 1, \dots, m$), the probability of common
 283 ancestor (*i.e.*, the probability of state 1) at t_m was larger than $(1 - 10^{-6})$. When the
 284 sample size was 10, the number of infinitesimal time intervals was 1,571,200. When
 285 the sample size was 200, the number of infinitesimal time intervals was 7,038,398.

286 Thus, each Δt was extremely small for precise calculation of expected branch length,
 287 and the time was partitioned to obtain $p_i(t)$ in order to calculate the expected branch
 288 length of type i .

289

290 **Tabulated FitCoal.** The expected branch length of each type can be calculated for
 291 arbitrary time intervals according to the procedure described above. Considering
 292 another tabulated time partition $\{t_0, t_1, \dots, t_m\}$ ($0 = t_0 < t_1 < \dots < t_m$), the expected
 293 branch length of a type is equal to the sum of the expected branch length of this type
 294 during each tabulated time interval, thus the latter can be rescaled and tabulated.

295 The scaled expected branch length $BL_{i,t}$ of type i during 0 and t is

296
$$BL_{i,t} = \sum_{l=2}^{n-i+1} \int_0^t p_s(l) l \frac{\binom{n-i-1}{l-2}}{\binom{n-1}{l-1}} ds, \text{ where } i = 1, \dots, n-1. \text{ For the tabulated time}$$

297 partition $\{t_0, t_1, \dots, t_m\}$, BL_{i,t_0} , BL_{i,t_1} , \dots , and BL_{i,t_m} are tabulated. When $n = 10$,
 298 $m = 231$. When $n = 200$, $m = 529$.

299 $BL_{i,t}$ is used to calculate the expected branch lengths under arbitrary
 300 demographic histories. When $\tilde{t} \in [t_{k-1}, t_k)$,

301
$$BL_{i,\tilde{t}} \approx \frac{t_k - \tilde{t}}{t_k - t_{k-1}} BL_{i,t_{k-1}} + \frac{\tilde{t} - t_{k-1}}{t_k - t_{k-1}} BL_{i,t_k}.$$

302 If $N(t)$ is a piecewise constant, that is, there exist a demographic time partition
 303 $\{\tilde{t}_0, \tilde{t}_1, \dots, \tilde{t}_{\tilde{m}}\}$, $N(t) = N_k$ for $t \in [\tilde{t}_k, \tilde{t}_{k+1})$, $k = 0, \dots, \tilde{m}$. Then, the expected
 304 branch length of type i is calculated as

305
$$BL_i(N(\cdot)) = \sum_{k=1}^{\tilde{m}} 2N_k (BL_{i,\tilde{t}_k} - BL_{i,\tilde{t}_{k-1}}).$$

306 When $N(t)$ is complex, the population size can be approximated by a piecewise
 307 constant function.

308

309 **Composite likelihood.** The mutation rate per base pair per generation is denoted μ ,
 310 and $\vec{\xi} = (\xi_i)$ is the observed number of SNPs of n sequences with σ base pairs,
 311 where $i = 1, \dots, n-1$. The expected SFS is $\vec{\lambda} = (\lambda_i)$, where $\lambda_i = \mu\sigma BL_i(N(\cdot))$.
 312 Following the Poisson probability and previous studies⁸, the composite likelihood is
 313 calculated as follows:

314
$$L_{\mu,i}(\vec{\xi}, N(\cdot)) = \prod_{i=1}^{n-1} \frac{\lambda_i^{\xi_i} e^{-\lambda_i}}{\xi_i!}.$$

315 The likelihood is extended to missing data and truncated SFS (see Supplemental
 316 materials).

317

318 **Demographic inference.** The number of demographic time intervals is variable.
 319 FitCoal first fits the observed SFS using a constant size model with one demographic
 320 time interval, and the number of time intervals is increased by one at a time to
 321 generate more complex models. The Local Unimodal Sampling (LUS) algorithm⁵² is
 322 used to maximize the likelihood and estimate demographic parameters. A
 323 log-likelihood promotion rate is used to determine the best model to explain the
 324 observed SFS, and 20% is used as the threshold.

325 A series of demography with m pieces is denoted by a set $S(m)$, where $S(m)$
 326 contains all of the following m pieces of population size:

$$327 \quad N(t|N_0 > 0, N_{(m)}, t_{(m)}, c_{(m)}) = \begin{cases} N_m N_0 & t \geq t_m \\ N_k N_0 & t_k \leq t < t_{k+1}, c_k \in \mathcal{C}, k = 1, \dots, m-1, \\ \frac{(t_{k+1}-t_k)N_{k+1}N_k N_0}{(t-t_k)N_k + (t_{k+1}-t)N_{k+1}} & t_k \leq t < t_{k+1}, c_k \in \mathcal{E}, k = 1, \dots, m-1 \end{cases}$$

328 where $N_{(m)} = (N_1, \dots, N_m) \in N[m]$, $t_{(m)} = (t_1, \dots, t_m) \in t[m]$,
 329 $c_{(m)} = (c_1, \dots, c_m) \in c[m]$, $N[m] = \{(N_1, \dots, N_m) | N_1 = 1, N_i > 0 \text{ for } i > 1\}$,
 330 $t[m] = \{(t_1, \dots, t_m) | 0 = t_1 < \dots < t_m\}$, $c[m] = \{(c_1, \dots, c_m) | c_m \in \mathcal{C}, c_i \in \mathcal{C} \cup$
 331 $\mathcal{E} \text{ for } i = 1, \dots, m-1\}$, $\mathcal{C} = \{\text{constant}\}$, and $\mathcal{E} = \{\text{exponential}\}$.

332 The set $S(m)$ was used as the wide-range parameter space to determine the
 333 maximum likelihood. To find the best demographic history to explain the observed
 334 SFS, the following procedures were used:

- 335 (1) The number of inference time intervals (or pieces) m is initially set to 1, and the
 336 maximum likelihood $\max L_1$ is determined with the constant size model (model in
 337 $S(1)$).
- 338 (2) Increase m by 1. For each change of type $c_{(m)}$, parameters $N_{(m)} = (N_1, \dots, N_m)$
 339 and $t_{(m)} = (t_1 = 0, t_2, \dots, t_m)$ are searched to maximize the likelihood by LUS
 340 algorithm to fit the observed SFS. The maximum likelihood $\max L_m$ is calculated
 341 with models in $S(m)$ with all possible change types.
- 342 (3) Repeat step (2) until $(1 + \text{threshold}) \cdot \log(\max L_m) < \log(\max L_{m-1})$ is
 343 obtained. The best model corresponding $\max L_{m-1}$ is determined to explain the
 344 observed SFS.
- 345 (4) To avoid local optima, steps (1) – (3) are repeated K times to find the best model.
 346 $K = 10$ when analyzing simulated samples, and $K = 200$ when analyzing the

347 observed SFSs of the 1000GP and HGDP-CEPH populations.

348 To determine the threshold of log-likelihood promotion rate, a large number of
349 simulations were performed (Table S9). For each model, 200 replicates were
350 conducted, and the number of inference time intervals in the estimated demographic
351 history was determined for each replicate. If the estimated number of inference time
352 intervals was larger than the true number of inference time intervals, overfitting was
353 recorded. When the former was smaller than the latter, underfitting was considered.
354 The thresholds of 10%, 20%, and 30% were used. When 10% was used, the maximum
355 overfitting rate was 2%. When 20% was used, all cases examined were inferred
356 correctly. When 30% was used, the underfitting was observed in one of 20 examined
357 models. Therefore, 20% was used as the threshold of log-likelihood promotion rate in
358 subsequent analyses.

359

360 **Data simulation.** Data were simulated using *ms*⁵³ and *MaCS*⁵⁴ software. Unless
361 otherwise specified, a generation time was assumed to be 24 years^{15,30}, the mutation
362 rate μ was set for 1.2×10^{-8} per base per generation^{15,31-33}, and the recombination
363 rate was $r = 0.8\mu$. For each model, 200 SFSs were simulated to calculate the median
364 and 2.5 and 97.5 percentiles. When verifying the inferred demographic histories,
365 80,000 DNA fragments with the length of 10kb each were used for simulation, taking
366 into the consideration of small fragments split by sequencing mask in 1000GP and
367 HGDP-CEPH data sets. High frequency alleles of SFS (10% mutation types for
368 Bottleneck I, II, III, VII, VIII, IX, and 15% for Bottleneck IV, V, VI) were removed
369 when assessing models to verify the super bottleneck. Detailed simulation command
370 lines and demographic inference are presented in the Supplementary information.

371

372 **1000 Genomes Project data.** Sequences of autosomal SNPs in 1000GP phase 3 (ref.¹)
373 were downloaded from the 1000GP ftp server
374 (<ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/>), and 26 populations
375 were analyzed, including seven African populations (ACB, ASW, ESN, GWD, LWK,
376 MSL, and YRI), five European populations (CEU, FIN, GBR, IBS, and TSI), five
377 East Asian populations (CDX, CHB, CHS, JPT, and KHV), five South Asian
378 populations (BEB, GIH, ITU, PJL, and STU), and four American populations (CLM,
379 MXL, PEL, and PUR). The 1000 GP strict mask was used to exclude artifacts of SNP
380 calling. Noncoding regions except pseudogenes, defined by GENCODE release 35

381 (ref.³⁵), were examined to avoid potential effects of purifying selection. Sites without
382 high-confidence ancestral allele inference, according to 1000GP annotations, were
383 excluded. The number of bi-allelic sites that passed the filtering was 826,649,529. To
384 avoid the effect of positive selection, high frequency mutations were excluded, and
385 the truncated SFS was used to infer demographic history (Fig. S11; Table S10). The
386 average proportion of excluded high-frequency SNPs for all 1000GP populations was
387 4.40%.

388

389 **HGDP-CEPH data.** In total, 24 populations were analyzed, including three African
390 populations (Biaka, Mandeka, and Yoruba), five European populations (Adygei,
391 Basque, French, Russian, and Sardinian), four Middle East populations (Bedouin,
392 Druze, Mozabite, and Palestinian), three East Asian populations (Han, Japanese, and
393 Yakut), eight Central and South Asian populations (Balochi, Brahui, Burusho, Hazara,
394 Kalash, Makrani, Pathan, and Sindhi), and an American population (Maya). Only
395 bi-allelic SNPs locating in GENCODE non-coding regions³⁵ except pseudogenes that
396 passed HGDP-CEPH filtering were used. HGDP-CEPH accessible mask was also
397 used to filter SNPs². The number of sites that passed filtering was 791,999,125.
398 Missing data were allowed to avoid artifacts due to imputation. The proportion of
399 sites with two or more missing individuals was less than 3% for all populations (Table
400 S11). Each population had two SFSs, with one calculated from sites with no missing
401 data, and another from sites with one missing individual. Similarly, truncated SFSs
402 were used to avoid the effect of positive selection (Fig. S12, S13; Table S12). The
403 average proportion of excluded high-frequency SNPs for all HGDP-CEPH
404 populations was 7.18%.

405

406 **Data availability.**

407 The authors declare that all data supporting the findings of this study are included in
408 this paper and its supplementary information file.

409

410 **Code availability.**

411 FitCoal is a free plug-in of the eGPS software⁵⁵ and can be downloaded and run as an
412 independent package. FitCoal and its documentation are available via Zenodo at
413 <https://zenodo.org/record/4765447#.YKDt7aG-vuq>, our institute website at
414 <http://www.picb.ac.cn/evolgen/>, and eGPS website <http://www.egps-software.net/>.

415

416 **Acknowledgement**

417 We thank Daniel Zivković for sharing his codes to calculate the expected branch
418 length, and Xiaoming Liu for sharing his simulated results. This work was supported
419 by grants from the Strategic Priority Research Program of the Chinese Academy of
420 Sciences (XDB13040800), the National Natural Science Foundation of China (nos.
421 31100273, 31172073, 91131010), and National Key Research and Development
422 Project (No. 2020YFC0847000).

423

424 **Author contributions:** W.H., Z.H., Y.H.P., and H.L. conceived and designed the
425 research; W.H., Z.H., and H.L. wrote the code; W.H., Z.H., P.D., and Y.H.P. analyzed
426 the data; W.H., Z.H., P.D., Y.H.P., and H.L. wrote the paper.

427

428 **Competing interests:** All authors declare that they have no competing interests.

429

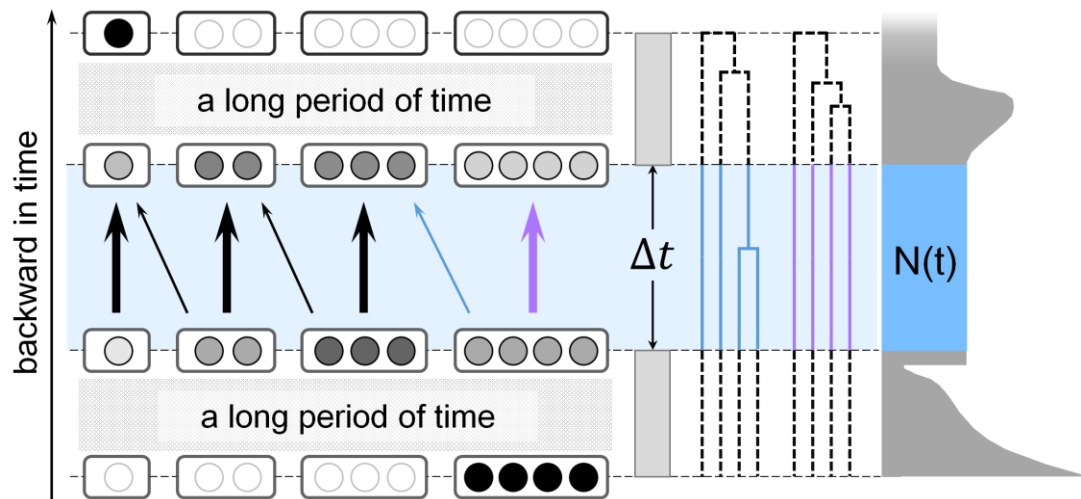
430 **References**

- 431 1 Altshuler, D. M. *et al.* A global reference for human genetic variation. *Nature*
432 **526**, 68-74 (2015).
- 433 2 Bergstrom, A. *et al.* Insights into human genetic variation and population
434 history from 929 diverse genomes. *Science* **367**, eaay5012 (2020).
- 435 3 Hu, J. Y. *et al.* Genomic consequences of population decline in critically
436 endangered pangolins and their demographic histories. *Natl. Sci. Rev.* **7**,
437 798-814 (2020).
- 438 4 Zeng, L. *et al.* Out of southern East Asia of the brown rat revealed by
439 large-scale genome sequencing. *Mol. Biol. Evol.* **35**, 149-158 (2018).
- 440 5 Fu, Q. *et al.* Genome sequence of a 45,000-year-old modern human from
441 western Siberia. *Nature* **514**, 445-449 (2014).
- 442 6 Raghavan, M. *et al.* Genomic evidence for the Pleistocene and recent
443 population history of Native Americans. *Science* **349**, aab3884 (2015).
- 444 7 Jinam, T. A. *et al.* Evolutionary history of continental southeast Asians: "early
445 train" hypothesis based on genetic analysis of mitochondrial and autosomal
446 DNA data. *Mol. Biol. Evol.* **29**, 3513-3527 (2012).
- 447 8 Li, H. & Stephan, W. Inferring the demographic history and rate of adaptive
448 substitution in *Drosophila*. *PLoS Genet.* **2**, e166 (2006).
- 449 9 Excoffier, L., Dupanloup, I., Huerta-Sanchez, E., Sousa, V. C. & Foll, M.
450 Robust demographic inference from genomic and SNP data. *PLoS Genet.* **9**,
451 e1003905 (2013).
- 452 10 Griffiths, R. C. & Tavaré S. Monte Carlo inference methods in population

- 453 genetics. *Math. Comput. Modell.* **23**, 141-158 (1996).
- 454 11 Gutenkunst, R. N., Hernandez, R. D., Williamson, S. H. & Bustamante, C. D.
455 Inferring the joint demographic history of multiple populations from
456 multidimensional SNP frequency data. *PLoS Genet.* **5**, e1000695 (2009).
- 457 12 Schaffner, S. F. *et al.* Calibrating a coalescent simulation of human genome
458 sequence variation. *Genome Res.* **15**, 1576-1583 (2005).
- 459 13 Heled, J. & Drummond, A. J. Bayesian inference of population size history
460 from multiple loci. *BMC Evol. Biol.* **8**, 289 (2008).
- 461 14 Li, H. & Durbin, R. Inference of human population history from individual
462 whole-genome sequences. *Nature* **475**, 493-496 (2011).
- 463 15 Liu, X. M. & Fu, Y. X. Exploring population size changes using SNP
464 frequency spectra. *Nat. Genet.* **47**, 555-559 (2015).
- 465 16 Schiffels, S. & Durbin, R. Inferring human population size and separation
466 history from multiple genome sequences. *Nat. Genet.* **46**, 919-925 (2014).
- 467 17 Terhorst, J., Kamm, J. A. & Song, Y. S. Robust and scalable inference of
468 population history from hundreds of unphased whole genomes. *Nat. Genet.* **49**,
469 303-309 (2017).
- 470 18 Liu, X. & Fu, Y. X. Stairway Plot 2: demographic history inference with
471 folded SNP frequency spectra. *Genome biology* **21**, 280 (2020).
- 472 19 Jouganous, J., Long, W., Ragsdale, A. P. & Gravel, S. Inferring the joint
473 demographic history of multiple populations: beyond the diffusion
474 approximation. *Genetics* **206**, 1549-1567 (2017).
- 475 20 Fu, Y. X. Statistical properties of segregating sites. *Theor. Popul. Biol.* **48**,
476 172-197 (1995).
- 477 21 Zivković, D. & Wiehe, T. Second-order moments of segregating sites under
478 variable population size. *Genetics* **180**, 341-357 (2008).
- 479 22 White, T. D. *et al.* Pleistocene *Homo sapiens* from Middle Awash, Ethiopia.
480 *Nature* **423**, 742-747 (2003).
- 481 23 Stoneking, M. & Krause, J. Learning about human population history from
482 ancient and modern genomes. *Nat. Rev. Genet.* **12**, 603-614 (2011).
- 483 24 Nielsen, R. *et al.* Tracing the peopling of the world through genomics. *Nature*
484 **541**, 302-310 (2017).
- 485 25 Manica, A., Amos, W., Balloux, F. & Hanihara, T. The effect of ancient
486 population bottlenecks on human phenotypic variation. *Nature* **448**, 346-348
487 (2007).
- 488 26 Ramachandran, S. *et al.* Support from the relationship of genetic and
489 geographic distance in human populations for a serial founder effect
490 originating in Africa. *Proc. Natl. Acad. Sci. USA* **102**, 15942-15947 (2005).
- 491 27 Stringer, C. The origin and evolution of *Homo sapiens*. *Philos Trans R Soc*
492 *Lond B Biol Sci* **371**, 20150237 (2016).
- 493 28 Coyne, J. A. & Allen Orr, H. *Speciation*. (Sinauer Associates, Inc., 2004).
- 494 29 Hudson, R. R. Two-locus sampling distributions and their application.
495 *Genetics* **159**, 1805-1817 (2001).
- 496 30 Scally, A. & Durbin, R. Revising the human mutation rate: implications for

- 497 understanding human evolution. *Nat. Rev. Genet.* **13**, 745-753 (2012).
- 498 31 Campbell, C. D. *et al.* Estimating the human mutation rate using autozygosity
499 in a founder population. *Nat. Genet.* **44**, 1277-1281 (2012).
- 500 32 Conrad, D. F. *et al.* Variation in genome-wide mutation rates within and
501 between human families. *Nat. Genet.* **43**, 712-714 (2011).
- 502 33 Kong, A. *et al.* Rate of de novo mutations and the importance of father's age to
503 disease risk. *Nature* **488**, 471-475 (2012).
- 504 34 Bhaskar, A. & Song, Y. S. Descartes' rule of signs and the identifiability of
505 population demographic models from genomic variation data. *Ann. Stat.* **42**,
506 2469-2493 (2014).
- 507 35 Frankish, A. *et al.* GENCODE reference annotation for the human and mouse
508 genomes. *Nucleic Acids Res.* **47**, D766-D773 (2019).
- 509 36 Fay, J. C. & Wu, C.-I. Hitchhiking under positive Darwinian selection.
510 *Genetics* **155**, 1405-1413 (2000).
- 511 37 Hsieh, P. *et al.* Model-based analyses of whole-genome data reveal a complex
512 evolutionary history involving archaic introgression in Central African
513 Pygmies. *Genome Res.* **26**, 291-300 (2016).
- 514 38 Schlebusch, C. M. & Jakobsson, M. Tales of human migration, admixture, and
515 selection in Africa. *Annu. Rev. Genomics Hum. Genet.* **19**, 405-428 (2018).
- 516 39 Skoglund, P. *et al.* Reconstructing prehistoric African population structure.
517 *Cell* **171**, 59-71 (2017).
- 518 40 Ramachandran, S. *et al.* Support from the relationship of genetic and
519 geographic distance in human populations for a serial founder effect
520 originating in Africa. *Proc. Natl. Acad. Sci. USA* **102**, 15942-15947 (2005).
- 521 41 Prugnolle, F., Manica, A. & Balloux, F. Geography predicts neutral genetic
522 diversity of human populations. *Curr. Biol.* **15**, R159-160 (2005).
- 523 42 Ayub, Q. *et al.* The Kalash genetic isolate: ancient divergence, drift, and
524 selection. *Am. J. Hum. Genet.* **96**, 775-783 (2015).
- 525 43 Li, H. *et al.* Large numbers of vertebrates began rapid population decline in
526 the late 19th century. *Proc. Natl. Acad. Sci. USA* **113**, 14079-14084 (2016).
- 527 44 Goren-Inbar, N. *et al.* Evidence of hominin control of fire at Gesher Benot
528 Ya'aqov, Israel. *Science* **304**, 725-727 (2004).
- 529 45 Alperson-Afil, N. *et al.* Spatial organization of hominin activities at Gesher
530 Benot Ya'aqov, Israel. *Science* **326**, 1677-1680 (2009).
- 531 46 Foley, R. & Gamble, C. The ecology of social transitions in human evolution.
532 *Philos. Trans. R. Soc. Lond., Ser. B: Biol. Sci.* **364**, 3267-3279 (2009).
- 533 47 Du, A. *et al.* Pattern and process in hominin brain size evolution are
534 scale-dependent. *Proc. Biol. Sci.* **285**, 20172738 (2018).
- 535 48 Green, R. E. *et al.* A draft sequence of the Neandertal genome. *Science* **328**,
536 710-722 (2010).
- 537 49 Reich, D. *et al.* Genetic history of an archaic hominin group from Denisova
538 Cave in Siberia. *Nature* **468**, 1053-1060 (2010).
- 539 50 Myers, S., Fefferman, C. & Patterson, N. Can one learn history from the allelic
540 spectrum? *Theor. Popul. Biol.* **73**, 342-348 (2008).

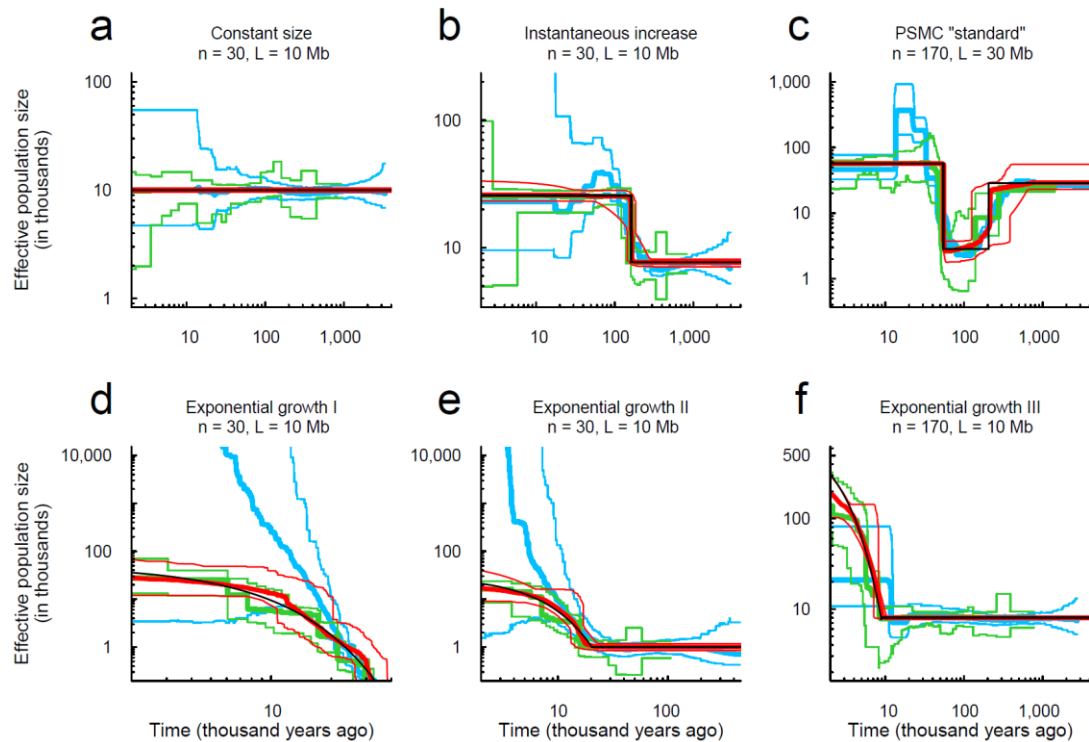
- 541 51 Chen, H. A computational approach for modeling the allele frequency
542 spectrum of populations with arbitrarily varying size. *Genomics Proteomics*
543 *Bioinf.* **17**, 635-644 (2019).
- 544 52 Pedersen, M. E. H. Tuning and simplifying heuristical optimization. *PhD*
545 *thesis, Univ. Southampton* (2010).
- 546 53 Hudson, R. R. Generating samples under a Wright-Fisher neutral model of
547 genetic variation. *Bioinformatics* **18**, 337-338 (2002).
- 548 54 Chen, G. K., Marjoram, P. & Wall, J. D. Fast and flexible simulation of DNA
549 sequence data. *Genome Res.* **19**, 136-142 (2009).
- 550 55 Yu, D. *et al.* eGPS 1.0: comprehensive software for multi-omic and
551 evolutionary analyses. *Natl. Sci. Rev.* **6**, 867-869 (2019).
- 552 56 Harpending, H. C. *et al.* Genetic traces of ancient demography. *Proc. Natl.*
553 *Acad. Sci. USA* **95**, 1961-1967 (1998).
- 554
- 555
- 556



557

558 **Figure 1. Illustration of the fast infinitesimal time coalescent (FitCoal) process.**

559 The left panel shows the backward process in which four lineages coalesce into one
560 after passing through millions of infinitesimal time intervals. The highlighted area
561 shows the backward transformation process of different states with tiny probability
562 changes in an infinitesimal time interval (Δt). Thick arrows indicate high
563 transformation probabilities, and thin arrows indicate low transformation probabilities.
564 Each state is indicated with a rounded rectangle, in which one circle indicates one
565 lineage. The rounded rectangles with black filled circles are the states with probability
566 1. The rounded rectangles with empty circles are the states with probability 0. The
567 probabilities between 0 and 1 are indicated by grey circles. The middle panel shows
568 branches of different states. The right panel shows the demographic history of a
569 population. The width of shadowed area indicate the effective population size, *i.e.*, the
570 number of breeding individuals⁵⁶. It is assumed that the effective population size
571 remains unchanged within Δt .



572

573 **Figure 2. Demographic histories estimated by FitCoal, stairway plot, and PSMC**

574 **using simulated samples. (a) Constant size model. (b) Instantaneous increase model.**

575 **(c) PSMC “standard” model. (d) Exponential growth I model. (e) Exponential growth**

576 **II model. (f) Exponential growth III model. These six models are the same as those of**

577 **the previous study by Liu and Fu¹⁵. Thin black lines indicate true models. Thick red**

578 **lines indicate the medians of FitCoal estimated histories; thin red lines are 2.5 and**

579 **97.5 percentiles of FitCoal estimated histories. Green and blue lines indicate the**

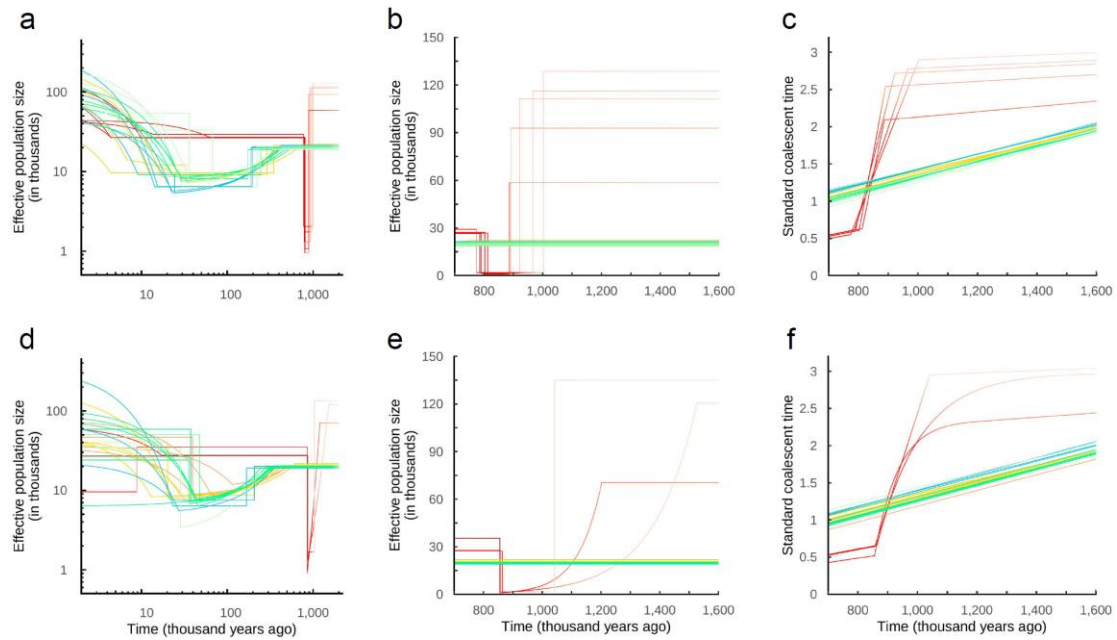
580 **results of stairway plot and PSMC, respectively, of the previous study¹⁵. The mutation**

581 **rate is assumed to be 1.2×10^{-8} per base per generation, and a generation time is**

582 **assumed to be 24 years. n is the number of simulated sequences, and L is the length of**

583 **simulated sequences.**

584



585

586

587 **Figure 3. FitCoal estimated histories of human populations using 1000GP and**

588 **HGPD-CEPH genomic data sets. (a) Estimated histories of 26 populations in**

589 **1000GP. (b) Linear-scaled estimation of histories of 1000GP populations during the**

590 **super bottleneck period. (c) Calendar time vs standard coalescent time of estimated**

591 **histories of 1000GP populations. (d) Estimated histories of 24 HGPD-CEPH**

592 **populations. (e) Linear-scaled estimation of histories of HGPD-CEPH populations**

593 **during the super bottleneck period. (f) Calendar time vs standard coalescent time of**

594 **estimated histories of HGPD-CEPH populations. Various color lines indicate the**

595 **following: red, African populations; yellow, European populations; brown, Middle**

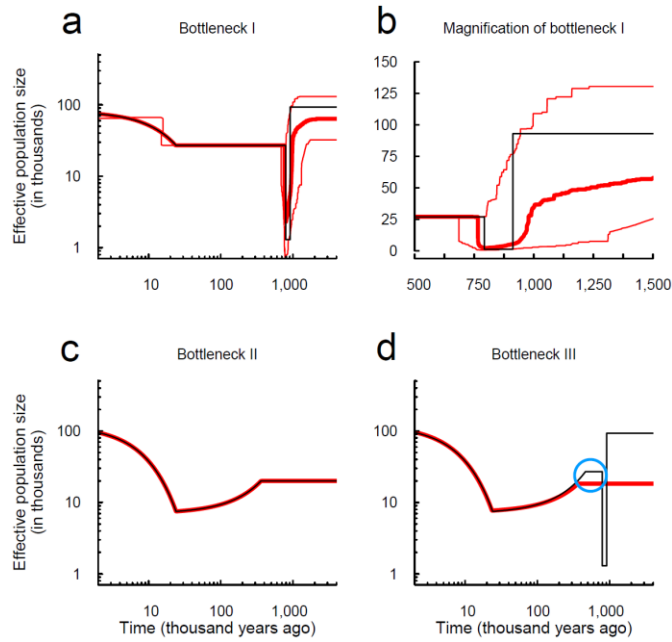
596 **East populations; blue, East Asian populations; green, Central or South Asian**

597 **populations; and dark sea green, American populations. The mutation rate is assumed**

598 **to be 1.2×10^{-8} per base per generation, and a generation time is assumed to be 24**

599 **years.**

600



601

602 **Figure 4. Verification of the super bottleneck.** (a) Bottleneck I model, mimicking
603 the demography of 1000GP African population and its estimated histories. (b)
604 Linear-scaled Bottleneck I model during the super bottleneck period. (c) Bottleneck II
605 model, mimicking the estimated demography of 1000GP non-African population and
606 its estimated histories. (d) Bottleneck III model, mimicking the true demography of
607 1000GP non-African population and its estimated histories. Thin black lines indicate
608 models. Thick red lines denote the medians of FitCoal estimated histories; thin red
609 lines represent 2.5 and 97.5 percentiles of FitCoal estimated histories. Blue circle
610 indicates the population size gap. The mutation rate is assumed to be 1.2×10^{-8} per
611 base per generation, and a generation time is assumed to be 24 years. The number of
612 simulated sequences is 202 in Bottleneck I and 200 in Bottleneck II and III. The
613 length of simulated sequence is 800 Mb.

614

615

616