



**HAL**  
open science

# Deep divergences of human gene trees and models of human origins

Michael Blum, Mattias Jakobsson

► **To cite this version:**

Michael Blum, Mattias Jakobsson. Deep divergences of human gene trees and models of human origins. Molecular Biology and Evolution, 2010, 28 (2), pp.889. 10.1093/molbev/MSQ265 . hal-00629968

**HAL Id: hal-00629968**

**<https://hal.science/hal-00629968>**

Submitted on 7 Oct 2011

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

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

# Deep divergences of human gene trees and models of human origins

Michael GB Blum\* and Mattias Jakobsson†

\*Laboratoire TIMC-IMAG, CNRS, Université Joseph Fourier, Grenoble, France

†Department of Evolutionary Biology, Uppsala University, Sweden

Corresponding author

Michael G B Blum

Laboratoire TIMC-IMAG, Faculty of Medicine, 38706 La Tronche, France

Phone +33 4 56 52 00 65

Fax +33 4 56 52 00 55

email: michael.blum@imag.fr

Keywords: human origins, time to the most recent common ancestor, TMRCA, archaic admixture, African bottleneck, coalescent

Running head: Deep divergences of human gene trees

Abbreviation: TMRCA, Time to the Most Recent Common Ancestor

## Abstract

Two competing hypotheses are at the forefront of the debate on modern human origins. In the first scenario, known as the recent Out-of-Africa hypothesis, modern humans arose in Africa about 100,000-200,000 years ago, and spread throughout the world by replacing the local archaic human populations. By contrast, the second hypothesis posits substantial gene flow between archaic and emerging modern humans. In the last two decades, the young time estimates – between 100,000 and 200,000 years – of the most recent common ancestors for the mitochondrion and the Y-chromosome provided evidence in favor of a recent African origin of modern humans. However, the presence of very old lineages for autosomal and X-linked genes has often been claimed to be incompatible with a simple, single origin of modern humans. Through the analysis of a public DNA sequence database, we find, similar to previous estimates, that the common ancestors of autosomal and X-linked genes are indeed very old, living, on average, respectively 1,500,000 and 1,000,000 years ago. However, contrary to previous conclusions, we find that these deep gene genealogies are consistent with the Out-of-Africa scenario provided that the ancestral effective population size was approximately 14,000 individuals. We show that an ancient bottleneck in the Middle Pleistocene, possibly arising from an ancestral structured population, can reconcile the contradictory findings from the mitochondrion on the one hand, with the autosomes and the X-chromosome on the other hand.

## Introduction

The process by which modern humans arose has been the subject of much debate in paleoanthropology (Stringer 2002). Especially the extent of admixture between anatomically modern humans and archaic populations of *Homo* has been vigorously debated (Wolpoff et al. 2000; Templeton 2002; Plagnol and Wall 2006; Garrigan and Hammer 2006; Fagundes et al. 2007). At one end of the spectrum, the recent Out-of-Africa hypothesis posits that a group of modern humans, arising in Africa about 100,000-200,000

years ago, spread throughout the world by replacing, without admixture, the local archaic human populations (Mellars 2006). At the other end, the Multi-Regional model posits substantial gene flow between local archaic humans and the emerging modern humans, even though it does not exclude Africa as the cradle of modern humans (Wolpoff et al. 2000). In between these two hypotheses, alternative scenarios assume archaic admixture restricted to Africa before the emerging modern humans eventually colonized the globe (Harding and McVean 2004; Gunz et al. 2009).

Substantial genetic evidence has been put forward in support of the recent Out-of-Africa hypothesis. For example, a continuous decrease of genetic diversity with increasing distance from Africa has been observed for autosomal microsatellites (Prugnolle et al. 2005; Ramachandran et al. 2005) and SNPs (Li et al. 2008), as well as a continuous increase of linkage disequilibrium with distance from Africa (Jakobsson et al. 2008). These world-wide patterns are consistent with a serial founder model in which migrant populations are formed from a subset of the previous population in the migration wave outward from Africa (DeGiorgio et al. 2009). Sex-linked markers have also provided evidence for a recent African origin since the common ancestor of all contemporary mitochondrial haplotypes existed as recently as  $\sim 200,000$  years ago (Cann et al. 1987), and the ancestor of all Y-chromosomes lived  $\sim 100,000$  years ago (Thomson et al. 2000; Wilder et al. 2004).

Direct evidence against the recent Out-of-Africa hypothesis can potentially come from comparisons between ancient DNA of archaic humans, such as Neanderthals, and DNA of present-day modern humans (Noonan 2010). Recently, when releasing the draft sequence of a Neanderthal genome, Green et al. (2010) found 1 – 4% of Neanderthal introgression in the gene pool of non-Africans; however, previous comparisons involving ancient Neanderthal DNA did not provide evidence in favor of admixture between humans and Neanderthals (Krings et al. 1997; Noonan et al. 2006; Wall and Kim 2007). Genetic evidence that does not involve ancient DNA, in particular elevated levels of Linkage Disequilibrium (LD) in modern humans, was also found to be indicative of ancient admixture (Plagnol and Wall 2006; Wall et al. 2009). It has additionally

been argued that the presence of very old lineages, or deep divergences, for autosomal genes and genes on the X-chromosome is incompatible with a simple, single origin of modern humans, and that these deep divergences are instead evidence in favor of archaic admixture (Harris and Hey 1999; Harding and McVean 2004; Garrigan et al. 2005; Evans et al. 2006; Garrigan and Hammer 2006; Hayakawa et al. 2006; Patin et al. 2006; Cox et al. 2008; Kim and Satta 2008).

A measure of the (deepest) divergence of a gene tree is the Time to the Most Recent Common Ancestor (TMRCA). The TMRCA is the time at which the Most Recent Common Ancestor (MRCA) of all existing copies of a given gene lived. A genome-wide frequency distribution of the TMRCAs has been reported by curating the literature (Garrigan and Hammer 2006) but no systematic and consistent analysis has been performed in a single genome-wide data set. We report the first genome-wide estimation of the TMRCAs of anatomically modern humans, and we investigate if different scenarios of human evolutionary history are supported by this estimate. In particular, we investigate to what extent the ages of the autosomal and X-linked lineages are compatible with the recent Out-of-Africa hypothesis.

## Materials and Methods

### Sequence data.

The data comprises of 40 re-sequenced independent intergenic regions from the autosomes and the X-chromosome (Wall et al. 2008). Each region encompasses approximately 20kb and consists of three 2kb sequence fragments, separated by 7kb of unsequenced DNA (“locus trio design”, see Garrigan et al. 2005). We present the results for 78 individuals from three African populations (Mandenka, Biaka, and San from Namibia), one European population (Basque), one Asiatic population (Han) and one population from Oceania (Melanesia). Two common chimpanzee sequences, available in the database, were used as outgroups.

## TMRCA estimation.

We use the method of Thomson et al. (2000) for estimating TMRCAs and we reiterate the basic motivation for the estimator (for an extensive description, see Thomson et al. 2000; Hudson 2007). Thomson’s estimator requires an outgroup that provides information on the ancestral state at every polymorphic position and assumes an infinite sites model.

Let  $T$  be the time to the MRCA for a sample of  $n$  lineages, and let  $x_i$  be the number of mutations that have occurred between the MRCA and lineage  $i$ . We assume that the number of mutations along a branch,  $x_i$ , is Poisson distributed with mean equal to the product of the mutation rate ( $u$ ) and the branch length ( $T$ ). Then, Thomson’s estimator of TMRCA is

$$\hat{T} = \sum_{i=1}^n x_i / (nu).$$

To estimate the mutation rate  $u$ , we assume a molecular divergence of 6 millions years between human and chimpanzee (Glazko and Nei 2003), and a generation time of 25 years. Computing the mean number of nucleotide differences between two chimp sequences and the human sequences, we find a mean mutation rate of  $9.90 \times 10^{-10}$ /bp/year, on the same order as Fagundes et al. (2007).

To compute Thomson’s estimator  $\hat{T}$ , we reduce the data set to human polymorphic sites. Since we only consider sites that are polymorphic within humans, the mutations unique to the outgroup sequences are excluded from the computation of the TMRCA estimator (and have no effect on the TMRCA estimates). We reconstruct the ancestral state of each polymorphic site and, under the infinitely-many-sites mutation model, the state of the most recent common ancestor (MRCA) of the within-species sample is assumed to be equal to the outgroup allelic state. For 20 out of 1588 sites, the ancestral state could not be deduced, and, for these sites, the most frequent allele was assumed to be the ancestral variant. The choice of the infinitely-many-sites model is motivated by the small estimated mutation rate ( $u = 9.90 \times 10^{-10}$ /bp/year), which makes the probability of two mutations hitting the same site negligible.

We investigate the effect of recombination on Thomson’s estimator using simulations, and we consider an effective recombination rate ranging from 0 to 10. This range encompasses previous estimates of the recombination rate in humans (The International HapMap Consortium 2005; Voight et al. 2005; Coop et al. 2008b). Using, for instance, a standard estimate of 1cM/Mb (i.e.  $10^{-8}$ /bp) for the recombination rate (The International HapMap Consortium 2005) gives an effective recombination rate of  $\rho = 8$  for a 20kb region, assuming an effective population size of 10,000 individuals. To quantify the effect of recombination on  $\hat{T}$ , we compute the relative bias and mean square-error of  $\hat{T}$  as a function of the recombination rate (Figure 1).

We also investigate if Thomson’s estimator  $\hat{T}$  remains accurate when the population experienced demographic changes, and we consider three different demographic models: a) a population with a constant size of  $N = 10,000$  individuals; b) an expanding population where the population experienced a 10-fold expansion starting at a time distributed uniformly between 0 and 200,000 years ago from a population-size of  $N = 10,000$  individuals; and c) a population splitting model. For the last model c, we assume that an ancestral population (of size  $N = 10,000$  individuals) split at a time distributed uniformly between 0 and 200,000 years ago into two subpopulations. One of the two subpopulations contains  $N = 10,000$  individuals and the other population contains  $N = 10,000 \times p$  individuals with  $p$  chosen from a uniform distribution between 0 and .5.

## Approximate Bayesian Computation (ABC).

We use ABC to find the range of demographic parameters that yield TMRCAs similar to the empirical estimates of TMRCAs. For each scenario of human evolution (described in the Results and Discussion Section), the ABC statistical procedure can be described as follows:

- Generate the demographic parameters according to the prior distributions given in Supplementary Table 1.

- Simulate sequence data with the software *ms* (Hudson 2002). One simulation comprises of generating 20 autosomal and 20 X-linked sequence-regions with the same number of samples and the same sequence lengths as in the empirical data.
- Compute the summary statistics (see below) for the simulated sequences and compute the Euclidean distance between observed and simulated summary statistics.

After performing a total of 100,000 simulations, we retain the 500 simulations that provide the best match to the data. We use an Epanechnikov kernel to assign larger weights to simulations that provide the best match (Beaumont et al. 2002). To account for the imperfect match between simulated and observed summary statistics, we then use regression adjustment as described by Blum and François (2010). After completion of the algorithm, the posterior distribution of the parameters consists of the set of accepted parameters after adjustment.

To compute the summary statistics, we consider, for each of the 40 sequence-regions, the mean number of mutations,  $\sum_{i=1}^n x_i/n$ , between the set of sequences and the ancestral sequence. To reduce the number of summary statistics, we compute, separately for the X-linked and the autosomal markers, the three quartiles of the 20 mean numbers of mutations. This procedure results in a total of 6 summary statistics.

## **Choice of priors for the mutation and recombination rates.**

For the mutation rate, we choose an empirical Bayes approach in which the prior depends on the data (Casella 1985). We choose a Gamma distribution for the mutation rate and we estimate the parameters so that the Gamma distribution fits the empirical distribution of the 40 estimated mutation rates. We obtain a Gamma distribution with a shape parameter of 15.18 and a scale parameter of  $6.50 \times 10^{-11}$ . This results in a median mutation rate of  $9.65 \times 10^{-10}$  mutations/bp/year and 95% of the simulated mutation rates are between  $5.54 \times 10^{-10}$  and  $1.54 \times 10^{-9}$  mutations/bp/year. For the crossing-over rate, we consider a log-normal distribution with a mean and standard deviation (on a log scale) of  $-18.148$  and  $0.5802$  (Voight et al. 2005). We assume



homogeneous cross-over rate along each sequence-region.

## Computation of the relative model probabilities.

We introduce an indicator variable  $Y = \{1, 2, 3, 4\}$  that indicates, for each simulation, which one of the four different demographic scenarios (described in the Results and Discussion Section) generated the data. We perform the same number of simulations for each scenario. We then regress the indicator variable  $Y$  by the 6 summary statistics using local multinomial logistic regression to obtain the model probabilities  $P(Y|\mathbf{s})$  as a function of the summary statistics  $\mathbf{s}$  (Fagundes et al. 2007; Beaumont 2008). Local logistic regression differs from standard logistic regression because larger weights are given to the simulations for which the summary statistics are close to the observed summary statistics. By computing the logistic regression equation for the observed summary statistics  $P(Y|\mathbf{s} = \mathbf{s}_{\text{obs}})$ , we obtain the relative model probabilities. To perform local multinomial logistic regression, we use the R package VGAM (Yee 2010).

## Posterior predictive simulations.

To perform goodness-of-fit of the scenarios of human evolution (Figure 2), we perform posterior predictive checks (Gelman et al. 2003). We sample, with replacement, 10,000 multivariate demographic parameters at random from the posterior distribution obtained with the ABC algorithm. Using *ms* (Hudson 2002), we simulate, for each multivariate demographic parameter, gene trees along a 20kb sequence, and compute the median of the (potentially different) simulated TMRCAs found along the sequence. For each scenario, this results in a total of 10,000 median TMRCAs that are displayed in Figure 2.

# Results and Discussion

## Genome-wide estimation of the TMRCA.

The data comprise 40 re-sequenced independent intergenic regions from the autosomes and the X-chromosome provided by a public DNA sequence database that has been designed for the purpose of analyzing human prehistory (Wall et al. 2008). To compute genome-wide estimation of the TMRCA, we consider the statistic of Thomson et al. (2000), which has been used for instance to date the ancestor of the human Y-chromosome (Thomson et al. 2000) and the human FOXP2 gene (Coop et al. 2008a). For each DNA fragment, the TMRCA estimate is obtained by computing the number of mutations between each sample (gene-copy) and a reconstructed ancestral sequence, and averaging across the gene-copies of the sample. Generally, the bias of this estimator has been shown to be small (Hudson 2007), and we demonstrate that it is also robust to recombination (see Figure 1). To affix a time scale in years to the TMRCA estimates, we assume a molecular divergence of 6 millions years between human and chimpanzee (Glazko and Nei 2003), and a generation time of 25 years. For each region of 20kb, we compute one TMRCA estimate so that this estimate captures an average TMRCA for the possibly different TMRCA found along the 20kb region.

Figure 2 displays the distribution of the TMRCA for the 20 autosomal loci (Figure 2A) and the 20 loci located on the X-chromosome (Figure 2B). The median of the TMRCA is approximately 1,500,000 years for the autosomes (first quartile = 950,000, third quartile = 1,700,000) and approximately 1,000,000 years for the X-chromosome (first quartile = 700,000, third quartile = 1,350,000). These numbers are close to what Garrigan and Hammer (2006) found by collecting TMRCA estimates from the literature. However, they are at odds with numerical predictions obtained by Fagundes et al. (2007) for the recent Out-of-Africa model. They found that the distribution of autosomal TMRCA should peak around the time when modern humans emerged (100,000 to 200,000 years ago) and that 50% of the TMRCA should be more recent than 650,000 years. In contrast, we find that only two out of 20 autosomal TMRCA

are more recent than 650,000 years.

A number of authors have argued that deep genealogical histories are incompatible with the recent Out-of-Africa hypothesis, and instead claimed that these deep genealogies are evidence in favor of archaic admixture (Harris and Hey 1999; Harding and McVean 2004; Garrigan et al. 2005; Evans et al. 2006; Garrigan and Hammer 2006; Hayakawa et al. 2006; Patin et al. 2006; Cox et al. 2008; Kim and Satta 2008). In the following sections, we investigate to what extent the genome-wide distribution of the TMRCA is compatible with the recent Out-of-Africa model. We also consider alternative models that assume archaic admixture and check if they provide a better fit to the TMRCA distribution.

## **TMRCAs distributions predicted by different scenarios of human evolution.**

We compare observed TMRCA to simulated TMRCA for four different scenarios of modern human origin (Figure 3 and Supplementary Table 1). The two first scenarios are versions of the Out-of-Africa model, and the last two scenarios are versions of the Multi-Regional model:

1. ***'Single Origin Population'*** A recent Out-of-Africa scenario in which modern humans descended from one subpopulation of archaic humans that was a separate population for a long time in Africa,
2. ***'Multiple Archaic Populations'*** A recent Out-of-Africa scenario in which different archaic African populations were connected by gene flow, even though only one archaic population eventually colonized the globe (Harding and McVean 2004; Garrigan and Hammer 2006; Campbell and Tishkoff 2008),
3. ***'Recent Admixture'*** A Multi-Regional scenario in which archaic and modern humans were isolated during 300-600,000 years and admixed recently in Eurasia, 30-70,000 years ago (Plagnol and Wall 2006), and

4. *‘Long-Standing Admixture’* A Multi-Regional scenario with continuous and long-standing admixture between archaic and modern humans.

To find the range of demographic parameter-values that yield TMRCAs similar to the 40 estimates from the empirical data, we use Approximate Bayesian Computation (ABC) (Beaumont et al. 2002; Csilléry et al. 2010) and coalescent simulations (Hudson 2002). The empirical genetic data is summarized by six summary statistics that capture the divergence of gene trees. For each DNA fragment, we compute, for the human polymorphic sites, the mean number of mutations between the gene-copies of the sample and the reconstructed ancestral sequence (see Materials and Methods). We compute – separately for the X-chromosome and the autosomes – the three quartiles of these averaged number of mutations.

We estimate the ancestral effective population size of archaic humans (Figure 4 and Supplementary Table 2) and find, for the recent Out-of-Africa scenario, that the most likely value is 14,500 (95% CI = 12,000-17,000) similar to previous estimates (Takahata 1993; Harding et al. 1997; Wall 2003; Voight et al. 2005). In scenario 2, which assumes a structured population in Africa before the emergence of modern humans, we find a slightly lower estimate on the order of 8,000 individuals (95% CI = 5,000-15,000) reflecting that several archaic African populations contribute to the modern gene pool.

In addition to parameter inference, the ABC approach also offers a convenient way to assign a probability to each of the scenarios (Fagundes et al. 2007; François et al. 2008; Verdu et al. 2009). We find that the four different models are almost equally supported by the divergence of gene trees since the four posterior probabilities range between 20% and 30% (Supplementary Figure 1). These even probabilities reflect that the relatively ancient lineages found in the autosomes and X-linked genes neither favor nor disfavor the models with archaic admixture (models 2-4).

To check if the different scenarios of human evolution provide a good fit to the data, we compare the empirical TMRCAs estimates to the TMRCAs predicted by the different scenarios. All four models predict, on average, lineages as old as 1,500,000 years for autosomal fragments (Figure 2A) and as old as 1,000,000 years for X-linked fragments

(Figure 2B). In short, we find that both the simple replacement model and the models with archaic admixture are perfectly compatible with the deep divergences found in the empirical data.

However, not all aspects of the empirical TMRCAs are well captured by the modeled scenarios of human evolution. The variance of the empirical TMRCAs is larger than the variance predicted by three of the four different models of human evolution (see Figure 2 and Supplementary Table 3), and this large variance has been interpreted as the result of archaic sub-structure in Africa (Harding and McVean 2004). Indeed, the ‘Multiple Archaic Populations’ (scenario 2) shows similar variance of TMRCAs as the empirical data, but the inflated variance of the empirical TMRCAs estimates can also be due to variation in mutation or recombination rate across the 40 sequence-regions (McVean et al. 2004).

## **The mitochondrion and the Y-chromosome.**

We also investigate the distribution of TMRCAs that is expected for the Y-chromosome and the mitochondrial chromosome (Figure 2C). The models of human evolution typically predict older TMRCAs compared to the estimated 170,000 years for mtDNA (Ingman et al. 2000) and the upper estimate of 100,000 years for the Y-chromosome (Tang et al. 2002; Wilder et al. 2004; Shi et al. 2010). For mtDNA, a TMRCAs of 170,000 years is within the range of values predicted by the ‘Multiple Archaic Populations’ scenario ( $P(\text{TMRCAs} < 170,000) = 0.21$ ), but the mitochondrial TMRCAs estimate is difficult to reconcile with the remaining three scenarios ( $P < 4 \times 10^{-2}$ ). For the Y-chromosome, a TMRCAs of 100,000 years is clearly at odds with three of the models ( $P < 6 \times 10^{-4}$ ), but for the ‘Multiple Archaic Populations’ scenario with archaic African admixture, the proportion of simulated gene trees with TMRCAs younger than 100,000 years is larger than for the other three models, albeit quite small ( $P = 1.5 \times 10^{-2}$ ). Although assuming an archaic structured population in Africa, as in scenario 2, reduces the Y-chromosome and mtDNA TMRCAs, the model cannot fully explain a Y-chromosome ancestor living as recently as 100,000 years ago. In scenario 2, we consider three archaic populations,

and increasing the number of archaic populations will further decrease the effective population size of each sub-population, which will decrease the predicted haploid TMRCA, bringing them in line with the empirical estimates. However, since the mtDNA and the Y-chromosome are non-recombining units, their young TMRCA can also be explained by recent selective sweeps, caused by directional selection at any gene within the non-recombining units (Kreitman 2000; Jobling and Tyler-Smith 2003).

## A bottleneck when anatomically modern humans emerged

An alternative explanation for young haploid TMRCA involves a demographic bottleneck concomitant with the emergence of anatomically modern humans. The earliest known suite of derived morphological traits associated with modern humans appears in fossil remains from Ethiopia dated to 150-190 kya (White et al. 2003; McDougall et al. 2005). A model of the origin and spread of modern humans proposed by Lahr and Foley (1994) assumes that the emerging modern humans experienced bottlenecks within Africa during the penultimate glacial age (130-190 kya) when cold, dry climates caused isolation of populations within Africa (see also Ambrose 1998). Here we consider a bottleneck that occurred 150,000 years ago, and we measure the bottleneck intensity by the inbreeding coefficient during the bottleneck

$$F = 1 - \left(1 - \frac{1}{2bN_A}\right)^D,$$

where  $bN_A$  is the diploid population size during the bottleneck and  $D$  corresponds to the duration (in generations) of the bottleneck. To give a reference value, the inbreeding coefficient  $F$  corresponding to the Out-of-Africa bottleneck was inferred at  $F = 0.175$  (Akey et al. 2004; Voight et al. 2005). Considering an uniform prior between 0 and 1 for  $F$ , we find that there is a large range of parameter values compatible with the autosomal and X-linked TMRCA (Figure 5). For instance, both pairs of parameter values ( $F = 0, N_A = 14,000$ ) and ( $F = 0.4, N_A = 18,000$ ) are clearly within the range of the bivariate posterior distribution found with our ABC approach. The relatively

large bivariate posterior range of the inbreeding coefficient  $F$  and of the ancestral size  $N_A$  shows that a strong bottleneck can still produce an average autosomal TMRCA of 1.5 million year provided that the ancestral size was large enough before the bottleneck.

To investigate if a bottleneck 150,000 years ago in Africa can account for both recent haploid and old autosomal ancestors, we simulate TMRCA for the different chromosomes by sampling the demographic parameters from the posterior distribution that was obtained using ABC. We find that the recent mtDNA TMRCA is clearly within the range of predicted values ( $P=0.37$ , Figure 2C) as well as the old X-linked and autosomal TMRCA (Figure 2A-B). The fact that a bottleneck can accommodate an 8-fold discrepancy between autosomal and mtDNA TMRCA can be seen by plotting TMRCA as a function of  $F$  and  $N_A$  (Figure 6). Setting the inbreeding coefficient at  $F = 0.4$  for instance, Figure 6 shows that the mean mtDNA TMRCA is smaller than 200,000 years, for a large range of values of  $N_A$ , whereas the 20kb-averages of TMRCA, for autosomal and X-linked markers, are older than 1.5 and 1 million years when  $N_A > 16,000$  individuals. Finally, as for the other models of human evolution, the TMRCA of 100,000 years for the Y-chromosome remains unexpectedly young ( $P = 4 \times 10^{-4}$ ).

We find that both the ‘multiple archaic population’ model and a sudden bottleneck, 150,000 years ago in Africa, can account for the 8-fold discrepancy between autosomal and mtDNA TMRCA. Although modeling different patterns of human evolution, these two scenarios are different versions of a bottleneck in the human lineage that arose before the succeeding migration out of Africa. Previous attempts to detect an ancestral African bottleneck have often been inconclusive and the genetic evidence in favor of the Out-of-Africa bottleneck are more salient (Marth et al. 2004; Voight et al. 2005). However, it is more difficult to detect this potential middle-Pleistocene bottleneck compared to the Out-of-Africa bottleneck since it is more ancient (Depaulis et al. 2003). Additionally, the unexpected large levels of African LD that has been interpreted as evidence of archaic admixture (Plagnol and Wall 2006; Wall et al. 2009) may potentially be explained by an ancient bottleneck (Schmegner et al. 2005). We therefore suggest

that an ancestral bottleneck, possibly following ancestral sub-structure and admixture, is important to consider when investigating demographic models of human evolution (Schaffner et al. 2005; Gutenkunst et al. 2009; Laval et al. 2010).

## A simple population genetic prediction

Although we perform coalescent simulations, standard population genetic theory can predict the TMRCA found for the autosomes and the X-chromosome in the Out-of-Africa model when there is no ancestral bottleneck or admixture. For a population of constant diploid size  $N$  in which the effective number of males and females is the same, the expected waiting time (in generations) before the coalescence of a (reasonably large) sample of genes is approximately  $4N$  for the autosomes and  $3N$  for the X-linked genes (see e.g. Hein et al. 2005). Using a generation time of 25 years and an effective population size of 14,000 individuals, the computation leads to an average TMRCA of 1,400,000 years for autosomal genes and 1,050,000 years for X-linked genes, which are both very close to our estimates from the sequence data. This theoretical argument shows that the difference between the TMRCA of the autosomes and the X-linked genes is easily explained by the difference in effective population size. The same argument yields an average TMRCA of 350,000 years for both the mtDNA and Y-chromosome, clearly deviating from the estimates based on empirical data (Ingman et al. 2000; Wilder et al. 2004; Shi et al. 2010).

## Conclusion

We provide 20 autosomal and 20 X-linked estimates of TMRCA of a sample of contemporary humans and find that the autosomal ancestors of modern humans lived  $\sim 1,500,000$  years ago and that the X-linked ancestors lived  $\sim 1,000,000$  years ago. The ranges of values for the TMRCA are quite large: 450,000-2,400,000 years for the autosomes, and 380,000-2,000,000 for the X-chromosome. These values are in the same range as previous estimates for autosomal and X-linked genes (see e.g. Templeton 2002;



Garrigan and Hammer 2006; Tishkoff and Gonder 2006). We investigate to what extent the recent Out-of-Africa model reproduces the pattern of estimated TMRCAs, and when setting the ancestral effective size of humans to  $\sim 14,000$ , this model reproduces the old TMRCAs of the empirical data. Deep divergences in human gene trees are therefore not incompatible with the recent Out-of-Africa hypothesis and the observation of deep gene genealogies should not be taken as evidence for the Multi-Regional hypothesis. Finally, we show that an ancestral bottleneck in Africa, possibly arising in a structured population, can account for the unexpectedly large discrepancy between young mtDNA and Y-chromosome ancestors and old autosomal and X-linked ancestors.

## Acknowledgments

This work was supported by a Grant from the Swedish Foundation for International Cooperation in Research and Higher Education (STINT) to MJ and MGBB. MJ would like to thank the Swedish Research Council and the Erik Philip-Sørensen's foundation for financial support.

## References

- Akey, J., M. Eberle, M. Rieder, C. Carlson, M. Shriver, D. Nickerson and L. Kruglyak. 2004. Population history and natural selection shape patterns of genetic variation in 132 genes. *PLoS Biology* **2**:e286.
- Ambrose, S. 1998. Late Pleistocene human population bottlenecks, volcanic winter, and differentiation of modern humans. *Journal of Human Evolution* **34**:623–651.
- Beaumont, M. 2008. Simulation, genetics and human prehistory. In S. Matsumura, P. Forster and C. Renfrew, eds., *Joint determination of topology, divergence time, and immigration in population trees*, 134–154. Cambridge: McDonald Institute for Archaeological Research.
- Beaumont, M. A., W. Zhang and D. J. Balding. 2002. Approximate Bayesian computation in population genetics. *Genetics* **162**:2025–2035.
- Blum, M. G. B. and O. François. 2010. Non-linear regression models for approximate Bayesian computation. *Statistics and Computing* **20**:63–73.
- Campbell, M. C. and S. A. Tishkoff. 2008. African genetic diversity: Implications for human demographic history, modern human origins, and complex disease mapping. *Annual Review of Genomics and Human Genetics* **9**:403–433.

- Cann, R. L., M. Stoneking and A. C. Wilson. 1987. Mitochondrial DNA and human evolution. *Nature* **325**:31–36.
- Casella, G. 1985. An introduction to empirical Bayes data analysis. *The American Statistician* **39**:83–87.
- Coop, G., K. Bullaughey, F. Luca and M. Przeworski. 2008a. The timing of selection at the human FOXP2 gene. *Mol Biol Evol* **25**:1257–1259.
- Coop, G., X. Wen, C. Ober, J. Pritchard and M. Przeworski. 2008b. High-resolution mapping of crossovers reveals extensive variation in fine-scale recombination patterns among humans. *Science* **319**:1395–1398.
- Cox, M. P., F. L. Mendez, T. M. Karafet, M. M. Pilkington, S. B. Kingan, G. Destro-Bisol, B. I. Strassmann and M. F. Hammer. 2008. Testing for archaic hominin admixture on the x chromosome: Model likelihoods for the modern human RRM2P4 region from summaries of genealogical topology under the structured coalescent. *Genetics* **178**:427–437.
- Csilléry, K., M. G. B. Blum, O. E. Gaggiotti and O. François. 2010. Approximate Bayesian Computation in practice. *Trends in Ecology & Evolution* **25**:410–418.
- DeGiorgio, M., M. Jakobsson and N. A. Rosenberg. 2009. Explaining worldwide patterns of human genetic variation using a coalescent-based serial founder model of migration outward from Africa. *Proc. Natl. Acad. Sci. USA* **106**:16057–16062.
- Depaulis, F., S. Mousset and M. Veuille. 2003. Power of neutrality tests to detect bottlenecks and hitchhiking. *J. Mol. Evol.* **57**:S190–S200.
- Evans, P., N. Mekel-Bobrov, E. Vallender, R. Hudson and B. Lahn. 2006. Evidence that the adaptive allele of the brain size gene microcephalin introgressed into homo sapiens from an archaic homo lineage. *Proceedings of the national academy of sciences USA* **103**:18178–18183.
- Fagundes, N. J. R., N. Ray, M. Beaumont, S. Neuenschwander, F. M. Salzano, S. L. Bonatto and L. Excoffier. 2007. Statistical evaluation of alternative models of human evolution. *Proc. Natl. Acad. Sci. USA* **104**:17614–17619.
- François, O., M. G. B. Blum, M. Jakobsson and N. A. Rosenberg. 2008. Demographic History of European Populations of *Arabidopsis thaliana*. *PLoS Genetics* **4**:e1000075.
- Garrigan, D. and M. F. Hammer. 2006. Reconstructing human origins in the genomic era. *Nature Reviews Genetics* **7**:669–680.
- Garrigan, D., Z. Mobasher, S. B. Kingan, J. A. Wilder and M. F. Hammer. 2005. Deep haplotype divergence and long-range linkage disequilibrium at Xp21.1 provide evidence that humans descend from a structured ancestral population. *Genetics* **177**:1849–1856.
- Gelman, A., J. B. Carlin, H. S. Stern and D. B. Rubin. 2003. *Bayesian Data Analysis, Second Edition (Texts in Statistical Science)*. Chapman & Hall/CRC, Boca Raton, 2nd edn.

- Glazko, G. and M. Nei. 2003. Estimation of divergence times for major lineages of primate species. *Molecular Biology and Evolution* **20**:424–434.
- Green, R. E., J. Krause, A. W. Briggs, T. Maricic, U. Stenzel, M. Kircher, N. Patterson, H. Li, W. Zhai, M. H.-Y. Fritz, N. F. Hansen, E. Y. Durand, A.-S. Malaspina, J. D. Jensen, T. Marques-Bonet, C. Alkan, K. Prufer, M. Meyer, H. A. Burbano, J. M. Good, R. Schultz, A. Aximu-Petri, A. Butthof, B. Hober, B. Hoffner, M. Siegemund, A. Weihmann, C. Nusbaum, E. S. Lander, C. Russ, N. Novod, J. Affourtit, M. Egholm, C. Verna, P. Rudan, D. Brajkovic, Z. Kucan, I. Gusic, V. B. Doronichev, L. V. Golovanova, C. Lalueza-Fox, M. de la Rasilla, J. Fortea, A. Rosas, R. W. Schmitz, P. L. F. Johnson, E. E. Eichler, D. Falush, E. Birney, J. C. Mullikin, M. Slatkin, R. Nielsen, J. Kelso, M. Lachmann, D. Reich and S. Paabo. 2010. A Draft Sequence of the Neandertal Genome. *Science* **328**:710–722.
- Gunz, P., F. L. Bookstein, P. Mitteroecker, A. Stadlmayr, H. Seidler and G. W. Weber. 2009. Early modern human diversity suggests subdivided population structure and a complex out-of-Africa scenario. *Proc. Natl. Acad. Sci. USA* **106**:6094–6098.
- 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**:e1000695.
- Harding, R. and G. McVean. 2004. A structured ancestral population for the evolution of modern humans. *Current Opinion in Genetics & Development* **14**:667–674.
- Harding, R. M., S. M. Fullerton, R. C. Griffiths, J. Bond, M. J. Cox, J. A. Schneider, D. S. Moulin and J. B. Clegg. 1997. Archaic African *and* Asian lineages in the genetic ancestry of modern humans. *Am. J. Hum. Genet.* **60**:772–789.
- Harding, R. M. and G. McVean. 2004. A structured ancestral population for the evolution of modern humans. *Curr. Op. Genet. Devel.* **14**:667–674.
- Harris, E. E. and J. Hey. 1999. X chromosome evidence for ancient human histories. *Proc. Natl. Acad. Sci. USA* **96**:3320–3324.
- Hayakawa, T., I. Aki, Y. S. A. Varki and N. T. 2006. Fixation of the human-specific CMP-N-acetylneuraminic acid hydroxylase pseudogene and implications of haplotype diversity for human evolution. *Genetics* **172**:1139–1146.
- Hein, J., M. H. Schierup and C. Wiuf. 2005. *Gene Genealogies, Variation and Evolution*. Oxford University Press, Oxford.
- Hudson, R. R. 2002. Generating samples under a Wright-Fisher neutral model of genetic variation. *Bioinformatics* **18**:337–338.
- . 2007. The variance of coalescent time estimates from DNA sequences. *Journal of Molecular Evolution* **64**:702–705.
- Ingman, M., H. Kaessmann, S. Pääbo and U. Gyllensten. 2000. Mitochondrial genome variation and the origin of modern humans. *Nature* **408**:708–713.

- Jakobsson, M., S. W. Scholz, P. Scheet, J. R. Gibbs, J. M. VanLiere, H. . Fung, Z. A. Szpiech, J. H. Degnan, K. Wang, R. Guerreiro, J. M. Bras, J. C. Schymick, D. G. Hernandez, B. J. Traynor, J. Simon-Sanchez, M. Matarin, A. Britton, J. Van De Leemput, I. Rafferty, M. Bucan, H. M. Cann, J. A. Hardy, N. A. Rosenberg and A. B. Singleton. 2008. Genotype, haplotype and copy-number variation in worldwide human populations. *Nature* **451**:998–1003.
- Jobling, M. and C. Tyler-Smith. 2003. The human Y chromosome: an evolutionary marker comes of age. *Nature Reviews Genetics* **4**:598–612.
- Kim, H. and Y. Satta. 2008. Population genetic analysis of the n-acylsphingosine amidohydrolase gene associated with mental activity in humans. *Genetics* **178**:1505–1515.
- Kreitman, M. 2000. Methods to detect selection in populations with applications to the human. *Annu. Rev. Genomics Hum. Genet.* **1**:539–559.
- Krings, M., A. Stone, R. W. Schmitz, H. Krainitzki, M. Stoneking and S. Pääbo. 1997. Neandertal DNA sequences and the origin of modern humans. *Cell* **90**:19–30.
- Lahr, M. and R. Foley. 1994. Multiple dispersals and modern human origins. *Evolutionary Anthropology* **3**:48–60.
- Laval, G., E. Patin, L. B. Barreiro and L. Quintana-Murci. 2010. Formulating a historical and demographic model of recent human evolution based on resequencing data from noncoding regions. *PLoS ONE* **5**:e10284.
- Li, J. Z., D. M. Absher, H. Tang, A. M. Southwick, A. M. Casto, S. Ramachandran, H. M. Cann, G. S. Barsh, M. Feldman, L. L. Cavalli-Sforza and R. M. Myers. 2008. Worldwide human relationships inferred from genome-wide patterns of variation. *Science* **319**:1100–1104.
- 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**:351–372.
- McDougall, I., F. Brown and J. G. Fleagle. 2005. Stratigraphic placement and age of modern humans from kibish, Ethiopia. *Nature* **433**:733–736.
- McVean, G. A. T., S. R. Myers, S. Hunt, P. Deloukas, D. R. Bentley and P. Donnelly. 2004. The fine-scale structure of recombination rate variation in the human genome. *Science* **304**:581–584.
- Mellars, P. 2006. Why did modern humans populations disperse from Africa ca. 60,000 years ago. *Proceedings of the national academy of sciences USA* **103**:9381–9386.
- Noonan, J. P. 2010. Neanderthal genomics and the evolution of modern humans. *Genome Research* **20**:547–553.
- Noonan, J. P., G. Coop, S. Kudaravalli, D. Smith, J. Krause, J. Alessi, D. Platt, S. Paabo, J. K. Pritchard and E. M. Rubin. 2006. Sequencing and analysis of Neanderthal genomic DNA. *Science* **314**:1113–1118.

- Patin, E., L. B. Barreiro, P. C. Sabeti, F. Austerlitz, F. Luca, A. Sajantila, D. M. Behar, O. Semino, A. Sakuntabhai, N. Guiso, B. Gicquel, K. McElreavey, R. M. Harding, E. Heyer and L. Quintana-Murci. 2006. Deciphering the ancient and complex evolutionary history of human arylamine N-acetyltransferase genes. *American Journal of Human Genetics* **78**:423–436.
- Plagnol, V. and J. D. Wall. 2006. Possible ancestral structure in human populations. *PLoS Genet.* **2**:110.1371/journal.pgen.0020105.eor.
- Prugnolle, F., A. Manica and F. Balloux. 2005. Geography predicts neutral genetic diversity of human populations. *Curr. Biol.* **15**:159–160.
- Ramachandran, S., O. Deshpande, C. C. Roseman, N. A. Rosenberg, M. W. Feldman and L. L. Cavalli-Sforza. 2005. Support from the relationship of genetic and geographic distance in human populations for a serial founder effect originating in Africa. *Proc. Natl. Acad. Sci. USA* **102**:15942–15947.
- Schaffner, S. F., C. Foo, S. Gabriel, D. Reich, M. J. Daly and D. Altshuler. 2005. Calibrating a coalescent simulation of human genome sequence variation. *Genome Res.* **15**:1576–1583.
- Schmegner, C., J. Hoegel, W. Vogel and G. Assum. 2005. Genetic variability in a genomic region with long-range linkage disequilibrium reveals traces of a bottleneck in the history of the european population. *Human Genetics* **118**:276–286.
- Shi, W., Q. Ayub, M. Vermeulen, R.-G. Shao, S. Zuniga, K. van der Gaag, P. de Knijff, M. Kayser, Y. Xue and C. Tyler-Smith. 2010. A Worldwide Survey of Human Male Demographic History Based on Y-SNP and Y-STR Data from the HGDP-CEPH Populations. *Molecular Biology and Evolution* **27**:385–393.
- Stringer, C. 2002. Modern human origins - progress and prospects. *Philosophical transactions of the royal society, London (B)* **357**:563–579.
- Takahata, N. 1993. Allelic genealogy and human evolution. *Mol. Biol. Evol.* **10**:2–22.
- Tang, H., D. O. Siegmund, P. Shen, P. J. Oefner and M. W. Feldman. 2002. Frequentist estimation of coalescence times from nucleotide sequence data using a tree-based partition. *Genetics* **161**:447–459.
- Templeton, A. 2002. Out of Africa again and again. *Nature* **416**:45–51.
- The International HapMap Consortium. 2005. A haplotype map of the human genome. *Nature* **437**:1299–1319.
- Thomson, R., J. K. Pritchard, P. Shen, P. J. Oefner and M. W. Feldman. 2000. Recent common ancestry of human Y chromosomes: evidence from DNA sequence data. *Proc. Natl. Acad. Sci. USA* **97**:7360–7365.
- Tishkoff, S. and M. Gonder. 2006. Human origins within and out of Africa. In M. Crawford, ed., *Anthropological Genetics: Theory Methods and Applications*, 337–379. Cambridge University Press.

- Verdu, P., F. Austerlitz, A. Estoup, R. Vitalis, M. Georges, S. Théry, A. Froment, S. Le Bomin, A. Gessain, J.-M. Hombert, L. Van Der Veen, L. Quintana Murci, S. Bahuchet and E. Heyer. 2009. Origins and genetic diversity of pygmy hunter-gatherers from Western Central Africa. *Current Biology* **19**:312–8.
- Voight, B. F., A. M. Adams, L. A. Frisse, Y. Qian, R. R. Hudson and A. Di Rienzo. 2005. Interrogating multiple aspects of variation in a full resequencing data set to infer human population size changes. *Proc. Natl. Acad. Sci. USA* **102**:18508–18513.
- Wall, J. 2003. Estimating ancestral population sizes and divergence times. *Genetics* **163**:395–404.
- Wall, J., M. Cox, F. Mendez, W. A., T. Severson and M. Hammer. 2008. A novel DNA sequence database for analyzing human demographic history. *Genome Research* **18**:1354–1361.
- Wall, J. D. and S. K. Kim. 2007. Inconsistencies in neanderthal genomic dna sequences. *PLoS Genet* **3**:e175.
- Wall, J. D., K. E. Lohmueller and V. Plagnol. 2009. Detecting ancient admixture and estimating demographic parameters in multiple human populations. *Mol Biol Evol* **26**:1823–1827.
- White, T. D., B. Asfaw, D. DeGusta, H. Gilbert, G. Richards, G. Suwa and F. C. Howell. 2003. Pleistocene *homo sapiens* from middle awash, Ethiopia. *Nature* **423**:742–747.
- Wilder, J. A., Z. Mobasher and M. F. Hammer. 2004. Genetic evidence for unequal effective population sizes of human females and males. *Mol. Biol. Evol.* **21**:2047–2057.
- Wolpoff, M., J. Hawks and R. Caspari. 2000. Multiregional, not multiple origins. *American Journal of Physical Anthropology* **112**:129–136.
- Yee, T. W. 2010. The VGAM package for categorical data analysis. *Journal of Statistical Software* **32**:1–34.

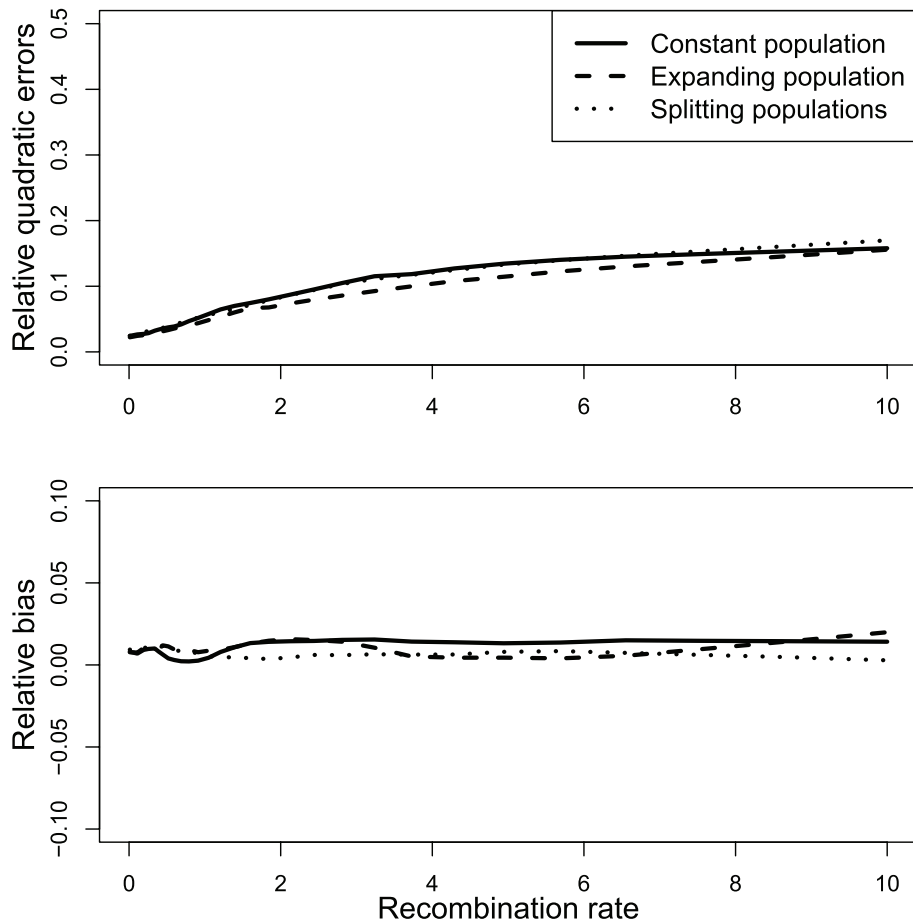


Figure 1: **Average errors of the TMRCA estimate for  $n = 100$  individuals.** We compare estimated TMRCA  $\hat{T}$  to true averaged TMRCA  $T_{av}$ . The true TMRCA are averaged using the median along the 20kb sequences. The relative bias is defined as the average over simulations of  $(\hat{T} - T_{av})/T_{av}$  and the relative quadratic error is the average of  $(\hat{T} - T_{av})^2/T_{av}^2$ . The effective recombination rate is equal to 4 times the ancestral population size times the recombination rate. In all simulations, we use the mutation rate estimated from the data.

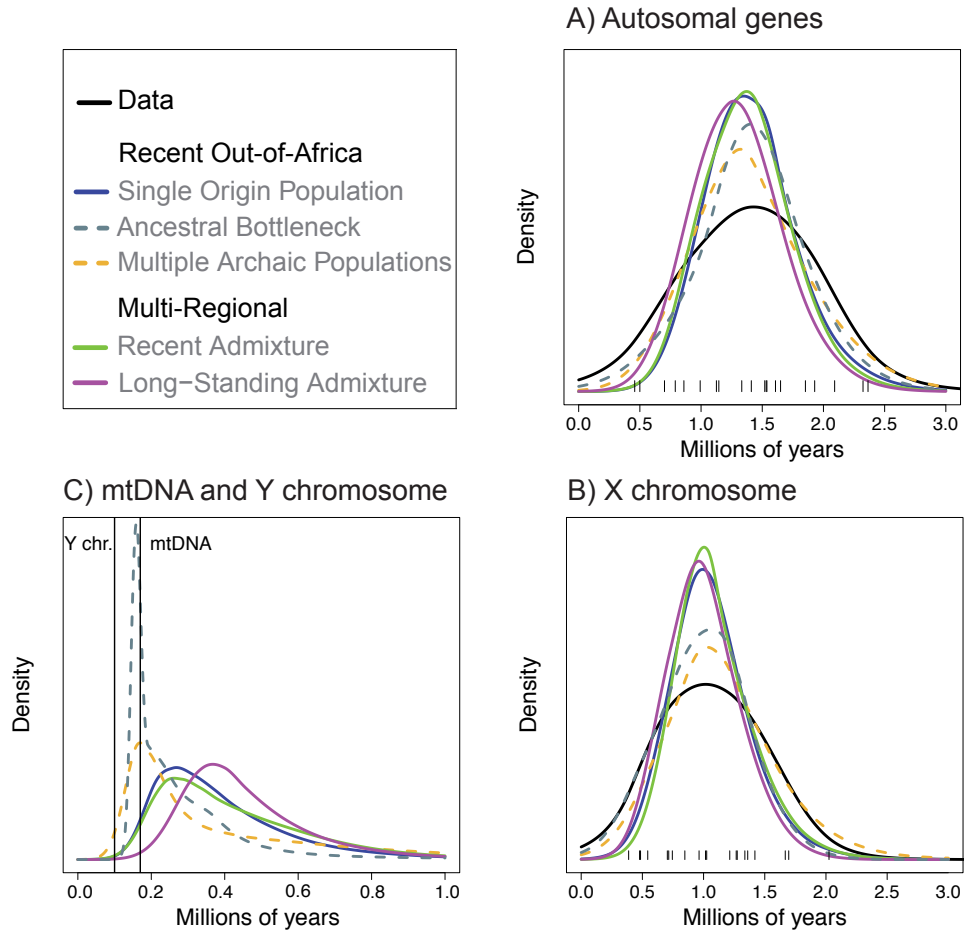


Figure 2: **Estimates of the TMRCA and posterior predictive distribution for the different models of human evolution.** A) The autosomes, B) the X-chromosome, and C) mtDNA and the Y-chromosome. The thick vertical lines correspond to TMRCA estimates, from the literature, for the Y-chromosome and the mtDNA (Ingman et al. 2000; Wilder et al. 2004). The simulated TMRCA have been obtained by computing the median of the TMRCA along the 20kb simulated sequences.



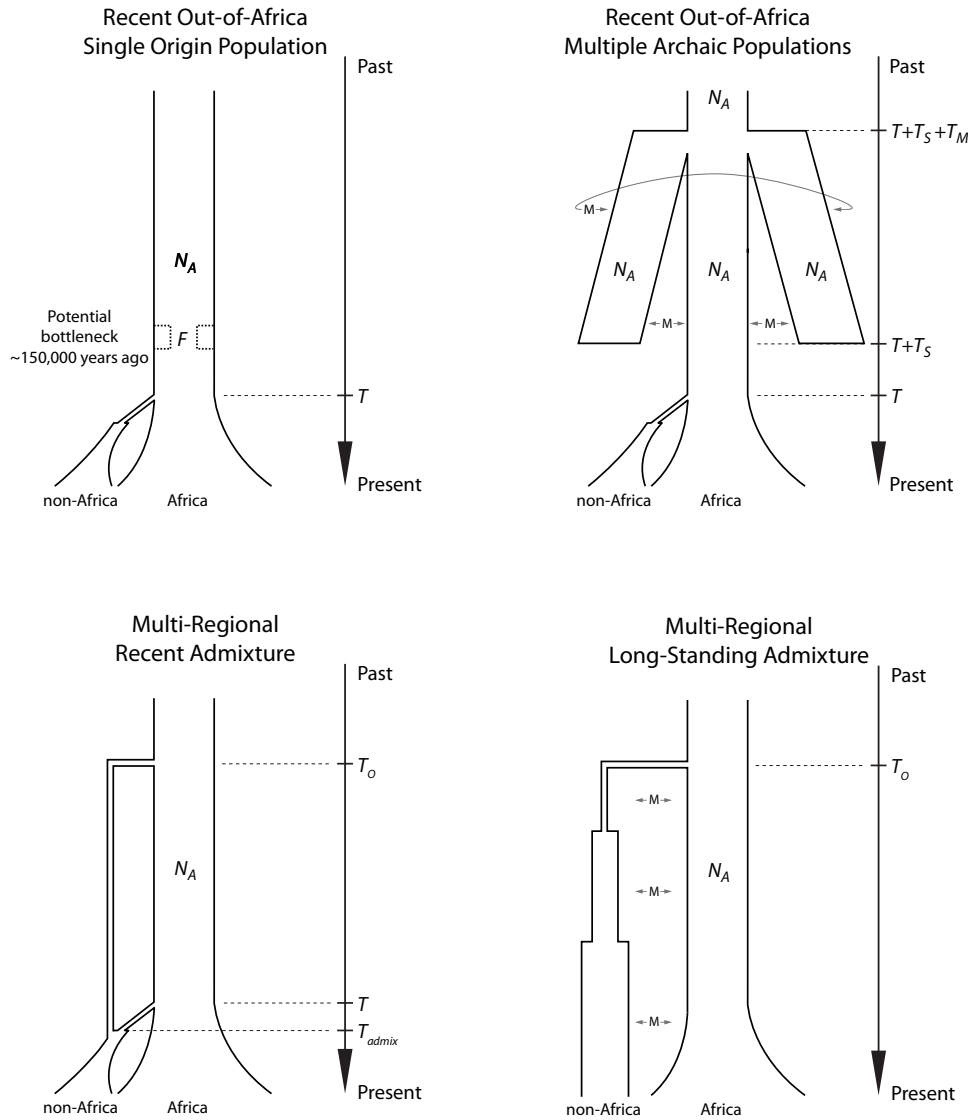


Figure 3: **Four different scenarios of human evolution.** 1) **'Single Origin Population'** A recent Out-of-Africa scenario in which modern humans descended from one subpopulation of archaic humans that was a separate population for a long time in Africa. The recent Out-of-Africa scenario potentially includes a bottleneck before the exodus from Africa, 150,000 years ago. 2) **'Multiple Archaic Populations'** A recent Out-of-Africa scenario in which different archaic African populations were connected by gene flow, even though only one archaic population eventually colonized the globe (Harding and McVean 2004; Garrigan and Hammer 2006). 3) **'Recent Admixture'** A Multi-Regional scenario in which archaic and modern humans were isolated during 300-600,000 years and admixed recently in Eurasia, 30-70,000 years ago (Plagnol and Wall 2006). 4) **'Long-Standing Admixture'** A Multi-Regional scenario with continuous and long-standing admixture between archaic and modern humans. Ancestral population size:  $N_A$ , time of the migration out of Africa:  $T$ , inbreeding coefficient during the bottleneck:  $F$ , time of structuring of archaic African population:  $T + T_S + T_M$ , ending of structured archaic African population:  $T + T_S$ , time of archaic humans exiting Africa:  $T_0$ , time of admixture:  $T_{admix}$ , and migration rate:  $M$ .

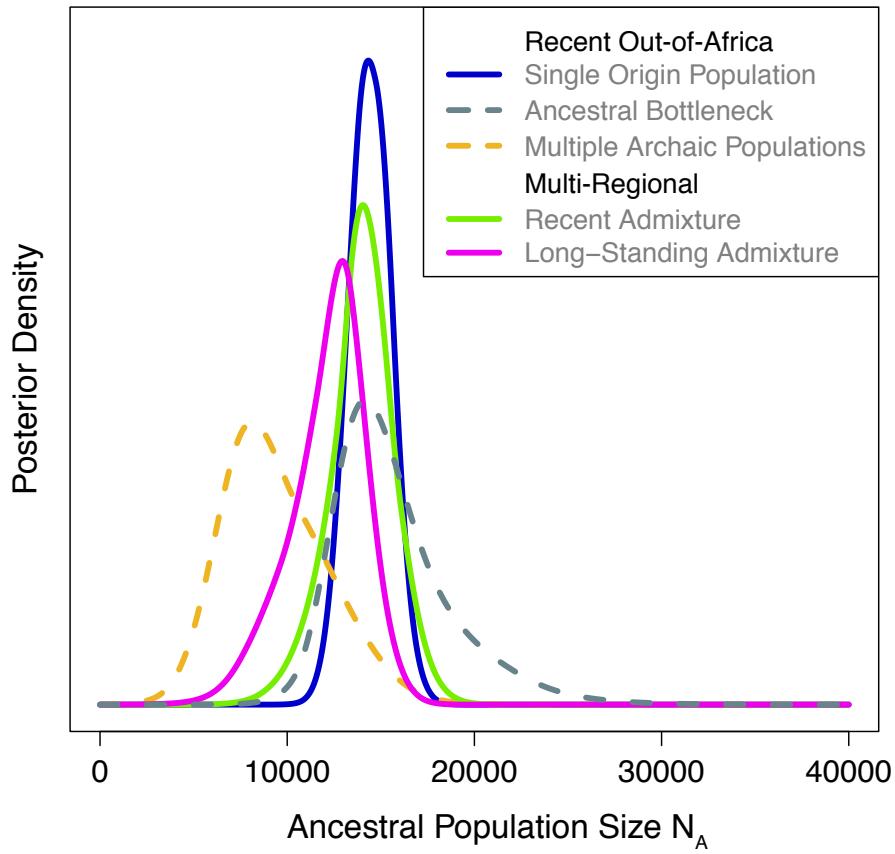


Figure 4: Posterior distribution of the ancestral effective population size  $N_A$ .

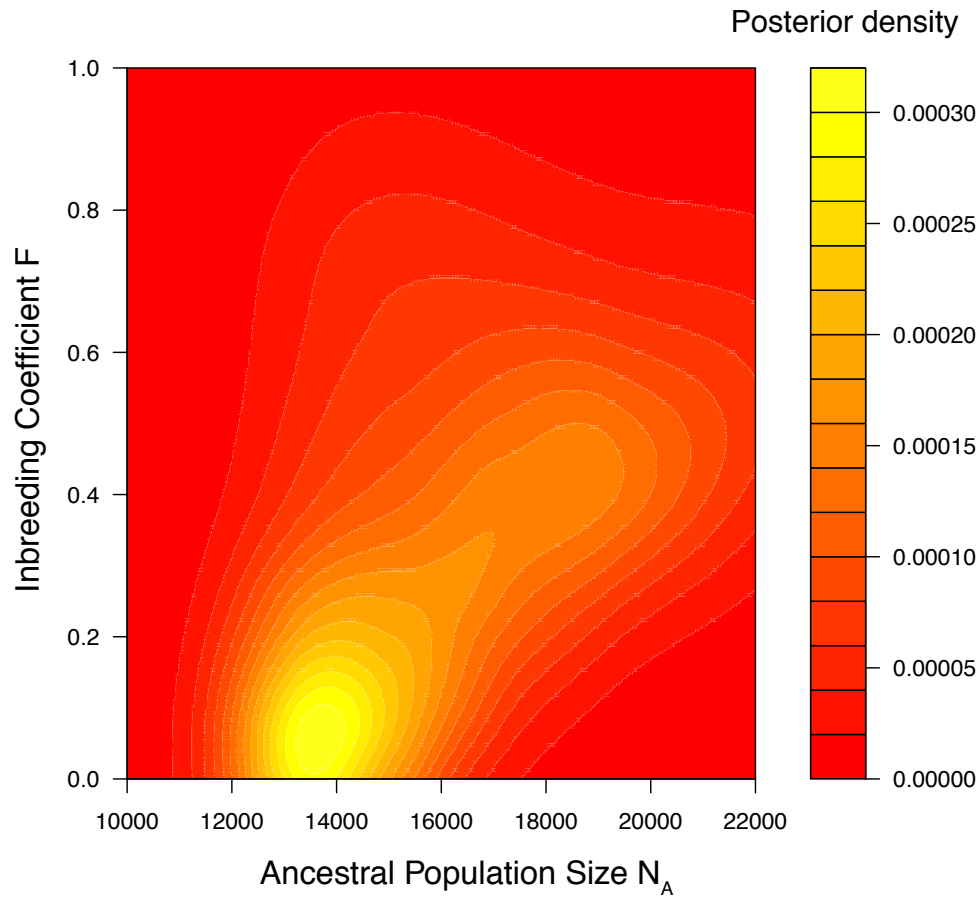


Figure 5: **Joint posterior distribution of the ancestral effective population size  $N_A$  and the inbreeding coefficient  $F$  during the ancestral bottleneck.**

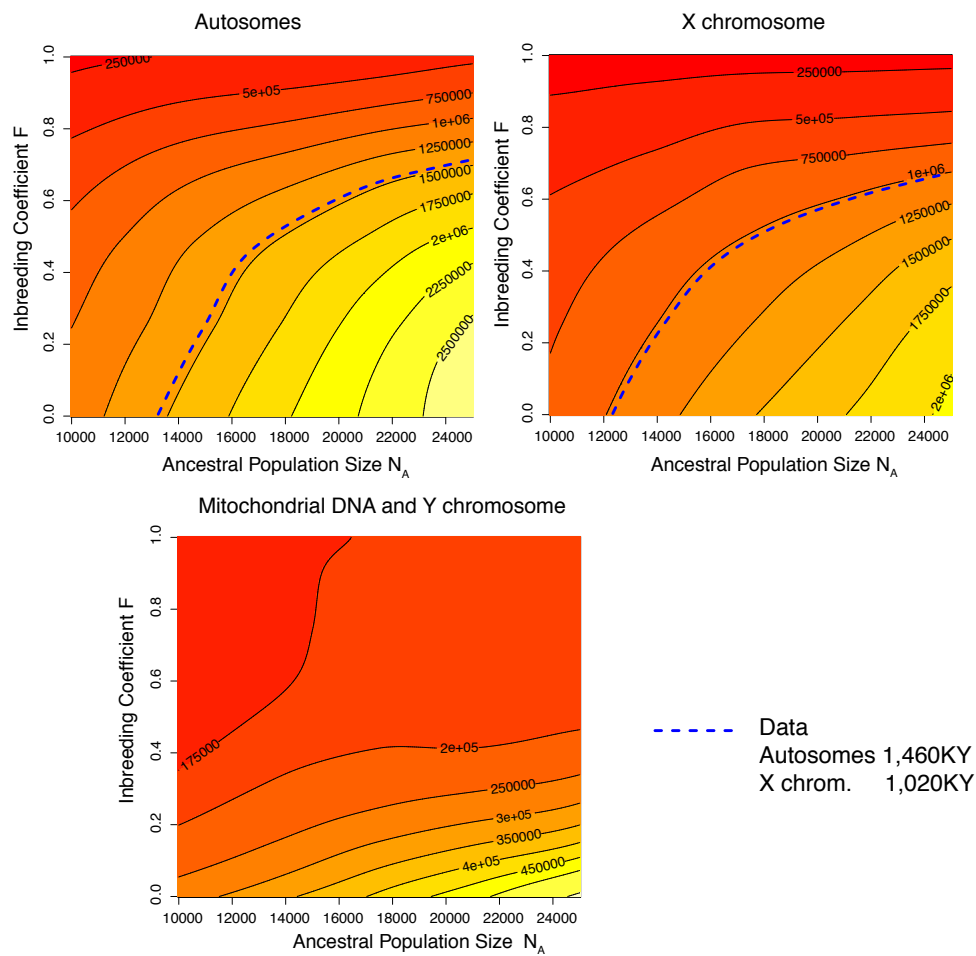


Figure 6: Mean TMRCA as a function of the inbreeding coefficient  $F$  and the ancestral population size  $N_A$  for autosomal, X-linked, and haploid genes.