

1 **Supplementary Information for**  
2 **Bayesian Inference Under the Multispecies Coalescent with ancient DNA sequences**

3 **Anna A. Nagel, Tomáš Flouri, Ziheng Yang, and Bruce Rannala**

4 **Corresponding Author: Anna A. Nagel, Bruce Rannala.**

5 **E-mail: [aanagel@ucdavis.edu](mailto:aanagel@ucdavis.edu), [brannala@ucdavis.edu](mailto:brannala@ucdavis.edu)**

6 **This PDF file includes:**

7     Supplementary text

8     Figs. S1 to S7

9     SI References

## 10 Supporting Information Text

### 11 1. Simulation method

12 The simulation method in BPP was modified to accommodate serial sampling. The dates are specified in units of expected  
13 number of substitutions and given in an input file. Simulation works similarly the standard MSC simulation with a few extra  
14 steps. When simulating the MSC without tip dates, times for the coalescent events are drawn from an exponential distribution  
15 with the rate determined by the number of lineages within a population. When a coalescent event occurs, two lineages are  
16 randomly chosen to coalesce and the number of lineages decreases by one. This continues until either there is only one sequence  
17 in the population or the time drawn is older than the population divergence time. In either case, the time is reset to be the  
18 population divergence time, the number of lineages from the two populations are combined and the simulation continues  
19 backward in time until the root population only has one lineage. With tip dates, the simulation starts with the youngest  
20 sample time, rather than at time zero. Every time a coalescent time is drawn, it must be checked if the time is older than  
21 either the population divergence time or next oldest sampling event. In the former case, the simulation proceeds in the same  
22 way as without tip dating. In the later case, the time is set to the next oldest sampling event to determine all of the lineages  
23 that the sampling event are added to the lineage count, and the simulation proceeds.

### 24 2. Bayesian Simulations

25 **A. MCMC settings.** Bayesian simulations were conducted with 3000 replicate datasets. The parameters are described in the  
26 main text. Each MCMC was sampled 400,000 times, sampling every 4 iterations with 80,000 iterations of burn-in.

27 **B. Convergence.** Two MCMCs were run for each dataset to check convergence. Convergence was checked by comparing  
28 posterior samples from the two MCMCs for each set of parameters. A two-sample t-test was used to compare the posterior  
29 means in the two chains.

$$t = \frac{\bar{X}_1 - \bar{X}_2}{s_p \sqrt{\frac{2}{n}}}$$

where

$$s_p = \sqrt{\frac{s_{X_1}^2 + s_{X_2}^2}{2}}$$

30 In the standard two-sample t-test,  $X_i$  are the sample means,  $s_{X_i}^2$  are the unbiased estimators of the variance and  $n$  is the  
31 sample size. There are  $2n - 2$  degrees of freedom. Since the samples were not independent, rather than using the total number  
32 of samples in the MCMC,  $n = 10,000$  was used as the sample size. The test was performed on estimates of all of the  $\theta$ s and  
33  $\tau$ s. If there was a significant difference between the samples for any variable, the run was considered to not have converged.  
34 Additionally, any pairs of MCMCs that had effective sample sizes lower than 200 for any the  $\theta$ s and  $\tau$ s were considered to not  
35 have converged. This resulted in 408 datasets with MCMCs that did not converge. All runs that did not converge were re-run  
36 with different seeds and a burnin of 200,000 iterations. Convergence was checked again using the same criteria. There were 213  
37 datasets that did not converge on this second analysis and these were excluded from the results (e.g. the plot summaries in Fig.  
38 S2, S3).

39 Assessing convergence of MCMC is non-trivial, and these methods of checking convergence were spot checked for MCMCs  
40 that did or did not converge. The trace plot, the effective sample sizes, and plots of kernel density estimation were further  
41 visually examined for these spot checked cases.

### 42 3. Simulations

43 **A. MCMC settings.** All MCMCs were sampled 400,000 of times, sampling every 4 of iterations. The burnin was 160,000  
44 iterations. Two independent MCMCs were run for each dataset. Convergence was checked comparing the results between the  
45 independent MCMCs. See Materials and Methods for details of simulations parameters.

46 **B. Convergence.** Convergence was checked using criteria similar to the Bayesian simulations, except that an  $n$  of 2000 was used  
47 in the two-sample t-test and differences in the means of all parameters ( $\theta$ s,  $\tau$ s,  $\tau^\Delta$ s, and  $\mu$ ) were required to not be significantly  
48 different between the replicate MCMCs. All parameters except  $\mu$  required an effective sample size of at least 200 in both  
49 MCMC replicates to be considered as converged. Runs that did not converge were re-run with different seeds and 600,000  
50 samples, sampled every 4th iteration. The burnin length was not changed. The same test was conducted after re-running the  
51 MCMCs, except that the ancestral population sizes and the root age in expected number of substitutions were not checked and  
52 a two-sample t-test sample size of  $n = 200$  was used. These parameters converged more slowly than other parameters, and  
53 were not central to the results. The root age in time before present was included in the convergence criteria and appeared to  
54 converge more quickly than root age in expected number of substitutions in some cases. The mitochondrial simulations and the  
55 recent population divergence simulations that did not meet the convergence criteria were removed from the results. These  
56 comprised no more than half of any set of 20 replicate simulations. The other MCMCs that did not meet the convergence  
57 criteria were re-run with different seeds and 1,200,000 samples, sampled every 4th iteration. These tended to be the larger

58 datasets with 500 or 2000 loci. The convergence was assessed again with the same test that was used for the first MCMC  
 59 re-runs (ancestral population sizes and the root age in expected number of substitutions were not checked and a two-sample  
 60 t-test sample size of  $n = 200$  was used). The simulations that did not meet the convergence criteria were removed from the  
 61 results. These simulations comprised no more than half of simulation replicates for any set of simulation parameters.

## 62 4. Empirical Analysis

63 **A. Priors for nuclear dataset.** To choose appropriate parameters for the root age prior, all of the loci were concatenated for  
 64 each species. The average pairwise divergence between sequences from the mammoth and elephant species and the mastodon  
 65 was calculated to specify a prior for the dataset with the mastodon. The average pairwise divergence between all pairs of  
 66 species that are not sisters was calculated to choose a prior for the dataset without the mastodon. Gaps were removed from  
 67 the two sequences being compared prior to calculating pairwise divergence and “n” was treated as a gap. When ambiguity  
 68 codes existed in the sequences, equal probability was given to all possible bases indicated in the ambiguity code. This method  
 69 will give an overestimate of root age, as the coalescent times must be older than the speciation time. However, this should give  
 70 a reasonable order of magnitude for the prior mean. The variance was chosen such that there was a broad distribution around  
 71 the mean, since there is not strong prior information about the speciation times in expected number of substitutions.

Species 1	Species 2	pairwise divergence
Asian	Forest	0.0072
Asian	Savannah	0.0070
Mammoth	Forest	0.0069
Mammoth	Savannah	0.0068
Asian	Mastodon	0.037
Forest	Mastodon	0.036
Mammoth	Mastodon	0.036
Savannah	Mastodon	0.036

73 To obtain a prior for  $\theta$ , the pairwise divergence between within a population was calculated for all populations with unphased  
 74 data. Sites with ambiguity codes were considered to be heterozygous in the individual and not due to sequencing error. As  
 75 before, concatenated sequences were used and all gaps were removed prior to calculating the pairwise divergence. The  $\mu$  prior  
 76 was chosen to have a mean of  $5 \times 10^{-9}$  based on the priors and justifications used in previous analyses of this dataset (1).

Species	pairwise divergence
Asian	0.0012
Forest	0.0024
Mammoth	0.0008
Savannah	0.0006

78 **B. Priors for mitochondrial dataset.** The  $\theta$  prior was determined by calculating the average pairwise divergence between all  
 79 contemporary samples within a species across all possible pairs. Gaps were removed prior to calculating pairwise divergence. A  
 80 relatively broad prior was chosen to reflect the large difference in average pairwise divergence in the different species.

Species	pairwise divergence
Asian	0.0034
Forest	0.013
Savannah	0.026

82 To find a prior for the root  $\tau$ , the average pairwise divergence was found between all pairs of Asian and Forest elephants  
 83 sequences and Asian and Savannah elephant sequences. The prior was chosen to have a mean close to the average pairwise  
 84 divergence, with a relatively large variance to reflect the prior uncertainty in the parameter value.

Species 1	Species 2	pairwise divergence
Asian	Forest	0.047
Asian	Savannah	0.048

86 **C. MCMC settings.** MCMCs for the empirical analyses of both the nuclear and mitochondrial datasets were sampled 400,000 of  
 87 times, sampling every 4 of iterations. The burnin was 160,000 iterations. Two and four independent MCMCs were run for each  
 88 nuclear and mitochondrial dataset, respectively.

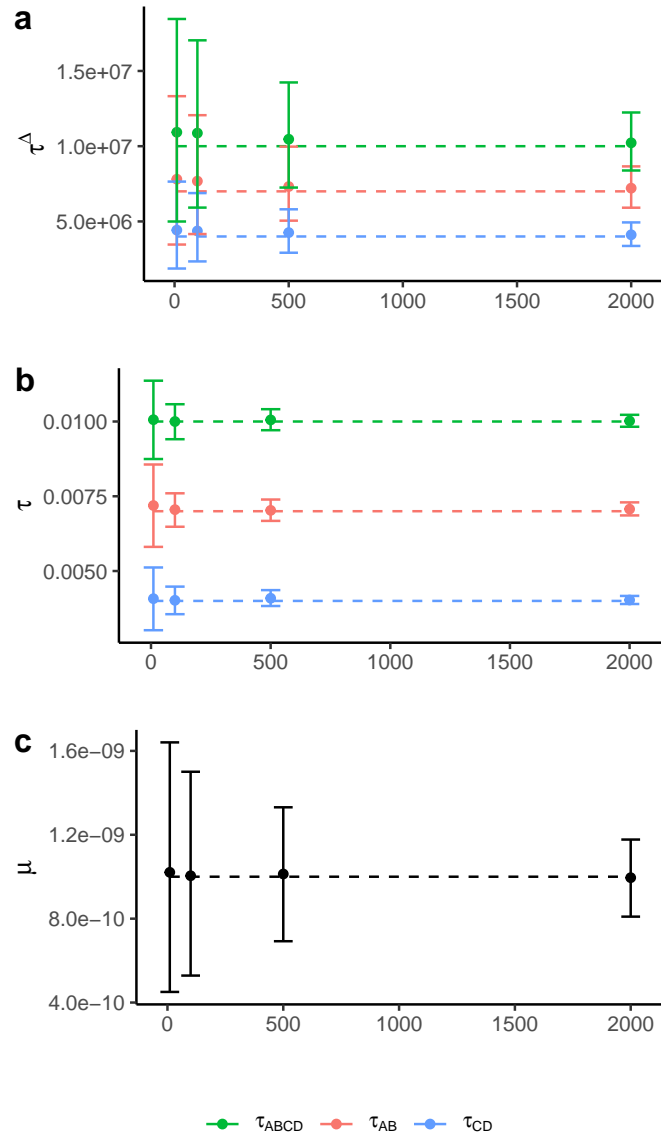
89 **D. Convergence.** Convergence was assessed in tracer by comparing the distributions of all parameters in the pairs of replicate  
 90 MCMCs and examining the trace plot.

Accession No.	Species	Age
KY616982.1	<i>Loxodonta africana</i>	modern
KY616977.1	<i>Loxodonta africana</i>	modern
KY616974.1	<i>Loxodonta africana</i>	modern
AB443879.1	<i>Loxodonta africana</i>	modern
MT636097.1	<i>Loxodonta cyclotis</i>	1533 (417 ybp)
MT636095.1	<i>Loxodonta cyclotis</i>	1533 (417 ybp)
MT636093.1	<i>Loxodonta cyclotis</i>	1533 (417 ybp)
KY616981.1	<i>Loxodonta cyclotis</i>	modern
KY616980.1	<i>Loxodonta cyclotis</i>	modern
KY616975.1	<i>Loxodonta cyclotis</i>	modern
KJ557423.1	<i>Loxodonta cyclotis</i>	modern
NC_020759.1	<i>Loxodonta cyclotis</i>	modern
DQ316068.1	<i>Elephas maximus</i>	modern
OP575307.1	<i>Elephas maximus</i>	modern
OL628830.1	<i>Elephas maximus</i>	modern

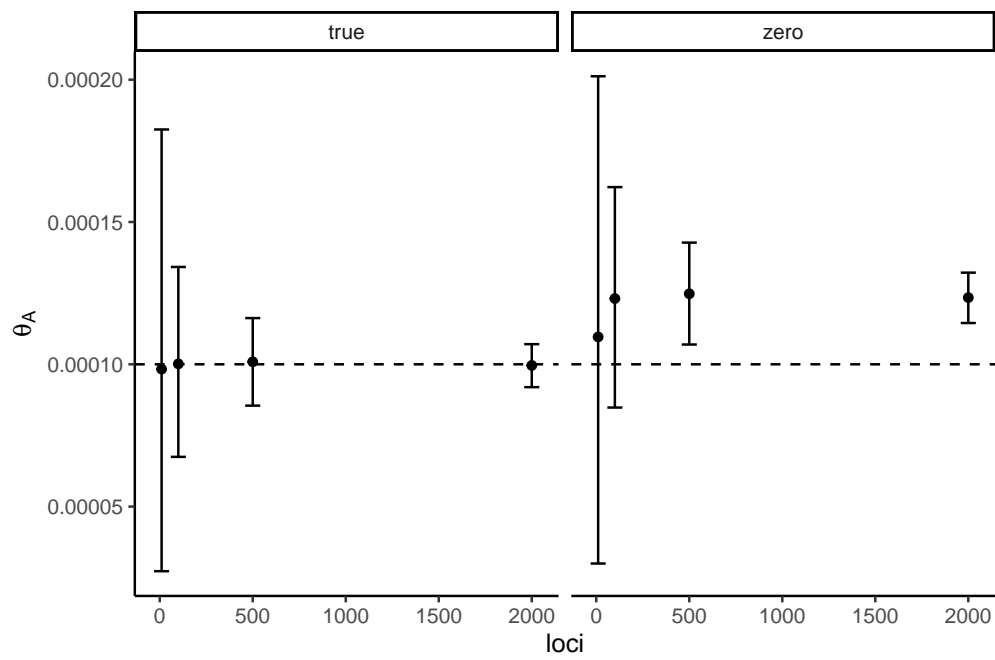
**Fig. S1.** Additional samples used in the mitochondrial analysis downloaded from GenBank.

91 **References**

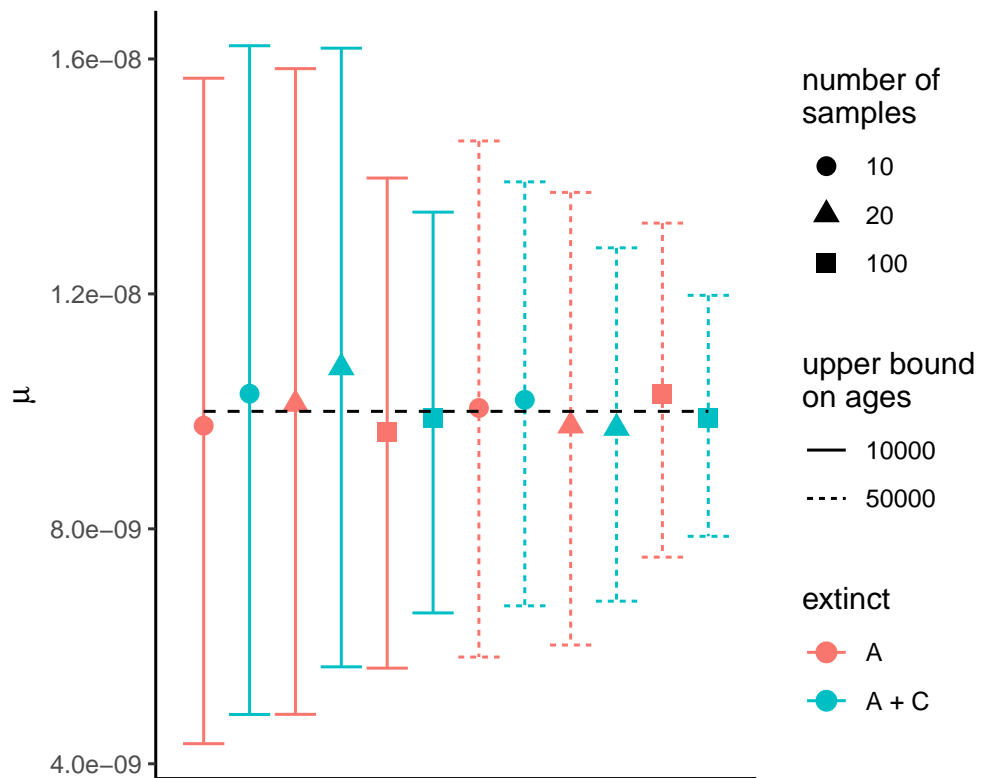
- 92 1. N Rohland, et al., Genomic DNA sequences from mastodon and woolly mammoth reveal deep speciation of forest and  
93 savanna elephants. *PLoS Biol.* **8**, e1000564 (2010).



**Fig. S2.** Average posterior means and 95% HPD CIs (bars), over 20 replicate datasets, of (a) divergence times in mutations, (b) divergence times in years, and (c) mutation rate. The data were simulated under the model of figure 1a with two extinct species (*A* and *C*), sample dates are between 5,000 and 50,000 years, and  $\theta = 0.001$ .

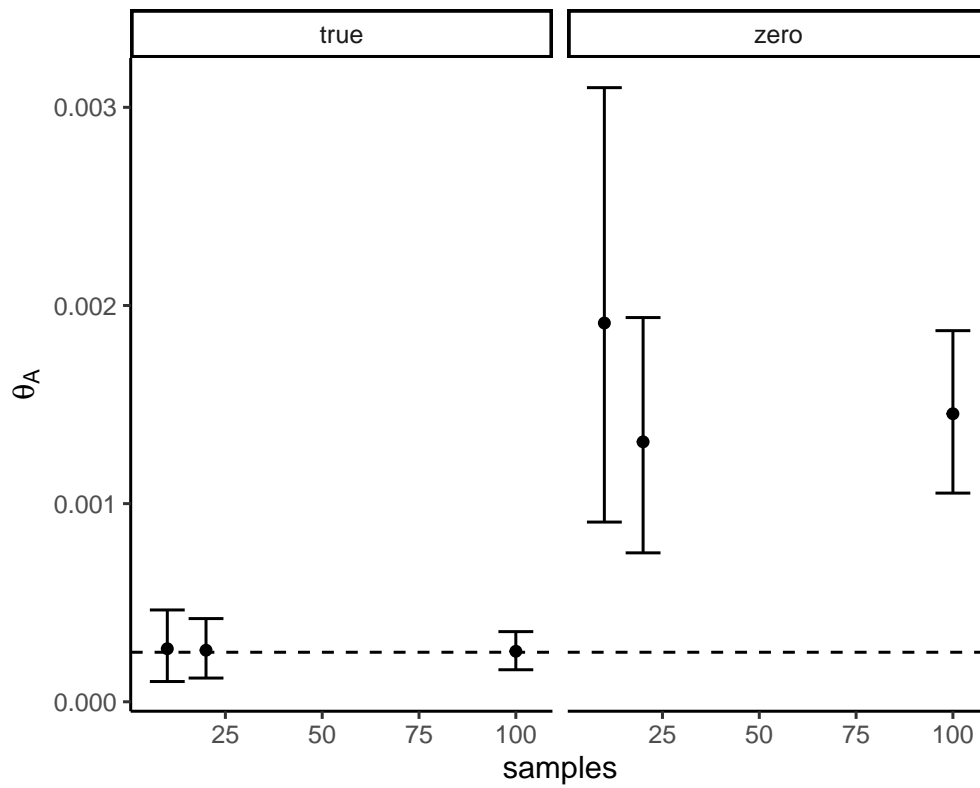


**Fig. S3.** Average posterior means and 95% HPD CIs (bars), over 40 replicate nuclear datasets, of  $\theta_A$  when sample dates were set to their true values (left) or zero (right). The datasets had 6 samples in each extinct species and the upper bound on the sample dates equal to 50,000 ybp. The dashed lines show the true values of the  $\theta_A$ .

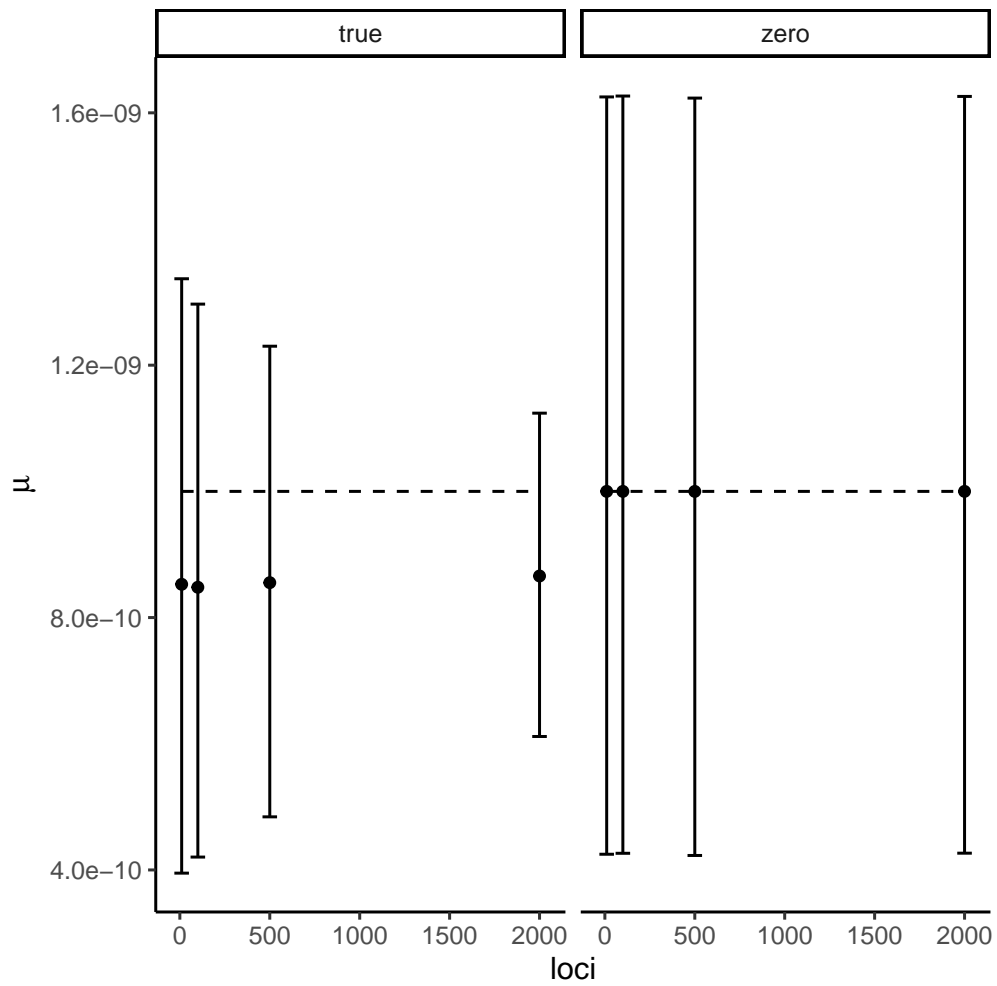


**Fig. S4.** Average posterior means and 95% HPD CIs (bars) of the mutation rate over 20 replicate datasets, simulated under the model of figure 1a with  $\theta = 0.00025$ . Solid lines are for sample dates between 5,000 to 10,000 ybp while dashed lines are for sample dates between 5,000 and 50,000 ybp. Either species A (red) or both A and C (teal) are extinct, and from each extinct species either 10 (circle), 20 (triangle), or 100 (square) samples are taken. The dashed line shows the true values of the  $\mu$ .

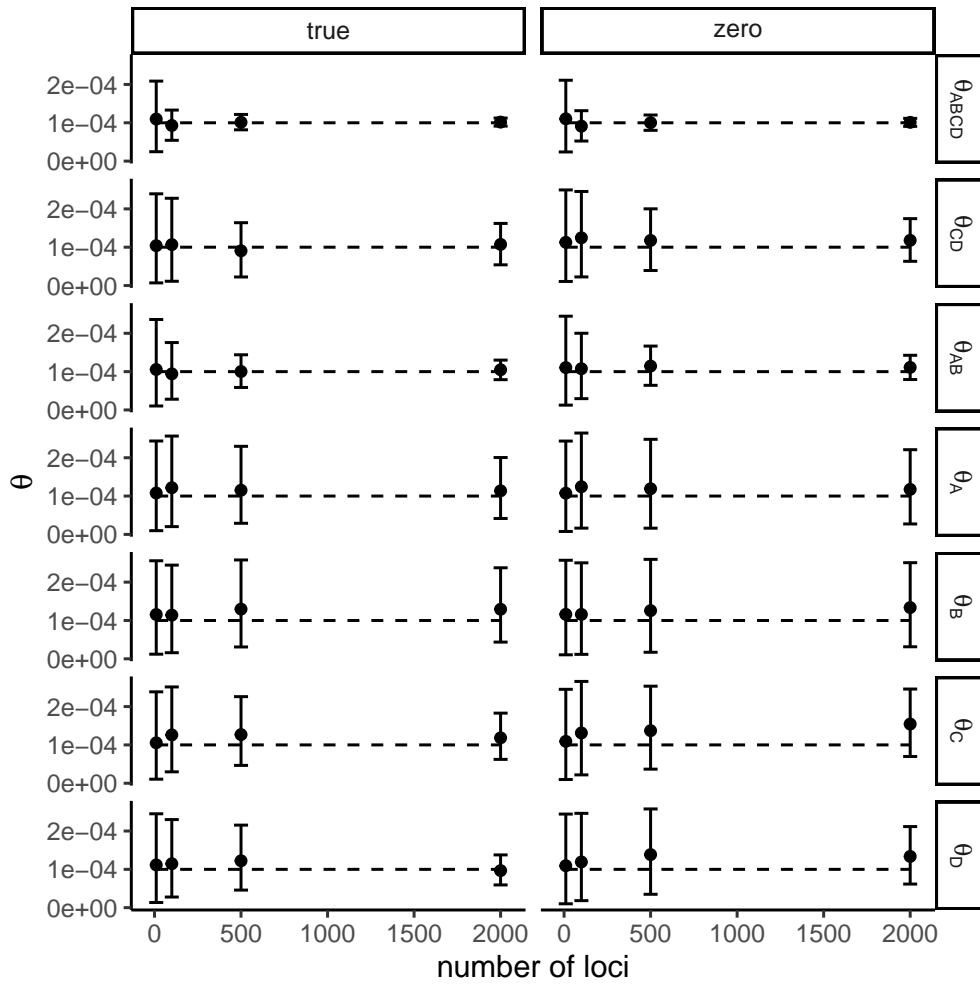




**Fig. S5.** Average posterior means and 95% HPD CIs (bars), over 40 replicate mitochondrial datasets, of  $\theta_A$  when sample dates were set to their true values (left) or zero (right). The datasets had the upper bound on the sample dates equal to 50,000 ybp. The dashed lines show the true values of the  $\theta_A$ .



**Fig. S6.** Each point shows the average of the posterior mean mutation rate and mean 95% credible set averaged across inferences for 20 replicate datasets when the samples ages are set to their true value (left) or zero (right). For all datasets, there were 2000 loci and  $\theta$  is equal to 0.0001. The dashed line shows the true value of  $\mu$ .



**Fig. S7.** Average posterior mean and 95% HPD CIs (bars) for  $\theta$  across 20 replicate simulations for the recent population divergence analysis. The left and right plots show the inferences when the sample ages are set to their true values and zero, respectively. The dashed lines show the true value of  $\theta$ .