

Supplementary Method: Site Frequency Spectrum for Variable Population Size

In this note we determine the site frequency spectrum (SFS) for a variable population size model, and apply the method to investigate a specific demographic model for three subspecies of chimpanzees. The note falls in three parts: In Section 1 we develop the general theory for obtaining the SFS for a variable population size model. The main problem is to determine the mean time between coalescence events, and in Section 2 we describe a Monte Carlo procedure for obtaining these mean coalescence times. Finally in Section 3 we consider the demographic model and corresponding SFS for three chimpanzee subspecies.

1 Simulating waiting times between coalescence events

Consider a tree with n leaves and denote by ℓ the number of branches that are present in the tree at a given time point. We call ℓ ($2 \leq \ell \leq n$) the level of the tree.

The probability $P_{i|\ell}$ that a mutation at level ℓ is present in i samples (the probability that a branch at level ℓ is ancestral to i leaves) is given by (see e.g. Durrett, 2008, equation (2.1) page 54)

$$P_{i|\ell} = \frac{\binom{n-i-1}{\ell-2}}{\binom{n-1}{\ell-1}}, \quad i = 1, \dots, n - \ell. \quad (1)$$

Let ξ_i ($1 \leq i \leq n - 1$) be the number of sites with i derived alleles and $(n - i)$ ancestral alleles. Conditional on the times (T_2, \dots, T_n) between coalescence events we have

$$\mathbb{E}[\xi_i | (T_2, \dots, T_n)] = \sum_{\ell=2}^{n-i+1} \ell \frac{\theta}{2} T_\ell P_{i|\ell},$$

where $P_{i|\ell}$ is defined above and θ is the scaled mutation rate. We get

$$\mathbb{E}[\xi_i] = \theta \sum_{\ell=2}^{n-i+1} \ell \frac{\binom{n-i-1}{\ell-2}}{\binom{n-1}{\ell-1}} \mathbb{E}[T_\ell],$$

so all we need in order to determine the site frequency count for a variable population size model is to determine $\mathbb{E}[T_\ell]$ ($2 \leq \ell \leq n$).

Using a time change we can easily *simulate* the times between coalescence events in a variable population size model (e.g. Tavaré, 2004, Section 2.4, Algorithm 2.1). We therefore use a Monte Carlo estimate to determine the mean between coalescent times.

2 Simulating waiting times between coalescence events

For convenience we re-state Algorithm 2.1 in Tavaré (2004) below.

Algorithm: Simulating times between coalescence events.

Input: Sample size n and variable population size with integrated intensity function Λ .

Output: Sample of times T_2, \dots, T_n between coalescence events.

Step 1: Generate $t'_j = -2 \log(U_j)/(j(j-1))$, $j = 2, \dots, n$.

Step 2: Form $s'_j = t'_j + \dots + t'_n$ $j = 2, \dots, n$.

Step 3: Compute $t_n = \Lambda^{-1}(s'_n)$, $t_j = \Lambda^{-1}(s'_j) - \Lambda^{-1}(s'_{j+1})$, $j = n - 1, \dots, 2$.

Step 4: Return $T_j = t_j$, $j = 2, \dots, n$.

Here the integrated intensity function Λ is defined as follows. Let N_0 be the scaling factor such that the relative population size function in scaled time t (generation tN_0 ago) is given by

$$f(t) = N(tN_0)/N_0, \quad t \geq 0.$$

We then define the intensity function

$$\lambda(u) = \frac{1}{f(u)}, \quad u \geq 0,$$

and the integrated intensity function

$$\Lambda(t) = \int_0^t \lambda(u) du, \quad t \geq 0.$$

For the chimpanzee data we are particularly interested in variable population size models where the population size is constant for two or three epochs.

2.1 Two epochs of constant size

Consider the demographic scenario depicted in Figure 1. Present population size is N , and the population size stays at N until time aN in the past. From time aN and further back in time the population size is a fraction α of the present population size.

Choosing $N_0 = N$ we get

$$f(x) = \begin{cases} 1 & 0 \leq x < a \\ \alpha & x \geq a \end{cases}$$

and

$$\Lambda^{-1}(u) = \begin{cases} u & 0 \leq u < a \\ a + \alpha(u - a) & u \geq a. \end{cases}$$

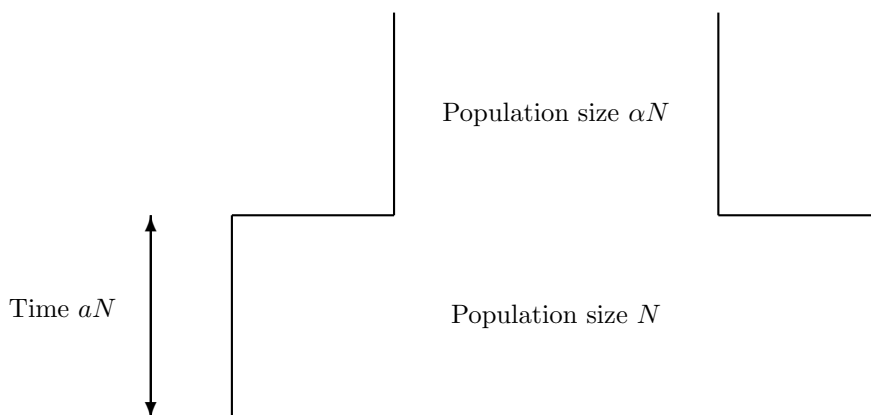


Figure 1: Backwards in time the population size is N in the first epoch and αN in the second epoch. The change in population size happens at time a in the past.

2.2 Three epochs of constant size

With three epochs of constant size the relative size function becomes

$$f(x) = \begin{cases} 1 & 0 \leq x < a \\ \alpha & a \leq x < b \\ \beta & x \geq b \end{cases}$$

and

$$\Lambda^{-1}(u) = \begin{cases} u & 0 \leq u < a \\ a + \alpha(u - a) & a \leq u < a + \frac{b-a}{\alpha} \\ b + \beta[(u - a) - \frac{(b-a)}{\alpha}] & u \geq a + \frac{b-a}{\alpha}. \end{cases}$$

3 Demographic model for chimpanzees

In Figure 2 we show the observed and estimated site frequency spectra for the three subspecies of chimpanzees. The estimated site frequency spectra are based on the demographic model obtained from the ABC procedure. In Figure 3 we show the observed and estimated folded site frequency spectra from the central, eastern and western chimpanzees.

References

- Tavare, S. (2004). *Ancestral inference in population genetics*. Lectures on Probability Theory and Statistics. Volume 1837, 2004, pp 1-188. Springer, Berlin Heidelberg.
- Durrett, R. (2008). *Probability Models for DNA Sequence Evolution*. 2nd edition. Springer, New York.

1 **List of supplementary tables & figures**

2

3 **Supplementary Table 1 Geographic origin and depth of sequencing of**
 4 **individuals analyzed.**

5

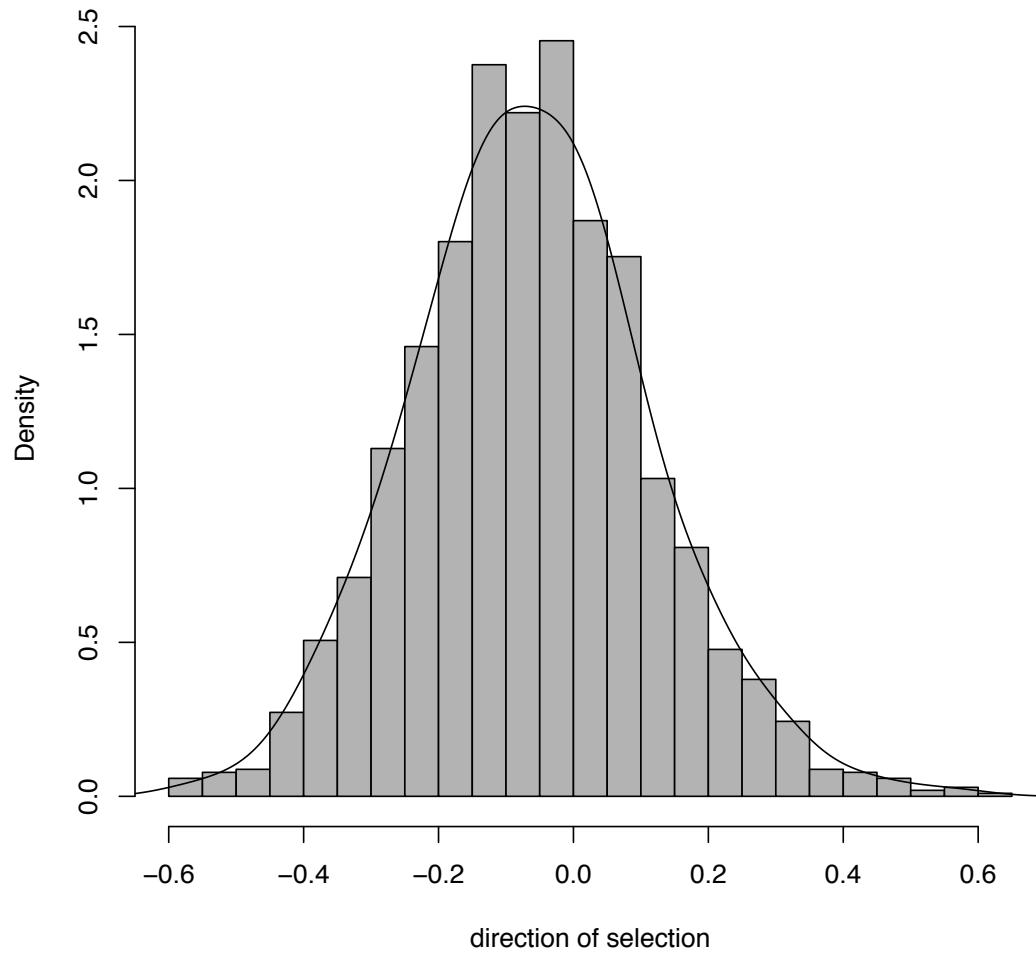
Individual ID	Sex	Origin	Subspecies	Avg. Seq. depth
Yoyo	F	Democratic Republic of Congo	<i>Pan troglodytes schweinfurthii</i>	32.9
Umugenzi	M	Democratic Republic of Congo	<i>Pan troglodytes schweinfurthii</i>	32.2
Ikuru	F	Democratic Republic of Congo	<i>Pan troglodytes schweinfurthii</i>	38.2
Tumbo	M	Kibale (Uganda)	<i>Pan troglodytes schweinfurthii</i>	32.2
Natasha	F	Kibale (Uganda)	<i>Pan troglodytes schweinfurthii</i>	35.5
Nkuumwa	F	Uganda	<i>Pan troglodytes schweinfurthii</i>	38.9
Pasa	F	Democratic Republic of Congo	<i>Pan troglodytes schweinfurthii</i>	33.1
Cindy	F	Budongo (Uganda)	<i>Pan troglodytes schweinfurthii</i>	28.8
Sunday	M	Democratic Republic of Congo	<i>Pan troglodytes schweinfurthii</i>	33.4
Exota (11785)	F	Wild caught	<i>Pan troglodytes schweinfurthii</i>	36.7
Paula (11784)	F	Wild caught	<i>Pan troglodytes schweinfurthii</i>	46.2
Susi (11043)	F	Wild caught	<i>Pan troglodytes troglodytes</i>	36.6
Cindy (11525)	F	Wild caught	<i>Pan troglodytes troglodytes</i>	46.2
Aboume	M	Gabon	<i>Pan troglodytes troglodytes</i>	
Amelie	F	Gabon	<i>Pan troglodytes troglodytes</i>	35.9
Ayrton	M	Moanda (Gabon)	<i>Pan troglodytes troglodytes</i>	37.3
Bakoumba	M	Gabon	<i>Pan troglodytes troglodytes</i>	37.6
Benefice	F	Makokou (Gabon)	<i>Pan troglodytes troglodytes</i>	40.1
Chiquita	F	Gabon	<i>Pan troglodytes troglodytes</i>	37.7
Lalala	F	Libreville (Gabon)	<i>Pan troglodytes troglodytes</i>	29.8
Makokou	F	Ogooué Ivindo (Gabon)	<i>Pan troglodytes troglodytes</i>	38.8
Masuku	F	Haut Ogooué (Gabon)	<i>Pan troglodytes troglodytes</i>	34.6
Noemie	F	Equatorial Guinea	<i>Pan troglodytes troglodytes</i>	34.4
Sita (11262)	F	1.gen of wild caught parents from Liberia	<i>Pan troglodytes verus</i>	37.4
Sepp-Toni (11300)	M	Wild Caught Liberia	<i>Pan troglodytes verus</i>	33.5
Olga (12314)	F	1.generation individual of wild caught parents from West Africa.	<i>Pan troglodytes verus</i>	38.4
Moreno (12341)	M	Wild caught West Africa	<i>Pan troglodytes verus</i>	32.7
Agneta (11758)	F	1.generation individual of wild caught parents from West Africa.	<i>Pan troglodytes verus</i>	34.2
Frits (11052)	M	Wild caught Sierra Leone	<i>Pan troglodytes verus</i>	35.1

6

7

8

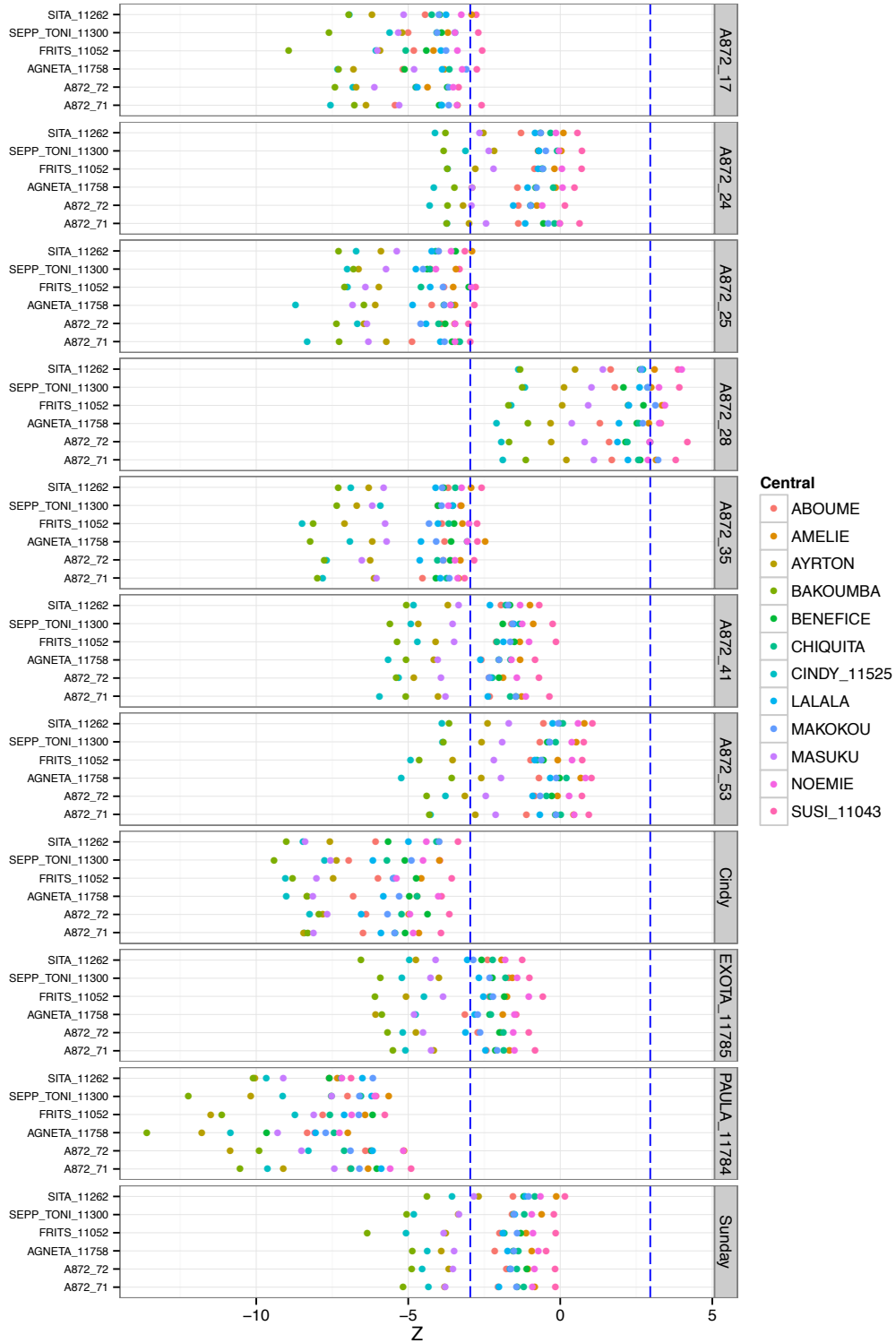
9 **Supplementary Figure 1: Non-smoothed histogram of DoS values**
10 **observed across autosomal windows in central chimpanzees.**



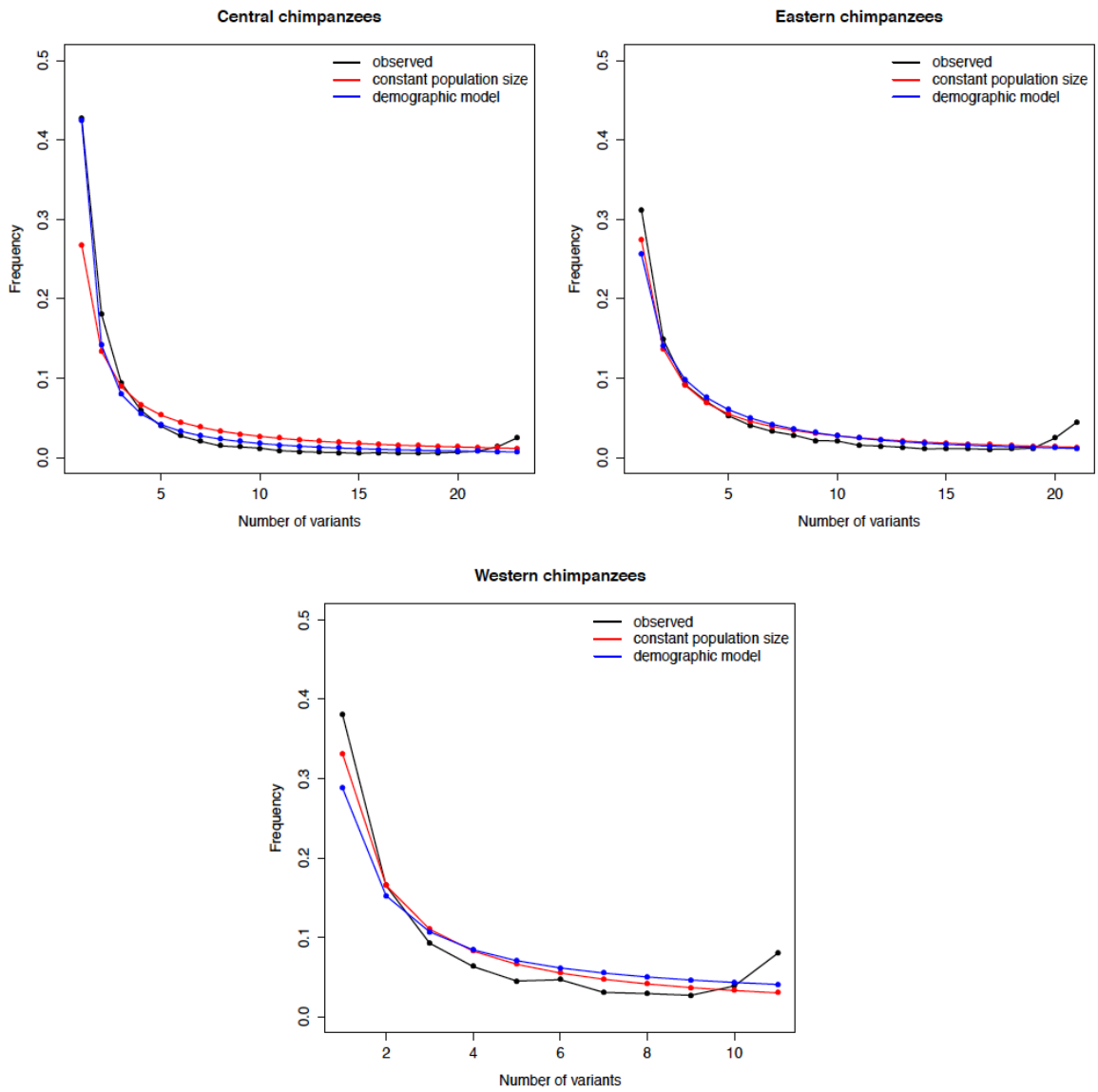
11
12

13 **Supplementary Fig. 2 Distribution of D statistics for all three-way**
 14 **comparison of individuals**

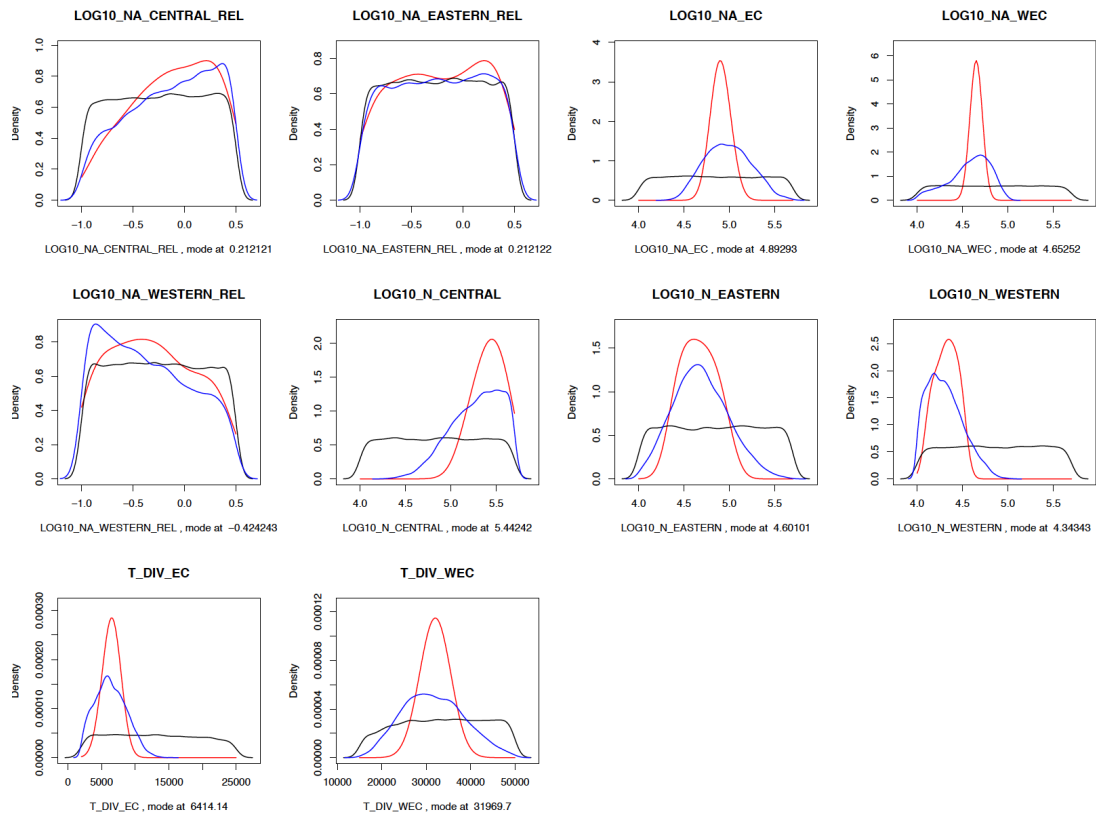
15 The plot shows all D statistics as Z scores with the significance threshold
 16 shown as vertical dashed lines. Each dot represents one three-way
 17 comparison. The different sub-plots correspond to different eastern
 18 chimpanzees, the different y-axis lines different western chimpanzees, and the
 19 different colors different central chimpanzees.
 20



22 **Supplementary Figure 3.** Observed versus predicted synonymous
23 autosomal SFS for each subspecies.
24
25



26
27 Predicted SFS under a constant population model (red) and the fitted
28 demographic model (blue) reported in Figure 3.

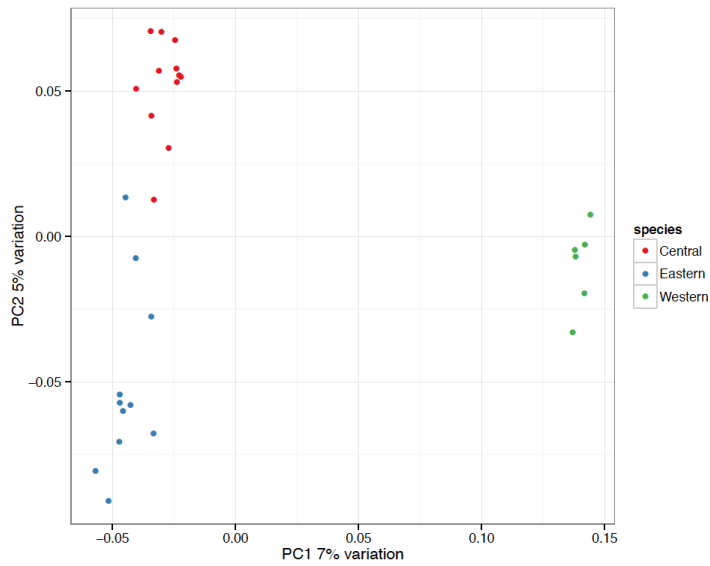


29
 30
 31
 32
 33
 34
 35
 36
 37

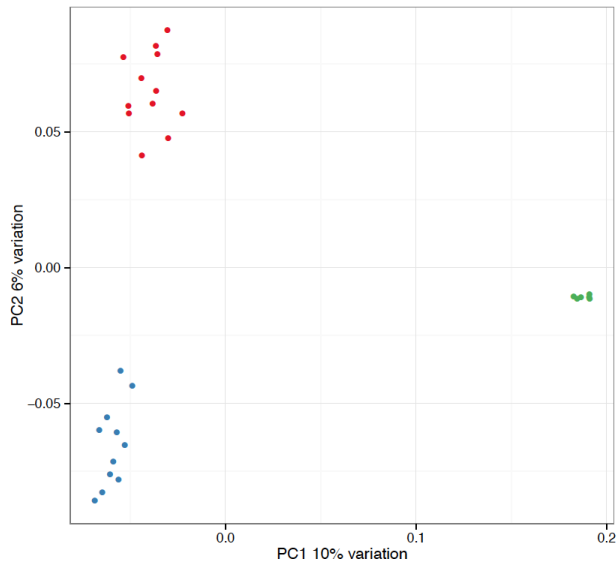
Supplementary Figure 4

Prior (black) distribution and posterior distribution before (blue) and after post sampling adjustment (red) of parameters of the demographic model fitted to autosomal SNP data using our ABC procedure.

38
39
40



41



42
43
44
45

46 **Supplementary Figure 5: comparison of PCA of indels versus SNPs.**

47
48
49
50
51
52
53
54

To factor out the effect of sample size on the differentiation of populations revealed by PCA analysis of SNPs and indels, we down sampled the SNP data set to the size of the indel data set (n=2073) and redid the PCA analysis. The analysis still shows a smaller degree of population differentiation in the indels (TOP) than the SNP (BOTTOM). Note that to make the plots directly comparable the resulting PCA coordinates were normalized with the sum of standard deviations of the principal components.