

Fine-scale Inference of Ancestry Segments without Prior Knowledge of Admixing Groups

Michael Salter-Townshend^{*,1} and Simon Myers[†]

^{*}University College Dublin, [†]University of Oxford and Wellcome Trust Centre for Human Genetics

ABSTRACT We present an algorithm for inferring ancestry segments and characterizing admixture events, which involve an arbitrary number of genetically differentiated groups coming together. This allows inference of the demographic history of the species, properties of admixing groups, identification of signatures of natural selection, and may aid disease gene mapping. The algorithm employs nested hidden Markov models to obtain local ancestry estimation along the genome for each admixed individual. In a range of simulations, the accuracy of these estimates equals or exceeds leading existing methods that return local ancestry. Moreover, and unlike these approaches, we do not require any prior knowledge of the relationship between sub-groups of donor reference haplotypes and the unseen mixing ancestral populations. Instead, our approach infers these in terms of conditional "copying probabilities". In application to the Human Genome Diversity Panel we corroborate many previously inferred admixture events (e.g. an ancient admixture event in the Kalash). We further identify novel events such as complex 4-way admixture in San-Khomani individuals, and show that Eastern European populations possess 1 – 5% ancestry from a group resembling modern-day central Asians. We also identify evidence of recent natural selection favouring sub-Saharan ancestry at the HLA region, across North African individuals. We make available an R and C++ software library, which we term MOSAIC (which stands for MOSAIC Organises Segments of Ancestry In Chromosomes).

KEYWORDS population genetics; admixture; drift; selection; demography

Admixture Occurs when reproductive isolation between groups allows genetic divergence via genetic drift and random mutation, followed by mixing of the diverged groups to form new populations. Such genetic admixture is near ubiquitous in observed human populations (Hellenthal *et al.* 2014; Loh *et al.* 2013; Patterson *et al.* 2012) and indeed other species including cattle (Upadhyay *et al.* 2016), bison (Musani *et al.* 2006), and wolves (Pickrell and Pritchard 2012).

Genome-wide summaries can reveal not only complex relationships between modern populations but also details of their demographic histories (Pickrell and Pritchard 2012; Hellenthal *et al.* 2014; Peter 2016) while accurate inference of local ancestry can be used to correct for population structure in association testing (Diao and Chen 2012; Xu and Guan 2014), detect selection (Zhou *et al.* 2016), and can be used for mapping disease loci (Zhang and Stram 2014).

Due to the process of recombination, contiguous chunks of

admixed individuals' genomes are inherited intact from one mixing population or another. In the second generation following the initial admixture, chromosomes from distinct ancestral groups begin to recombine, and so the expected length of these chunks (in units of Morgans) will be one (by definition) and (neglecting crossover interference) chunk lengths can be modelled using an exponential distribution with rate parameter 1. In each subsequent generation, recombination further breaks down these chunks so that the chunk lengths (if they could be observed) are distributed according to an exponential distribution with rate parameter one less than the number of generations since admixture.

To fully characterize admixture for the above purposes, we need to infer: (1) Whether a group of individuals are admixed (2) The component / mixing groups (3) The timing of the admixture event(s) (4) Which segments of the admixed genome are inherited from each mixing group. Typically we lack prior knowledge of each of these points and we do not have access to representative samples of the mixing groups, as these are often no longer present (without drift) in modern samples.

A wide variety of approaches to model admixture have been

¹School of Mathematics and Statistics, University College Dublin, Belfield, Dublin 4, Ireland. michael.salter-townshend@ucd.ie

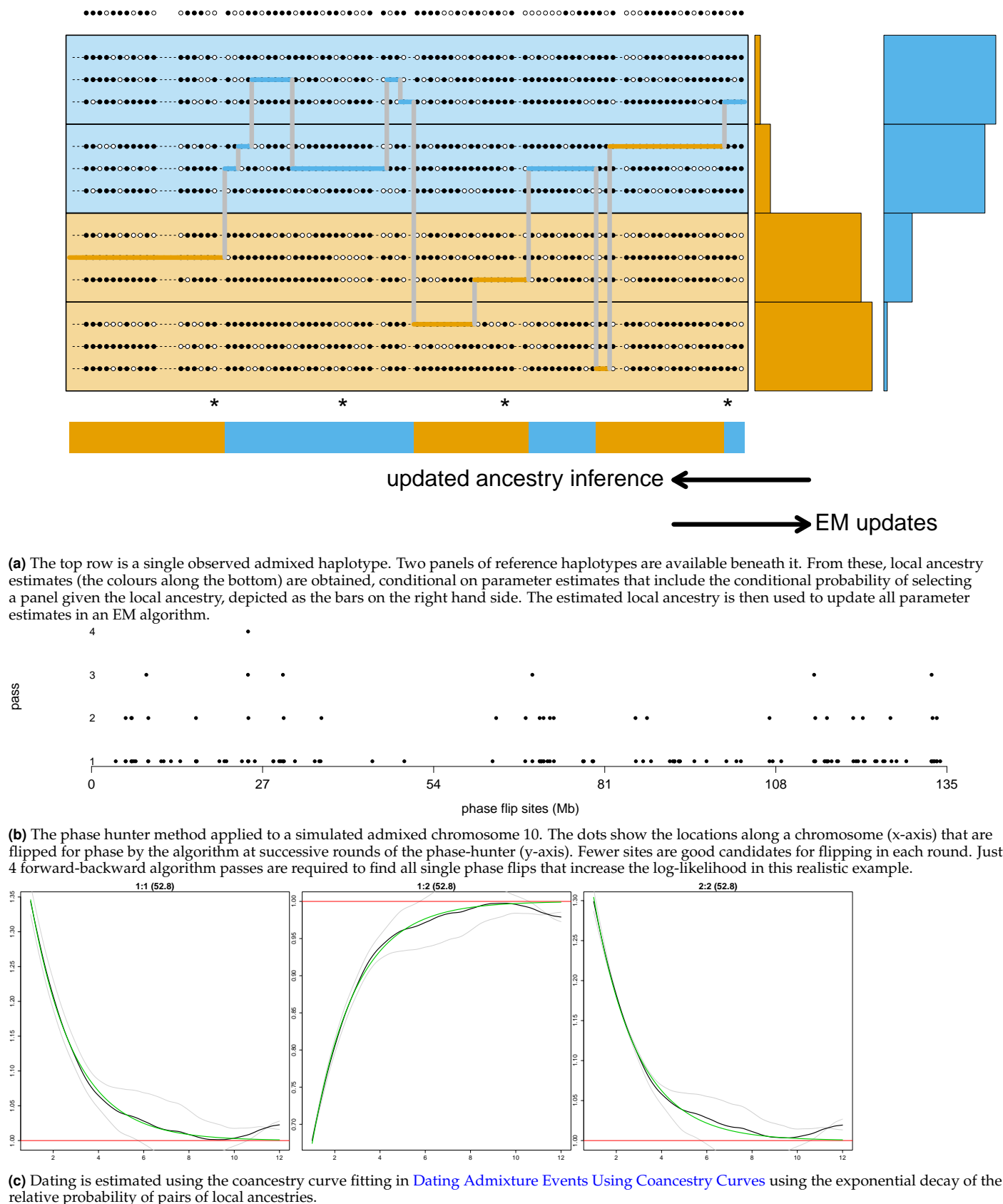


Figure 1 MOSAIC proceeds by rounds of thin (see [Thinning](#)), EM (see [EM Updates](#)), phasing (see [Re-Phasing](#)). **1a** is a cartoon version and **1b** and **1c** depict the realistic simulations used to test the approach in [Simulation Studies](#).

2012). We perform EM until convergence to estimate ρ, θ, \mathbf{M} , where the latter is simply a vector of copying probabilities for each group.

We now wish to initialise values for the ancestry aware part of the model i.e. an \mathbf{M} copying matrix with one column per hidden ancestry and the marginal probabilities of each ancestry α . We first impose windows of size $0.5cM$ across the genome and calculate the expected number of switches into each donor group in each window. This window size is chosen so that there is typically one latent ancestry per window but multiple donor groups are copied from (the number of generations undergoing recombination until we expect a single event in a window is 200). We then use EM to fit a mixture model where the number of mixtures is the number of hidden ancestries we wish to model. For haplotype h in window w , the expected number of switches $\hat{S}_{w}^{(h)}$ into the panels is modelled as a mixture of Multinomials.

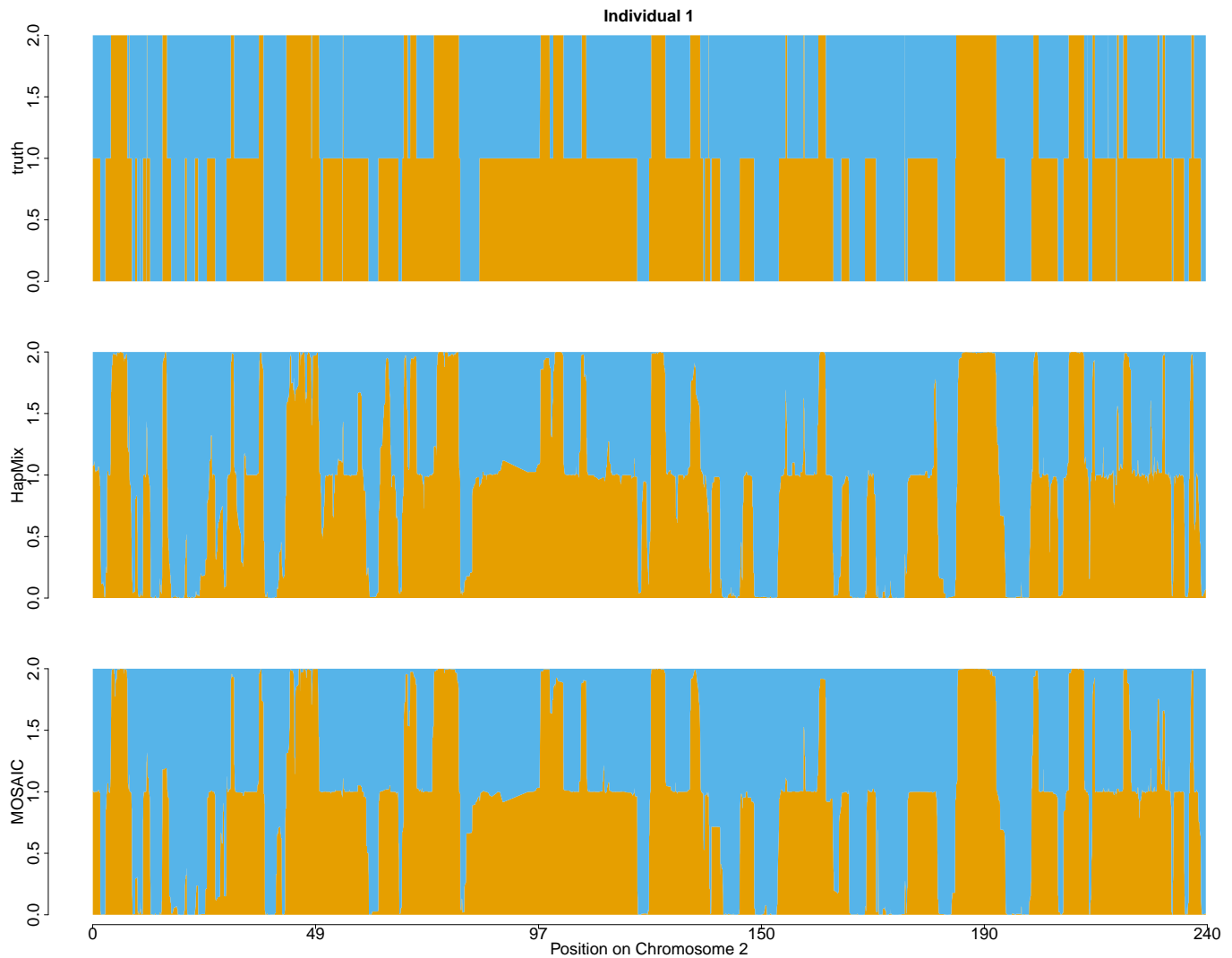
Appendix: Comparison with GLOBETROTTER for 2-way admixture events

As we regard GLOBETROTTER to be both the closest in spirit to our approach (in terms of genome-wide estimate) and the state-of-the-art in estimating time since admixture, we expand upon the comparison provided in Section [Two-way, Single Event](#). Table 4 compares results for the two-way admixture events inferred from the expanded HGDP dataset and comprises the values used to create Figure 3. The contribution of our method lies in accurate multiway local ancestry estimation within a single model and framework that provides these estimates.

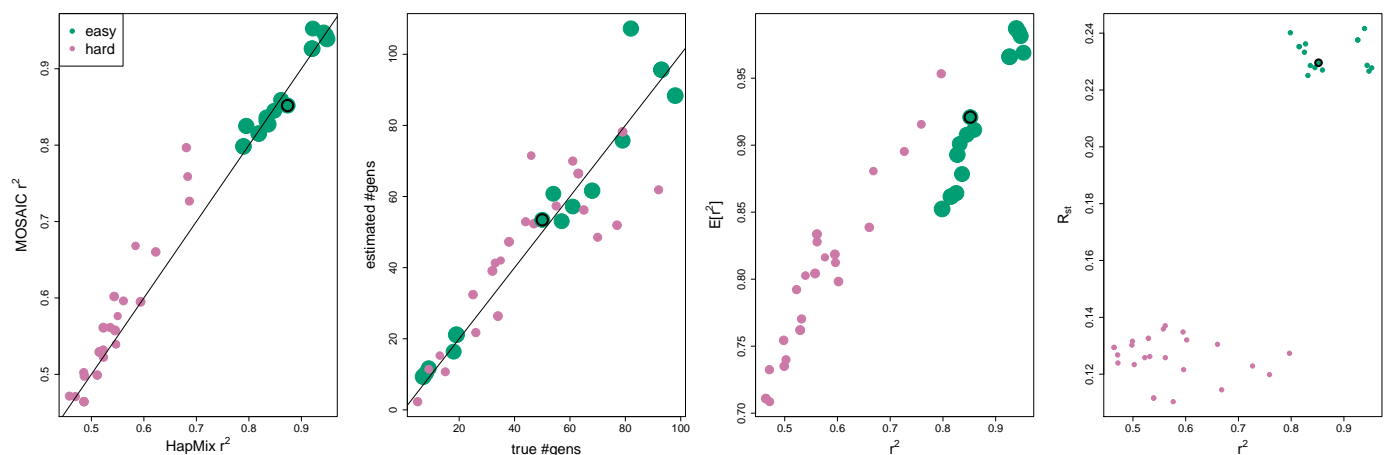
Appendix: Simulation of positive selection

We simulated positive selection in an admixed population for a single locus on Chromosome 6. We begin with ancestry proportions equal to those inferred in North Africa (minor ancestry of 15%, see Section [Selection Signal at the HLA in North Africa](#)). Using $N_e = 10,000$ diploid individuals, we simulated random recombinations along chromosomes of genetic length equal to Chromosome 6 (1.93 Morgans) for 31 generations. The number of recombinations is Poisson with rate 1.93 and the locations of the recombinations are uniform along genetic distance. We assume a Wright-Fisher with selection model of random mating amongst individuals and keep track of the simulated ancestry switch points continuously along each chromosome, as well as the ancestry of each segment. We set a non-zero selection coefficient at a single locus of $s = 0.035$ such that haplotypes containing ancestry of the minor type a at this locus are up-weighted with a relative weight of $1 + s$. i.e. when considering parents, each individual selected parents randomly with the probability of selecting a parent of ancestry a given by $\frac{(1+s)N_a}{2N_c+N_a}$, where N_a is the number of haplotypes of ancestry a at the selected locus. After 31 generations, we have 10^4 admixed individuals from which we sub-sample 220 diploid individuals (440 haplotypes). We then plot the average ancestry across these individuals against locus in Figure 9, noting the spike centred at the locus simulated to be under selection.

Importantly, when we then set real haplotypic data along the so-simulated ancestry segments but with $s = 0$ (no post-admixture selection effect) using donors from Southern Europe and sub-Saharan Africa and then run MOSAIC on this simulated data we do not infer a spike of African ancestry at the HLA.



(a) Example diploid local ancestry in simulated dataset. Top is truth, middle HapMix, bottom MOSAIC where re-phasing is as per [Re-Phasing](#). Both methods infer local ancestry that is highly consistent with the ground truth, however MOSAIC is more confident, extends to > 2-way events, and doesn't depend on prior knowledge of mixing groups.



(b) Easy (Yoruba and French) and harder (Pathan and Mongola) admixture simulations. **left:** r^2 for HapMix against MOSAIC, showing superior local ancestry estimation against the current state-of-the-art in 2-way admixture, even though HapMix is provided with known reference panels and MOSAIC is not. The plotting character is sized proportional to R_{st} . **left-centre:** true versus inferred dates. **right-centre:** estimated coefficient of determination of local ancestry $E[r^2]$ (Equation 2) versus squared correlation between true and estimated local ancestry. **right:** R_{st} (Equation 4) against squared correlation between true and estimated local ancestry, showing that R_{st} can be used to identify challenging cases. In all plots the black circle highlights the simulation shown in Figure 2a.

Figure 2 Comparison between MOSAIC and HapMix in 2-way admixture simulations, as described in [Simulation Studies](#).

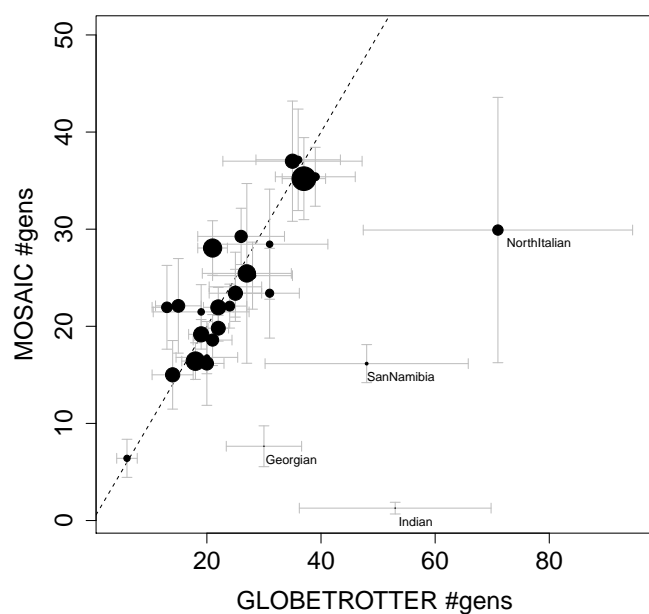


Figure 3 Inferred dates from MOSAIC are plotted against inferred dates from GLOBETROTTER, including bootstrapped ± 2 standard error bars for real data 2-way admixture events (see note in [Two-way, Single Event](#)). The GLOBETROTTER dates are from Table S14 of [Hellenthal *et al.* \(2014\)](#). The size of the central disc of each event is proportional to R_{st} (see [F_{st} Summaries](#)). Note that Melanesian is not shown as MOSAIC infers a very old mixing event of 221 generations ago.

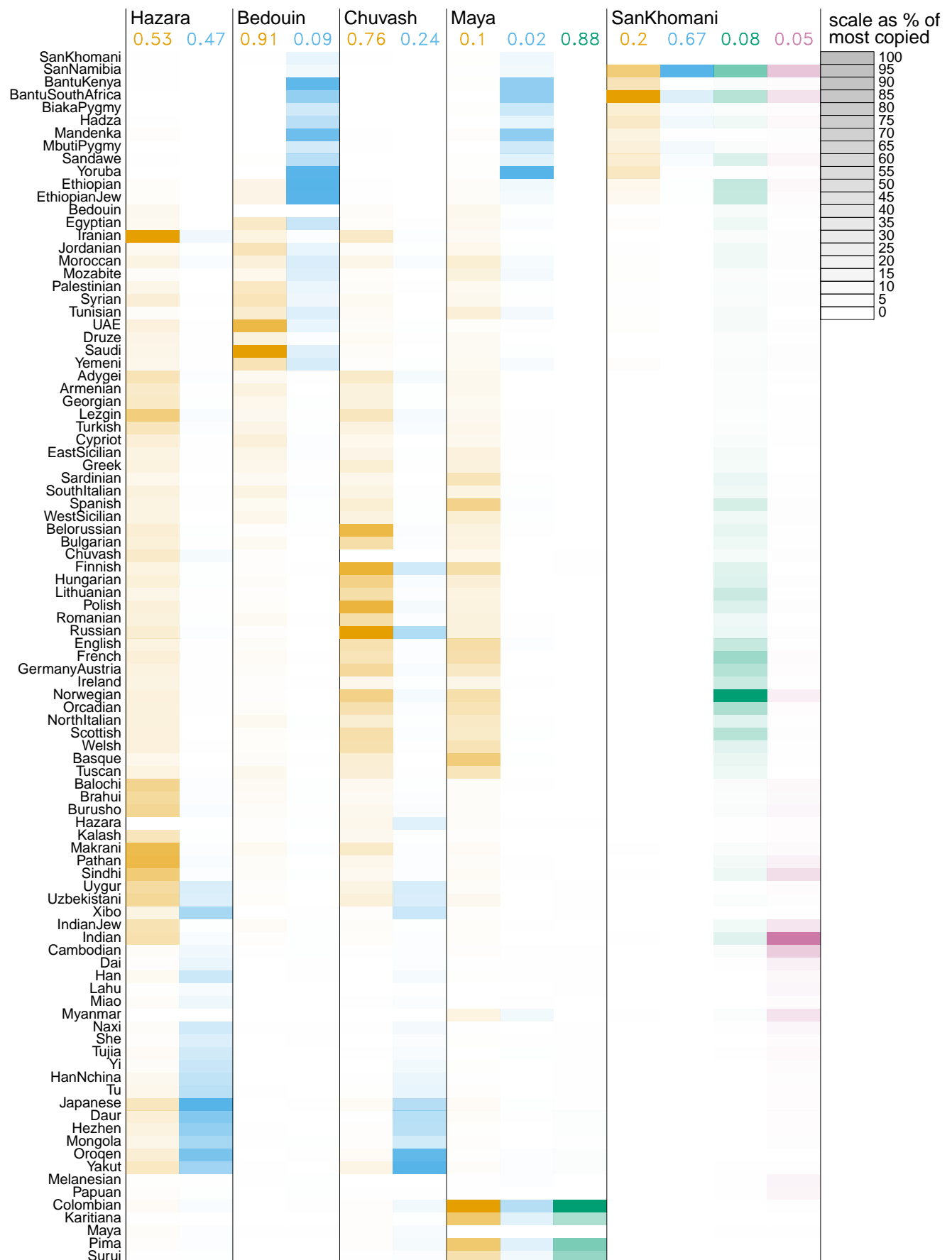


Figure 4 Inferred copying matrices for case studies of human admixture based on the HGDP dataset. The copying proportions μ_{pa} are scaled within columns to % of the most copied donor population so that each cell shading is equal to $100 \cdot \mu_{pa} / \arg \max_p \mu_{pa}$. Along the top are the marginal ancestry proportions for each admixed target.

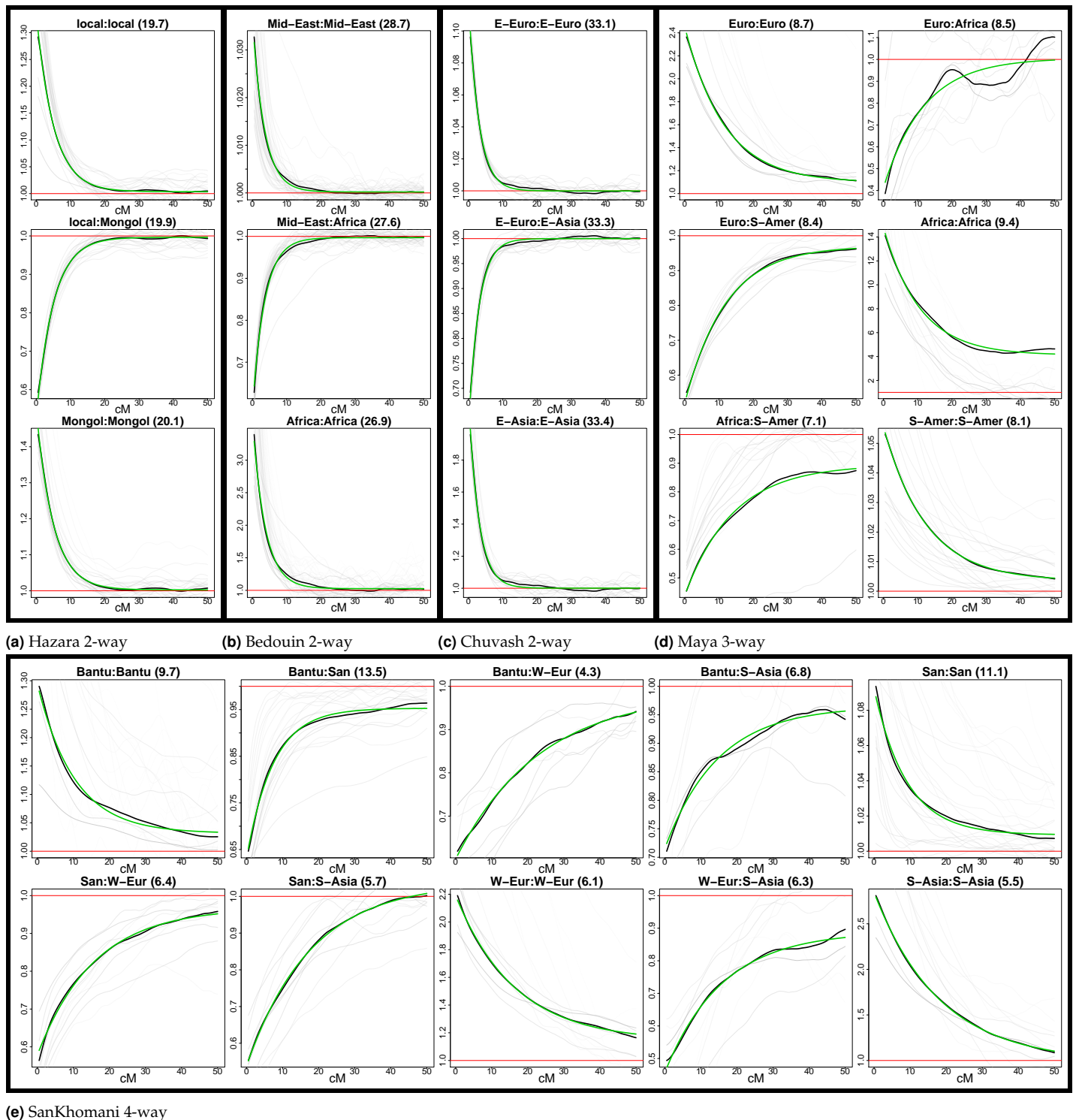


Figure 5 Coancestry curves for case studies of admixture within the HGDP dataset. On the top of each sub-plot the ancestry sides are labelled according to the closest donor panel as measured by \hat{F}_{st} (see Tables 1 to 3) and the estimated number of generations since admixture between each pair of ancestries is given in brackets.

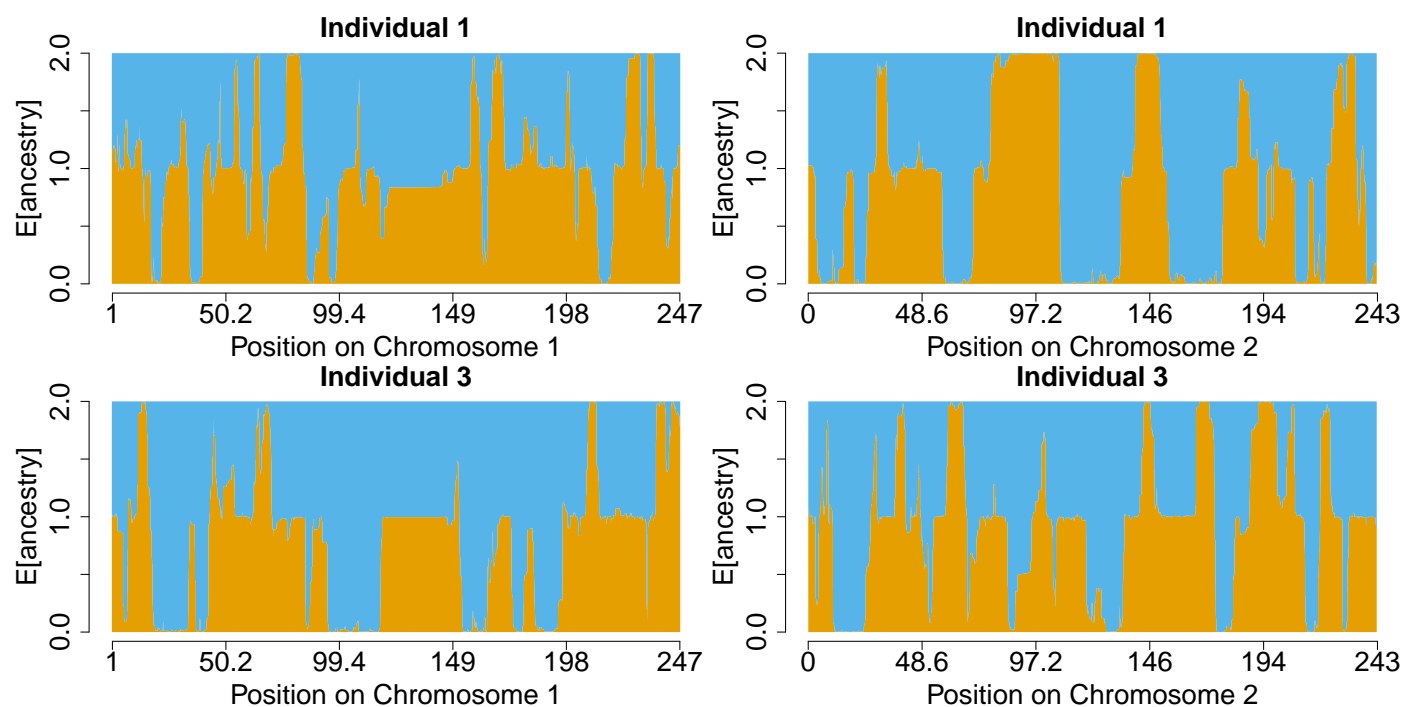
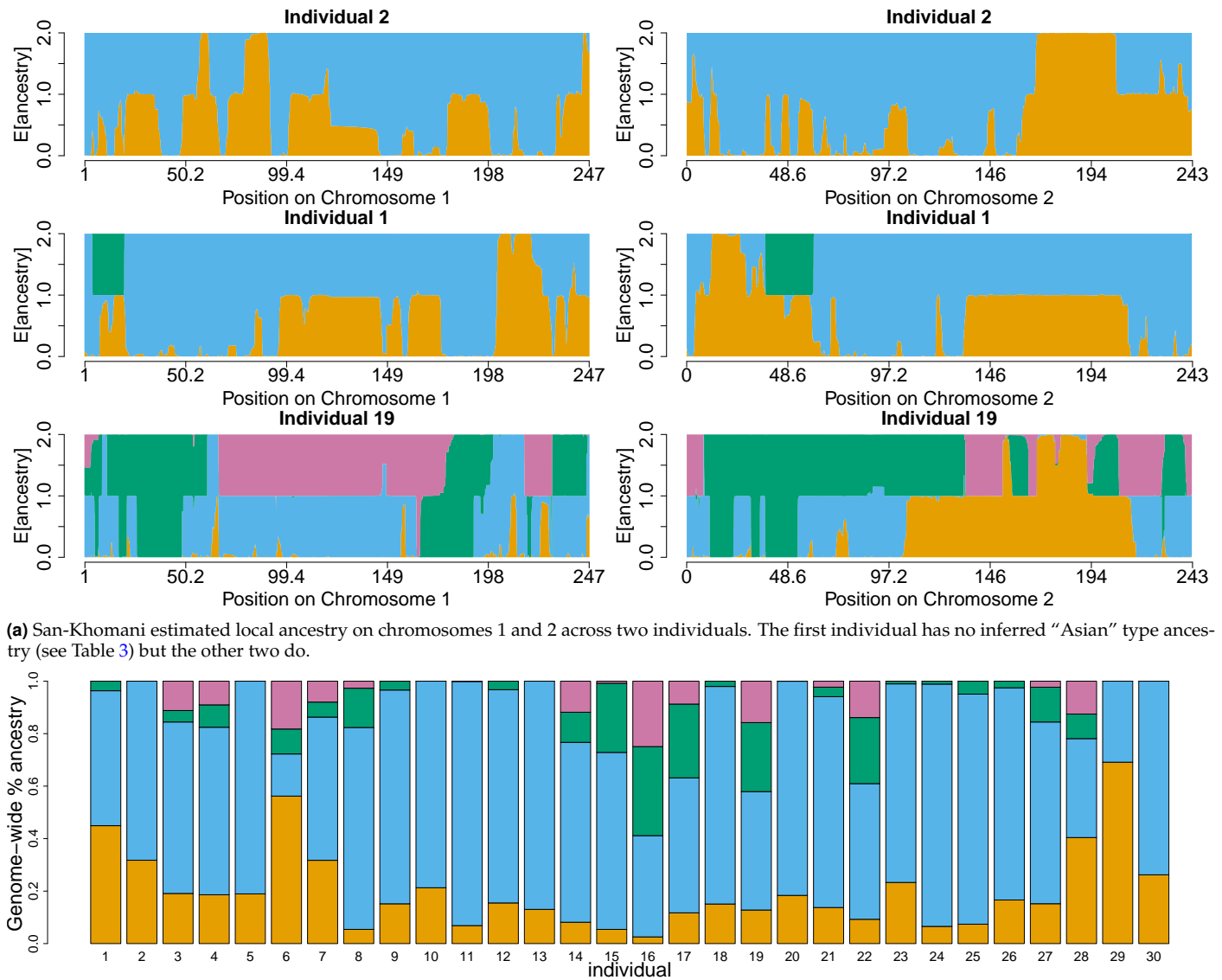
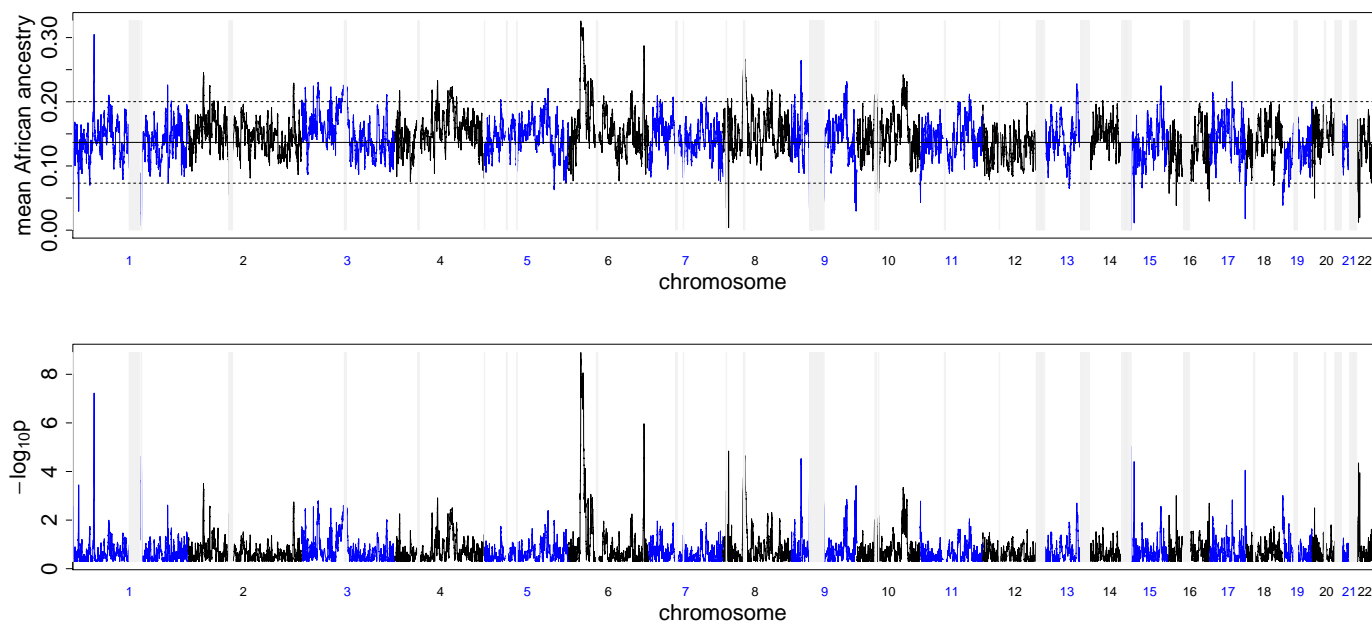


Figure 6 Hazara estimated local ancestry on chromosomes 1 and 2 across two individuals. There is a roughly 50-50 ancestry contribution from two sources approximately 20 generations ago. The orange source is Pathan like and the blue is Mongolian type (see Table 1 for details). The colours are consistent with Figure 4 which shows scaled copying proportions for each donor panel.

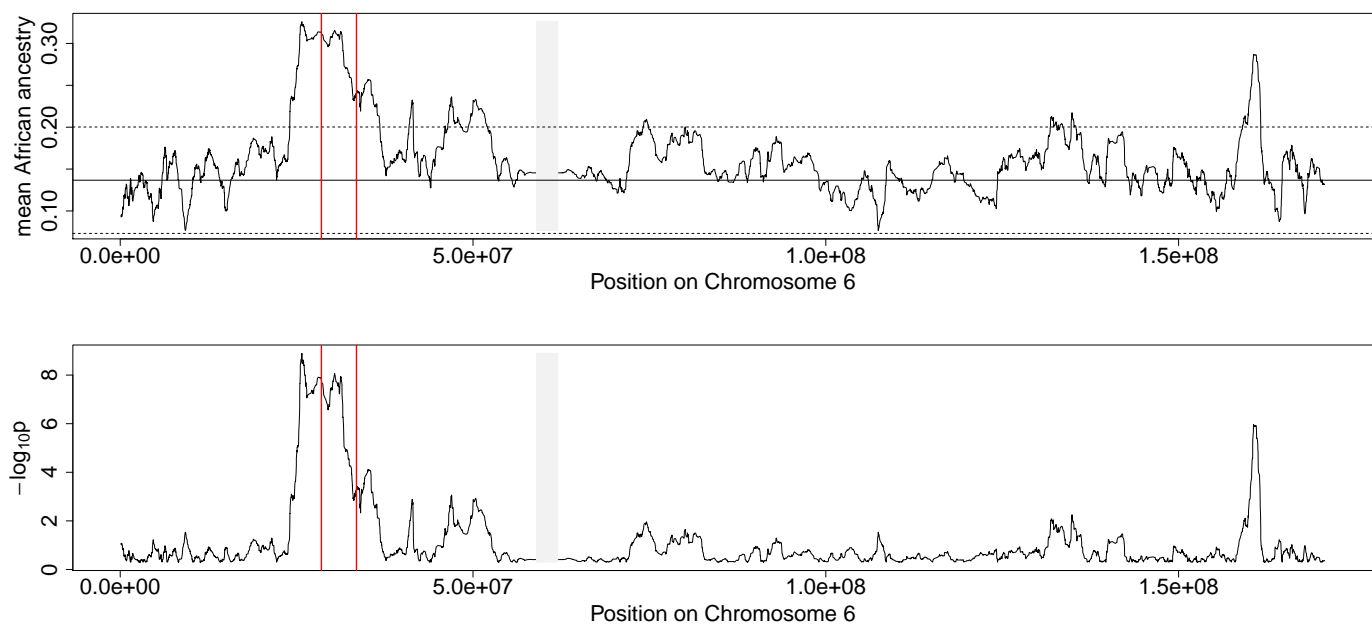


(b) Diagram showing the inferred proportions of the three ancestries in the San-Khomani.

Figure 7 Details of San-Khomani 4-way admixture model fit. The orange source is Bantu-like, blue is San, green is European, and purple is Asian (see Table 3 for details). The colours in both plots are consistent with Figure 4 which shows scaled copying proportions for each donor panel.



(a) Mean African Ancestry across all 220 individuals in North Africa against genome position.



(b) Mean African Ancestry across all 220 individuals in North Africa against Chromosome 6 position.

Figure 8 Mean African Ancestry. There is a high and wide spike at the HLA on Chromosome 6 at the HLA. Note that we have blocked out (in light grey) all 1Mb regions with fewer than 10 markers; this includes centromeres with low recombination rates and few SNPs.

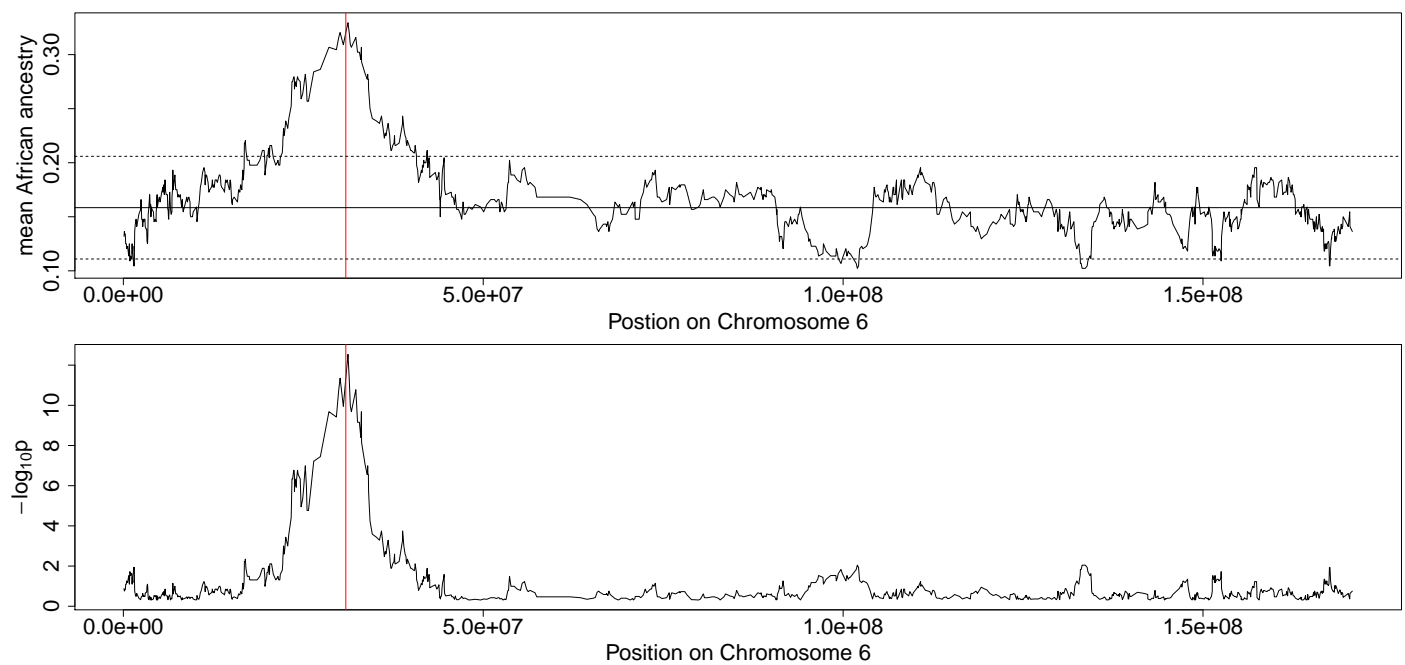


Figure 9 Mean Ancestry in a Wright-Fisher simulation on Chromosome 6 with positive selection at a single locus. There is a high and wide spike at the locus under selection. Note the width of the resulting spike due to hitchhiking of neighbouring loci. The solid vertical line is the mean ancestry outside the selected region and the dashed lines denote ± 2 standard deviations.