

ADMIXTURE INTO AND WITHIN SUB-SAHARAN AFRICA

GEORGE B.J. BUSBY^{1,†}, GAVIN BAND^{1,2}, QUANG SI LE¹, MUMINATOU JALLOW^{3,4}, EDITH BOUGAMA⁵, VALENTINA MANGANO⁶, LUCAS AMENGA-ETEGO⁷, ANTHONY ENIMIL⁸, TOBIAS APINJOH⁹, CAROLYNE NDILA¹⁰, ALPHAXARD MANJURANO^{11,12}, VYSAUL NYIRONGO¹³, OGOBARA DOUMBO¹⁴, KIRK, A. ROCKETT^{1,2}, DOMINIC P. KWIATKOWSKI^{1,2}, & CHRIS C.A. SPENCER^{1,†}
IN ASSOCIATION WITH THE MALARIA GENOMIC EPIDEMIOLOGY NETWORK¹⁵

¹Wellcome Trust Centre for Human Genetics, Roosevelt Drive, Oxford, OX3 7BN, UK

²Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge, CB10 1SA, UK

³Medical Research Council Unit, Serrekunda, The Gambia

⁴Royal Victoria Teaching Hospital, Banjul, The Gambia

⁵Centre National de Recherche et de Formation sur le Paludisme (CNRFP), Ouagadougou, Burkina Faso

⁶Dipartimento di Sanita Pubblica e Malattie Infettive, University of Rome La Sapienza, Rome, Italy

⁷Navrongo Health Research Centre, Navrongo, Ghana

⁸Komfo Anokye Teaching Hospital, Kumasi, Ghana

⁹Department of Biochemistry and Molecular Biology, University of Buea, Buea, Cameroon

¹⁰KEMRI-Wellcome Trust Research Programme, Kilifi, Kenya

¹¹Joint Malaria Programme, Kilimanjaro Christian Medical Centre, Moshi, Tanzania

¹²Faculty of Infectious and Tropical Diseases, London School of Hygiene and Tropical Medicine, London, UK

¹³Malawi-Liverpool Wellcome Trust Clinical Research Programme, College of Medicine, University of Malawi, Blantyre, Malawi

¹⁴Malaria Research and Training Centre, Faculty of Medicine, University of Bamako, Bamako, Mali

¹⁵MalariaGEN consortium; <http://www.malariagen.net/projects/host/consortium-members>

[†]corresponding authors: GBJB george@well.ox.ac.uk; CCAS spencer@well.ox.ac.uk

February 15, 2016

Abstract

Understanding patterns of genetic diversity is a crucial component of medical research in Africa. Here we use haplotype-based population genetics inference to describe gene-flow and admixture in a collection of 48 African groups with a focus on the major populations of the sub-Sahara. Our analysis presents a framework for interpreting haplotype diversity within and between population groups and provides a demographic foundation for genetic epidemiology in Africa. We show that coastal African populations have experienced an influx of Eurasian haplotypes as a series of admixture events over the last 7,000 years, and that Niger-Congo speaking groups from East and Southern Africa share ancestry with Central West Africans as a result of recent population expansions associated with the adoption of new agricultural technologies. We demonstrate that most sub-Saharan populations share ancestry with groups from outside of their current geographic region as a result of large-scale population movements over the last 4,000 years. Our in-depth analysis of admixture provides an insight into haplotype sharing across different geographic groups and the recent movement of alleles into new climatic and pathogenic environments, both of which will aid the interpretation of genetic studies of disease in sub-Saharan Africa.

Introduction

Genetic epidemiological studies aim to uncover novel relationships between genes, the environment, and disease [Malaria Genomic Epidemiology Network, 2015]. A critical component of this work is a description of patterns of genetic variation and an understanding of the genetic structure of populations. Whilst tens of thousands of genetic variants have been associated with different diseases in populations of European descent [Welter et al., 2014], despite the high burden of disease in Africa, medical genetic research on the continent has lagged behind [Need and Goldstein, 2009]. To alleviate this, several broad consortia are beginning to focus on understanding the genetic basis of infectious and non-communicable disease specifically in Africa [Malaria Genomic Epidemiology Network, 2008; H3Africa Consortium, 2014; Gurdasani et al., 2014; Malaria Genomic Epidemiology Network, 2015], and a number of recent studies have sought to describe patterns of genetic variation across the continent [Campbell and Tishkoff, 2008; Tishkoff et al., 2009; Gurdasani et al., 2014].

Genome-wide analyses of African populations are refining previous models of the continent's genetic history. One such emerging insight is the identification of clear, but complex, evidence for the movement of Eurasian ancestry back into the continent as a result of admixture over a variety of timescales [Pagani et al., 2012; Pickrell et al., 2014; Gurdasani et al., 2014; Hodgson et al., 2014a; Llorente et al., 2015]. Admixture occurs when genetically differentiated ancestral groups come together and mix, a process which is increasingly regarded as a common feature of human populations across the globe [Patterson et al., 2012; Hellenthal et al., 2014; Busby et al., 2015]. In a broad sample of 18 ethnic groups from eight countries, the African Genome Variation Project (AGVP) [Gurdasani et al., 2014] recreated previously published results to identify recent Eurasian admixture within the last 1.5 thousand years (ky) in the Fulani of West Africa [Tishkoff et al., 2009; Henn et al., 2012] and several East African groups from Kenya, older Eurasian ancestry (2-5 ky) in Ethiopian groups, consistent with previous studies of similar populations [Pagani et al., 2012; Pickrell et al., 2014], and a novel signal of ancient (>7.5 ky) Eurasian admixture in the Yoruba of Central West Africa [Gurdasani et al., 2014]. Comparisons of contemporary sub-Saharan African populations with the first ancient genome from within Africa, a 4.5 ky Ethiopian individual [Llorente et al., 2015], provide additional support for limited migration of Eurasian ancestry back into East Africa within the last 3,000 years.

Within this timescale, the major demographic change within Africa was the transition from hunting and gathering to pastoralist and agricultural lifestyles [Diamond and Bellwood, 2003; Smith, 2005; Barham and Mitchell, 2008; Li et al., 2014]. This shift was long and complex and occurred at different speeds, instigating contrasting interactions between the agriculturalist pioneers and the inhabitant people [Mitchell, 2002; Marks et al., 2014]. The change was initialised by the spread of pastoralism (i.e. the raising and herding of livestock) across Africa and the subsequent movement east and south from Central West Africa of agricultural technology together with the branch of Niger-Congo languages known as Bantu [Mitchell, 2002; Barham and Mitchell, 2008]. The AGVP also found evidence of widespread hunter-gatherer ancestry in African populations, including ancient (9 ky) Khoesan ancestry in the Igbo from Nigeria, and more recent hunter-gatherer ancestry in eastern (2.5-4.5 ky) and southern (0.9-4 ky) African populations [Gurdasani et al., 2014]. The identification of hunter-gatherer ancestry in non-

hunter-gatherer populations together with the timing of these latter events is consistent with the known expansion of Bantu languages across Africa within the last 3 ky [Mitchell, 2002; Diamond and Bellwood, 2003; Smith, 2005; Barham and Mitchell, 2008; Marks et al., 2014; Li et al., 2014].

These studies have described the novel and important influence of both Eurasian and hunter-gatherer ancestry on the population genetic history of sub-Saharan Africa. Nevertheless, there is still more to understand about the nature and timing of admixture across Africa. For example, do we observe Eurasian ancestry in additional African populations? Do alternative dating methods recreate a similar timescale for these admixtures? Can we understand more about how the Eurasian ancestry entered African populations? Uncertainty about the relationships between populations within Africa also remains. For example, can we use genetics to refine the origins, interactions, and timings of the Bantu expansion? What did the ancestry of pre-Bantu populations of east and southern Africa look like? What other historical relationships can be characterised between different African ethno-linguistic groups?

Here we review these questions using the same methods as previous authors, and provide additional novel inference by using methods that utilise haplotype information. To gain a detailed understanding of the population structure and history of Africa, we analyse genome-wide data from 14 Eurasian and 46 sub-Saharan African groups. Half (23) of the African groups represent subsets of samples collected from nine countries as part of the MalariaGEN consortium. Details on the recruitment of samples in relation to studying malaria genetics are published elsewhere [Malaria Genomic Epidemiology Network, 2014, 2015]. The remaining 23 groups are from publicly available datasets from a further eight sub-Saharan African countries [Pagani et al., 2012; Schlebusch et al., 2012; Petersen et al., 2013] and the 1000 Genomes Project (1KGP), with Eurasian groups from the latter included to help understand the genetic contribution from outside of the continent (Figure 1-figure supplement 1). Although we have representative groups from all four major African linguistic macro-families (Supplementary File 1), our sample represents a significant proportion of the sub-Saharan population in terms of number, but not does not equate to a complete picture of African ethnic diversity.

To study patterns of genetic diversity we created an integrated dataset at a common set of over 328,000 high-quality SNP genotypes and use established approaches for comparing population allele frequencies across groups to provide a baseline view of historical gene flow. We apply statistical approaches to phasing genotypes to obtain haplotypes for each individual, and use previously published methods to represent the haplotypes that an individual carries as a mosaic of other haplotypes in the sample (so-called chromosome painting [Li and Stephens, 2003]). We use these data to demonstrate that haplotype-based methods have the potential to tease apart subtle relationships between closely related populations. We present a detailed picture of haplotype sharing across sub-Saharan Africa using a model-based clustering approach that groups individuals using haplotype information alone. The inferred groups reflect broad-scale geographic patterns. At finer scales, our analysis reveals smaller groups, and often differentiates closely related populations consistent with self-reported ancestry [Tishkoff et al., 2009; Bryc et al., 2010; Hodgson et al., 2014a]. We describe these patterns by measuring gene flow between populations and relate them to potential historical movements of people into and within sub-Saharan Africa. Understanding the extent to which individuals share haplotypes (which we call *coancestry*), rather than independent markers, can provide a rich description of ancestral relationships and population history [Lawson et al.,

2012; Leslie et al., 2015]. For each group we use the latest analytical tools to characterise the populations as mixtures of haplotypes and provide estimates for the date of admixture events [Lawson et al., 2012; Hellenthal et al., 2014; Leslie et al., 2015; Montinaro et al., 2015]. As well as providing a quantitative measure of the coancestry between groups, we identify the detailed dominant events which have shaped current genetic diversity in sub-Saharan Africa. We discuss the relevance of these observations to studying genotype-phenotype associations in Africa, particularly in the context of infectious disease.

Results

Broad-scale population structure reflects geography and language

Throughout this article we use shorthand current-day geographical and ethno-linguistic labels to describe ancestry. For example we write “Eurasian ancestry in East African Niger-Congo speakers”, where the more precise definition would be “ancestry originating from groups currently living in Eurasia in groups currently living in East Africa that speak Niger-Congo languages” [Pickrell et al., 2014]. We also stress that the use of Khoesan in the current setting refers to groups with shared linguistic characteristics which does not necessarily imply shared close genealogical relationships [Güldemann and Fehn, 2014]. Our combined dataset included 3,283 individuals from 46 sub-Saharan African ethnic groups and 12 non-African populations (Figure 1A and Figure 1-figure supplement 1).

As in other regions of the world [Novembre et al., 2008; Reich et al., 2009; Behar et al., 2010], analyses of population structure using principal component analysis (PCA) show that genetic relationships are broadly defined by geographical and ethno-linguistic similarity (Figure 1B,C). The first two principal components (PCs) reflect ethno-linguistic divides: PC1 splits southern Khoesan speaking populations from the rest of Africa, and PC2 splits the Afroasiatic and Nilo-Saharan speakers from East Africa from sub-Saharan African Niger-Congo speakers. The third axis of variation defines east versus west Africa, suggesting that in general, population structure in Africa largely mirrors linguistic and geographic similarity [Tishkoff et al., 2009].

Using a previously published implementation of chromosome painting (CHROMOPAINTER [Lawson et al., 2012]), we reconstructed the genomes of each study individual as mosaics of haplotype segments (or chunks) from all other individuals. We used the clustering algorithm implemented in fineSTRUCTURE [Lawson et al., 2012] to cluster individuals into groups based purely on the similarity of these paintings (Figure 1 and Figure 1-figure supplement 3). More specifically, for a given recipient individual, we can estimate the amount of their genome that is shared with (or copied from) each of a set of donor individuals based on the paintings, which are summarised as *copying vectors*. These vectors are clustered hierarchically to form a tree which describes the inferred relationship between different groups (Figure 1-figure supplement 3). As such, this method uses chromosome painting to describe the coancestry between groups, which can then be visualised as a heatmap, as in Figure 1D.

Consistent with PCA, African populations tend to share more DNA with geographically proximate populations (dark colours on the diagonal; Figure 1D). Block structures on this diagonal indicate higher levels of haplotype sharing within groups, which is indicative of close genealogical relationships with other

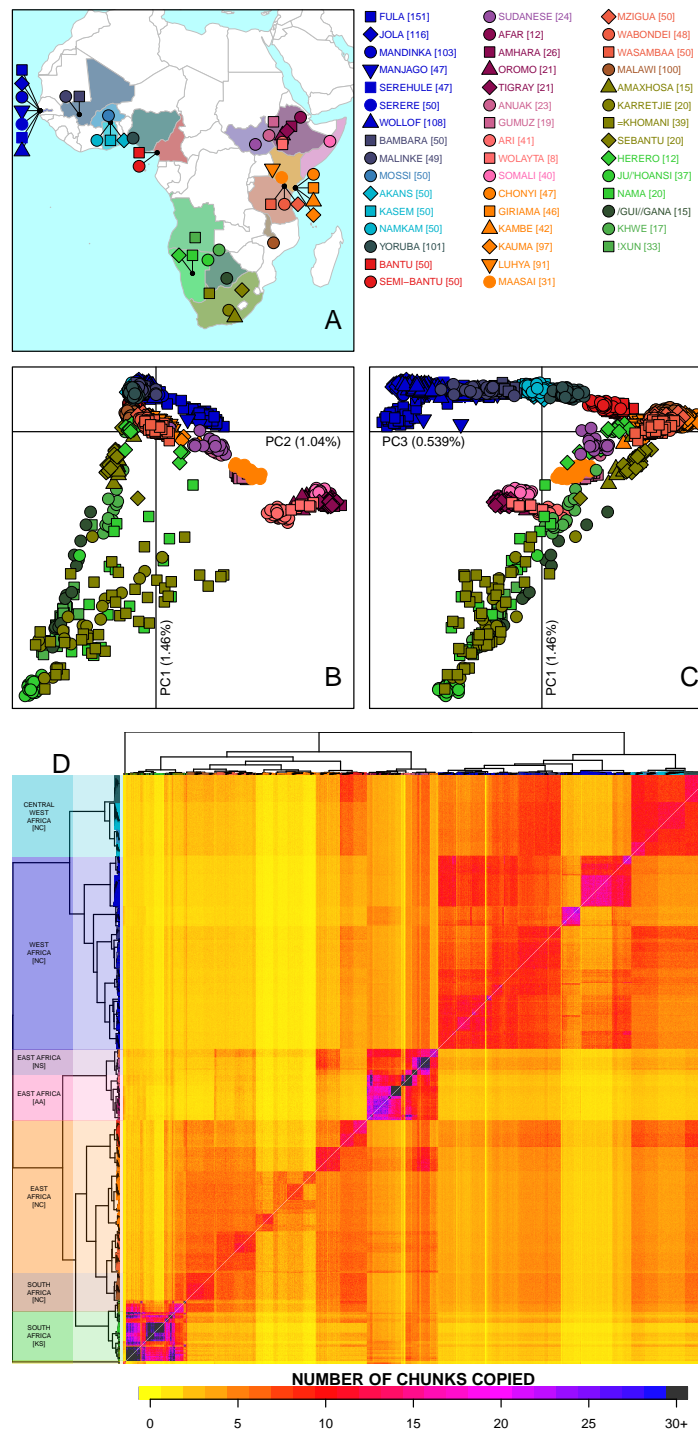


Figure 1. Sub-Saharan African genetic variation is shaped by ethno-linguistic and geographical similarity. (A) the origin of the 46 African ethnic groups used in the analysis; ethnic groups from the same country are given the same colour, but different shapes; the legend describes the identity of each point. Figure 1-figure supplement 1 and Figure 1-Source Data 1 provide further detail on the provenance of these samples. (B) PC1vPC2 shows that the first major axis of variation in Africa splits southern groups from the rest of Africa, each symbol represents an individual; PC2 reflects ethno-linguistic differences, with Niger-Congo speakers split from Afroasiatic and Nilo-Saharan speakers. (C) PC1vPC3 shows that the third principle component represents geographical separation of Niger-Congo speakers, forming a cline from west to east Africans. (D) results of the fineSTRUCTURE clustering analysis using copying vectors generated from chromosome painting; each row of the heatmap is a recipient copying vector showing the number of chunks shared between the recipient and every individual as a donor (columns); the tree clusters individuals with similar copying vectors together, such that block-like patterns are observed on the heatmap; darker colours on the heatmap represent more haplotype sharing (see text for details); individual tips of the tree are coloured by country of origin, and the seven ancestry regions are identified and labelled to the left of the tree; labels in parentheses describe the major linguistic type of the ethnic groups within: AA = Afroasiatic, KS = Khoesan, NC = Niger-Congo, NS = Nilo-Saharan.

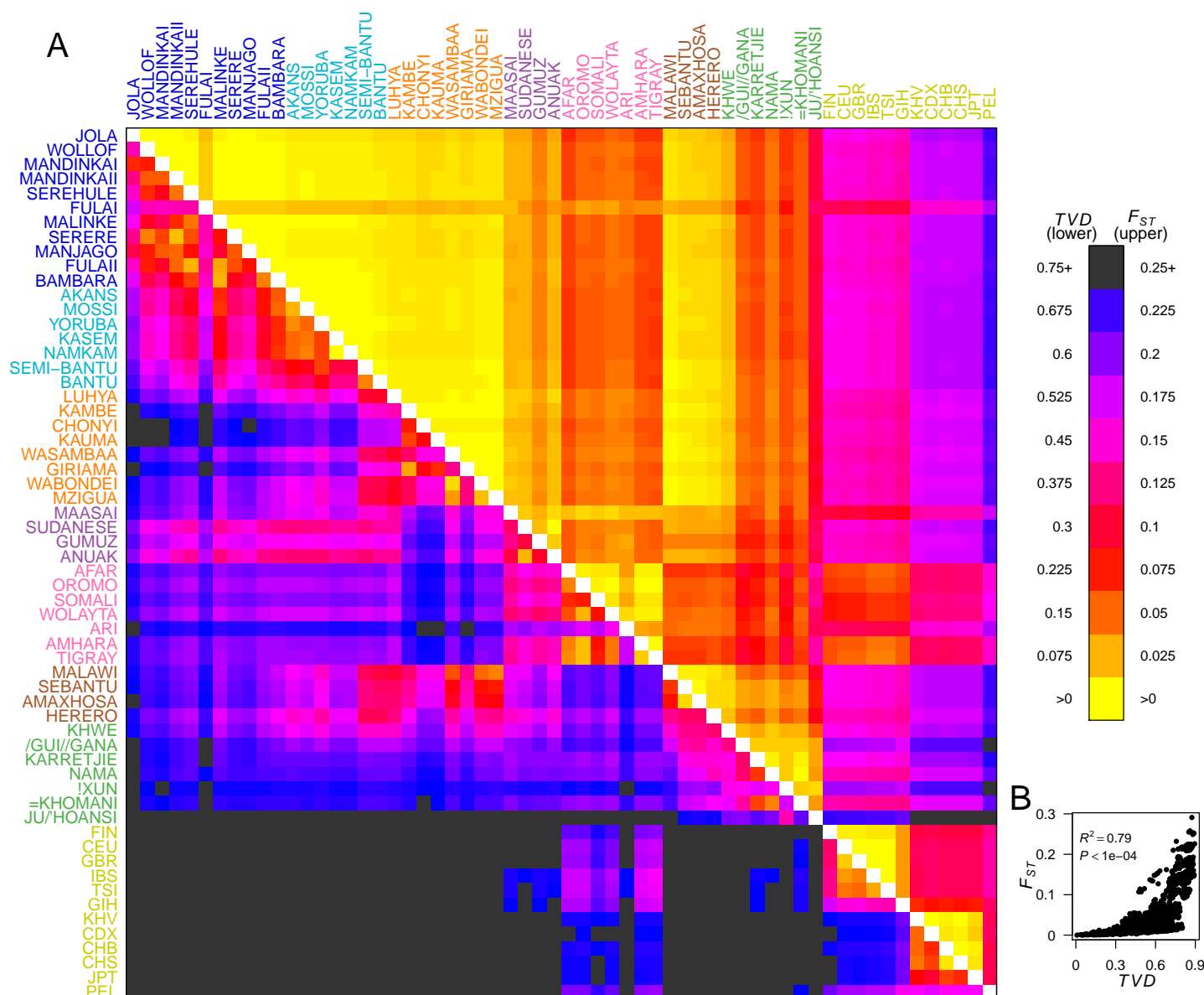
members of your group. These patterns can be seen in some of the Khoesan speaking individuals (eg. the Ju/'hoansi), several groups from the East Africa (Sudanese, Ari, and Somali groups), and the Fulani and Jola from the Gambia, each of which provides clues to the ancestral connections between the groups in our analysis. The heatmap also shows evidence for coancestry across regions (dark colours away from the diagonal), which can be informative of historical connections between modern-day groups. For example, east Africans from Kenya, Malawi and Tanzania tend to share more DNA with west Africans (lower right) which suggests that more haplotypes may have spread from west to east Africa (upper left) than vice versa. Using the results of the PCA and fineSTRUCTURE analyses together with ethno-linguistic classifications and geography, we defined seven groups of populations within Africa (Supplementary File 1), which we refer to as *ancestry regions* when describing gene-flow across Africa. Under this definition, there is widespread sharing of haplotypes within and between ancestry regions (Figure 1).

Haplotypes reveal subtle population structure

To quantify the extent of the genetic difference between groups we used two different metrics. First, we used the classical measure F_{ST} [Hudson et al., 1992; Bhatia et al., 2013] which measures the differentiation in allele frequencies between populations. It can be thought of as measuring the proportion of the heterozygosity at SNPs explained by the group labels. The second metric uses the similarity in copying patterns between two groups to estimate the total variation difference (TVD) at the haplotypic level. TVD takes advantage of the fact that recombination rates are faster than mutation rates, and so is expected to capture more variability than genotype counts. Figure 2A shows these two metrics side by side in the upper and lower diagonal. When compared to non-African populations, F_{ST} measured at our integrated set of SNPs is relatively low between many groups from West, Central, and East Africa (yellows on the upper right triangle), whereas TVD in the same populations can reveal haplotypic differences as strong as between Europe and Asia (pink and purples in lower left triangle). For example, the Chonyi from Kenya have relatively low F_{ST} but high TVD with West African groups, like the Jola (Chonyi-Jola $F_{ST} = 0.019$; Chonyi-Jola TVD = 0.803) suggesting that, whilst allele frequency differences between the two populations are relatively low, when we compare the populations' ancestry vectors, the haplotypic differences are some the strongest between sub-Saharan groups. In fact, whilst pairwise TVD tends to increase with pairwise F_{ST} (Pearson's correlation $R^2 = 0.79$) the relationship is not linear (Figure 2B) demonstrating that this discrepancy is not simply a function of TVD generating larger statistics. We note that high values of TVD are not sufficient to infer specific demographic events [van Dorp et al., 2015] and we therefore do not interpret these differences as necessarily resulting from admixture, but as motivation to use haplotype-based methods to characterise population relationships.

Widespread evidence for admixture

To check that our dataset was consistent with those used in other analyses of African variation, we next applied similar approaches to previous authors to infer admixture between populations using correlations in allele frequencies within and between populations [Pickrell et al., 2014; Gurdasani et al., 2014]. The first approach, the three-population test (f_3 statistic [Reich et al., 2009]), uses correlations in allele frequencies



between a target population and two potential source populations to identify significant departures from the null model of no admixture. Negative values are indicative of canonical admixture events where the allele frequencies in the target population are intermediate between the two source populations.

Consistent with recent research [Pickrell et al., 2014; Pickrell and Reich, 2014; Gurdasani et al., 2014; Llorente et al., 2015], the majority (80%, 40/48), but not all, of the African groups surveyed showed evidence of admixture ($f_3 < -5$), with the most negative f_3 statistic tending to involve either a Eurasian or African hunter-gatherer group [Gurdasani et al., 2014] (Supplementary Table 1).

The second approach, ALDER [Loh et al., 2013; Pickrell et al., 2014] (Supplementary Table 1) exploits the fact that correlations between allele frequencies along the genome decay over time as a result of recombination. Linkage disequilibrium (LD) can be generated by admixture events, and leaves detectable signals in the genome that can be used to infer historical processes from allele frequency data [Loh et al., 2013]. Following Pickrell et al. [2014] and the AGVP [Gurdasani et al., 2014], we computed weighted admixture LD curves using the ALDER [Loh et al., 2013] package to characterise the sources and timing of gene flow events. Specifically, we estimated the y-axis intercept (amplitude) of weighted LD curves for each target population using curves from an analysis where one of the sources was the target population (self reference) and the other was, separately, each of the other (non-self reference) populations. Theory predicts that the amplitude of these “one-reference” curves becomes larger the more similar the non-self reference population is to the true admixing source [Loh et al., 2013]. As with the f_3 analysis outlined above, for many of the sub-Saharan African populations, Eurasian and hunter-gatherer groups produced the largest amplitudes (Figure 3-figure supplement 1 and Figure 3-figure supplement 2), reinforcing the contribution of these ancestries to our broad set of African populations.

We investigated the evidence for more complex admixture using MALDER [Pickrell et al., 2014], an implementation of ALDER which fits a mixture of exponentials to weighted LD curves to infer multiple admixture events (Figure 3 and Table Figure 3-Source Data 1). In Figure 3A, for each target population, we show the ancestry region of the two populations involved in generating the MALDER curves with the greatest amplitudes, together with the date of admixture for at most two events. Throughout, we convert time since admixture in generations to a date by assuming a generation time of 29 years [Fenner, 2005]. In general, we find that groups from similar ancestry regions tend to have inferred events at similar times and between similar groups (Figure 3), which suggests that genetic variation in groups has been shaped by shared historical events. For every event, the curves with the greatest amplitudes involved a population from a (normally non-Khoesan) African population on the one side, and either a Eurasia or Khoesan population on the other. To provide more detail on the composition of the admixture sources, we compared MALDER curve amplitudes between curves involving populations from different ancestry regions (central panel Figure 3). In general, this analysis showed that, with a few exceptions, we were unable to precisely define the ancestry of the African source of admixture, as curves involving populations from multiple different regions were not statistically different from each other ($Z < 2$; SOURCE 1). Conversely, comparisons of MALDER curves when the second source of admixture was Eurasian (yellow) or Khoesan (green), showed that these groups were usually the single best surrogate for the second source of admixture (SOURCE 2).

We performed multiple MALDER analyses, varying the input parameters and the genetic map used. We observed a large amount of shared LD at short genetic distances between different African populations (Figure 3-figure supplement 3 and Figure 3-figure supplement 4). Such patterns may result from population genetic processes other than admixture, such as shared demographic history and population

bottlenecks [Loh et al., 2013]. In the main MALDER analysis we present, we have removed the effect of this short-range LD by generating curves only after ignoring SNPs at short genetic distances where frequencies of markers at such distances are correlated between target and reference population – which can be thought of as conservative – but provide supplementary analyses where this parameter was relaxed. The main difference between the two types of analysis is that we were only able to identify ancient events in West Africa if we force the MALDER algorithm to start computing LD decay curves from a genetic distance of 0.5cM, irrespective of any short-range correlations in LD between populations.

Many West African groups show evidence of recent (within the last 4 ky) admixture involving African and Eurasian sources. The Mossi from Burkina Faso have the oldest inferred date of admixture, at roughly 5000BCE, but we were unable to recreate previously reported ancient admixture events in other Central West African groups [Gurdasani et al., 2014] (although see Figure 3-figure supplement 5 for results where we use different MALDER parameters). Across East Africa Niger-Congo speakers (orange) we infer admixture within the last 4 ky (and often within the last 1 ky) involving Eurasian sources on the one hand, and African sources containing ancestry from other Niger-Congo speaking African groups from the west on the other. Despite events between African and Eurasian sources appearing older in the Nilo-Saharan and Afroasiatic speakers from East Africa, we see a similar signal of very recent Central West African ancestry in a number of Khoesan groups from Southern Africa, such as the Khwe, /Gui //Gana, and !Xun, together with Malawi-like (brown) sources of ancestry in recent admixture events in East African Niger-Congo speakers.

Inference of older events relies on modelling the decay of LD over short genetic distances because recombination has had more time to break down correlations in allele frequencies between neighbouring SNPs. We investigated the effect of using European (CEU) and Central West African (YRI) specific recombination maps [Hinch et al., 2011] on the dating inference. Whilst dates inferred using the CEU map were consistent with those using the HAPMAP recombination map (Figure 3B), when using the African map dates were consistently older (Figure 3C), although still generally still within the last 7ky. There was also variability in the number of inferred admixture events for some populations between the different map analyses (Figure 3-figure supplement 6 and Figure 3-figure supplement 7).

Most events involved sources where Eurasian (dark yellow in Figure 3A) groups gave the largest amplitudes. In considering this observation, it is important to note that the amplitude of LD curves will partly be determined by the extent to which a reference population has differentiated from the target. Due to the genetic drift associated with the out-of-Africa bottleneck and subsequent expansion, Eurasian groups will tend to generate the largest curve amplitudes even if the proportion of this ancestry in the true admixing source is small [Pickrell et al., 2014] (in our dataset, the mean pairwise F_{ST} between Eurasian and African populations is 0.157; Figure 2A and Figure 2-Source Data 1). To some extent this also applies to Khoesan groups (green in Figure 3A), who are also relatively differentiated from other African groups (mean pairwise F_{ST} between Ju/'hoansi and all other African populations in our dataset is 0.095; Figure 2A and Figure 2-Source Data 1). In light of this, and the observation that curves involving groups from different ancestry regions are often no different from each other, it is therefore difficult to infer the proportion or nature of the African, Khoesan, or Eurasian admixing sources, only that the sources themselves contained African, Khoesan, or Eurasian ancestry. Moreover, given uncertainty in the dating

of admixture when using different maps and MALDER parameters, these results should be taken as a guide to the general genealogical relationships between African groups, rather than a precise description of the gene-flow events that have shaped Africa.

Modelling gene flow with haplotypes

So far, our analyses have largely recapitulated recent studies of the genetic history of African populations, albeit across a broader set of sub-Saharan populations. However, we were interested in gaining a more detailed characterisation of historical gene-flow to provide an understanding and framework with which to inform genetic epidemiological studies in Africa. The chromosome painting methodology described above provides an alternative approach to inferring admixture events which directly models the similarity in haplotypes between pairs of individuals. Evidence of recent haplotype sharing suggests that the ancestors of two individuals must have been geographically proximal at some point in the past. More generally, the distance over which haplotype sharing extends along chromosomes is inverse to how far in the past coancestry events have occurred. We can use copying vectors inferred through chromosome painting to help identify those populations that share ancestry with a recipient group by fitting each vector as a mixture of all other population vectors (Figure 4A) [Leslie et al., 2015; Montinaro et al., 2015; van Dorp et al., 2015]. Figure 4A shows the contribution that each ancestry region makes to these mixtures (MIXTURE MODEL column). Almost all groups can best be described as mixtures of ancestry from different regions. For example, the copying vector of the Bantu ethnic group from Cameroon is best described as a combination of 40% Central West African Niger-Congo (sky blue), 30% Eastern Niger-Congo (orange), 25% Southern Niger-Congo (brown), and the remaining 5% coming from West African Niger-Congo (dark blue) and Khoesan-speaking (green) groups. The key insight from this analysis is that many African groups share fragments of haplotypes with groups from outside of their own ancestry region.

The mixture model approach is useful for describing and summarising the copying vectors of populations, but such summaries can result from both admixture and shared evolutionary history. We thus used GLOBETROTTER [Hellenthal et al., 2014], a method that explicitly tests for and characterises admixture as an extension of the mixture model approach described above, to gain a more detailed understanding of recent ancestry. Admixture inference can be challenging for a number of reasons: the true admixing source population is often not well represented by a single sampled population; admixture could have occurred in several bursts, or over a sustained period of time; and multiple groups may have come together as complex convolution of admixture events. GLOBETROTTER aims to overcome some of these challenges, in part by using painted chromosomes to explicitly model the correlation structure among nearby SNPs, but also by allowing the sources of admixture themselves to be mixed [Hellenthal et al., 2014]. In addition, the approach has been shown to be relatively insensitive to the genetic map used [Hellenthal et al., 2014], and therefore potentially provides a more robust inference of admixture events, the ancestries involved, and their dates. GLOBETROTTER uses the distance between chromosomal chunks of the same ancestry to infer the time since historical admixture has occurred.

Throughout the following discussion, we refer to target populations as *recipients*, any other sampled

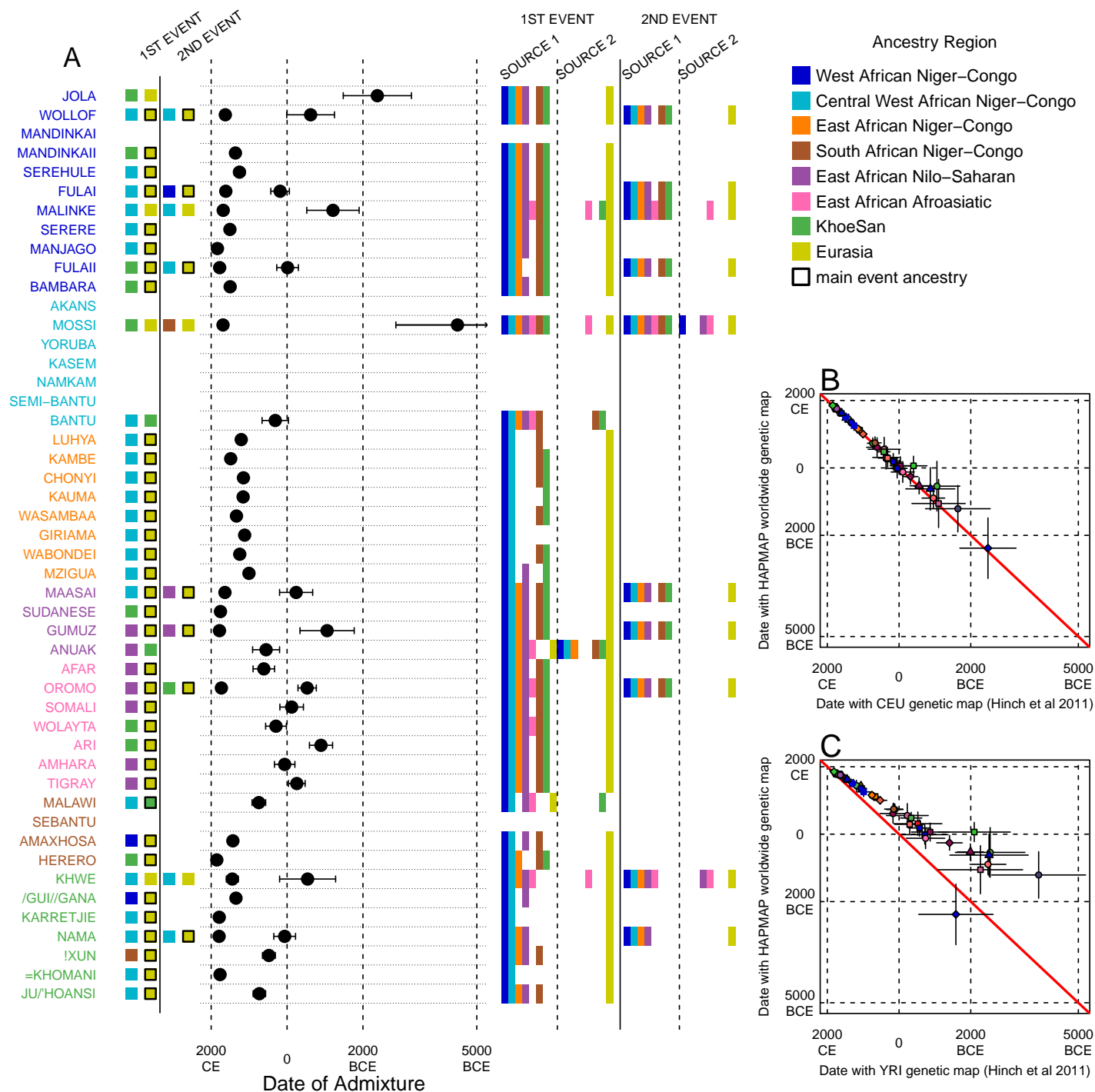


Figure 3. Inference of admixture in sub-Saharan Africa using MALDER. We used MALDER to identify the evidence for multiple waves of admixture in each population. (A) For each population, we show the ancestry region identity of the two populations involved in generating the MALDER curves with the greatest amplitudes (coloured blocks) for at most two events. The major contributing sources are highlighted with a black box. Populations are ordered by ancestry of the admixture sources and dates estimates which are shown ± 1 s.e. For each event we compared the MALDER curves with the greatest amplitude to other curves involving populations from different ancestry regions. In the central panel, for each source, we highlight the ancestry regions providing curves that are not significantly different from the best curves. In the Jola, for example, this analysis shows that, although the curve with the greatest amplitude is given by Khoesans (green) and Eurasian (yellow) populations, curves containing populations from any other African group (apart from Afroasiatic) in place of a Khoesan population are not significantly smaller than this best curve (SOURCE 1). Conversely, when comparing curves where a Eurasian population is substituted with a population from another group, all curve amplitudes are significantly smaller ($Z < 2$). (B) Comparison of dates of admixture ± 1 s.e. for MALDER dates inferred using the HAPMAP recombination map and a recombination map inferred from European (CEU) individuals from [Hinch et al., 2011]. We only show comparisons for dates where the same number of events were inferred using both methods. Point symbols refer to populations and are as in Figure 1. (C) as (B) but comparison uses an African (YRI) map. Source data can be found in Figure 3-Source Data 1

populations used to describe the recipient population’s admixture event(s) as *surrogates*, and populations used to paint both target and surrogate populations as *donors*. Including closely related individuals in chromosome painting analyses can cause the resulting painted chromosomes to be dominated by donors from these close genealogical relationships, which can mask signals of admixture in the genome [Hellenthal et al., 2014; van Dorp et al., 2015]. To help ameliorate this, we painted chromosomes for the GLOBETROTTER analysis by performing a fresh run of CHROMOPAINTER where we painted each individual from a recipient group with a set of donors which did not include individuals from within their own ancestry region. We additionally painted all (59) other surrogate populations with the same set of non-local donors, and used these copying vectors, together with the non-local painted chromosomes, to infer admixture. Using this approach, we found evidence of recent admixture in all African populations (Figure 4A). To summarise these events, we show the composition of the admixing source groups as barplots for each population coloured by the contribution from each African ancestry regions and Eurasia, alongside the inferred date with confidence interval determined by bootstrapping, and the estimated proportion of admixture (Figure 4). For each event we also identify the best matching donor population to the admixture sources. We describe observations from these analyses below, and note that the dates of admixture we infer indicate when gene-flow occurred between source populations and not the arrival of groups into an area, which may often be several generations earlier.

Direct and indirect gene flow from Eurasia back into Africa

We did not find evidence for recent Eurasian admixture in every African population (Figure 4). In particular, in several groups from South Africa and all from the Central West African ancestry region, which includes populations from Ghana, Nigeria, and Cameroon, we infer admixture between groups that are best represented by contemporary populations residing in Africa. As GLOBETROTTER is designed to identify the most recent admixture event(s) [Hellenthal et al., 2014], this observation does not rule out gene-flow from Eurasia back into these groups, but does suggest that subsequent movements between African groups were important in generating the contemporary ancestry of Central West and Southern African Niger-Congo speaking groups. With some exceptions that we describe below, we also do not observe Eurasian ancestry in all East African Niger-Congo speakers, instead finding more evidence for coancestry with Afroasiatic speaking groups. As we show later, Afroasiatic populations have a significant amount of genetic ancestry from outside of Africa, so the observation of this ancestry in several African groups identifies a route by which Eurasian ancestry may have indirectly entered the continent [Pickrell et al., 2014].

In fact, characterising admixture sources as mixtures allows us to infer whether Eurasian haplotypes are likely to have come directly into sub-Saharan Africa – in which case the admixture source will contain only Eurasian surrogates – or whether Eurasian haplotypes were brought indirectly together with sub-Saharan groups. In West African Niger-Congo speakers from The Gambia and Mali, we infer admixture involving minor admixture sources which contain mostly Eurasian (dark yellow) and Central West African (sky blue) ancestry, which most closely match the contemporary copying vectors of northern European populations (CEU and GBR) or the Fulani (FULAI, highlighted in gold in Figure 4A). The Fulani, a

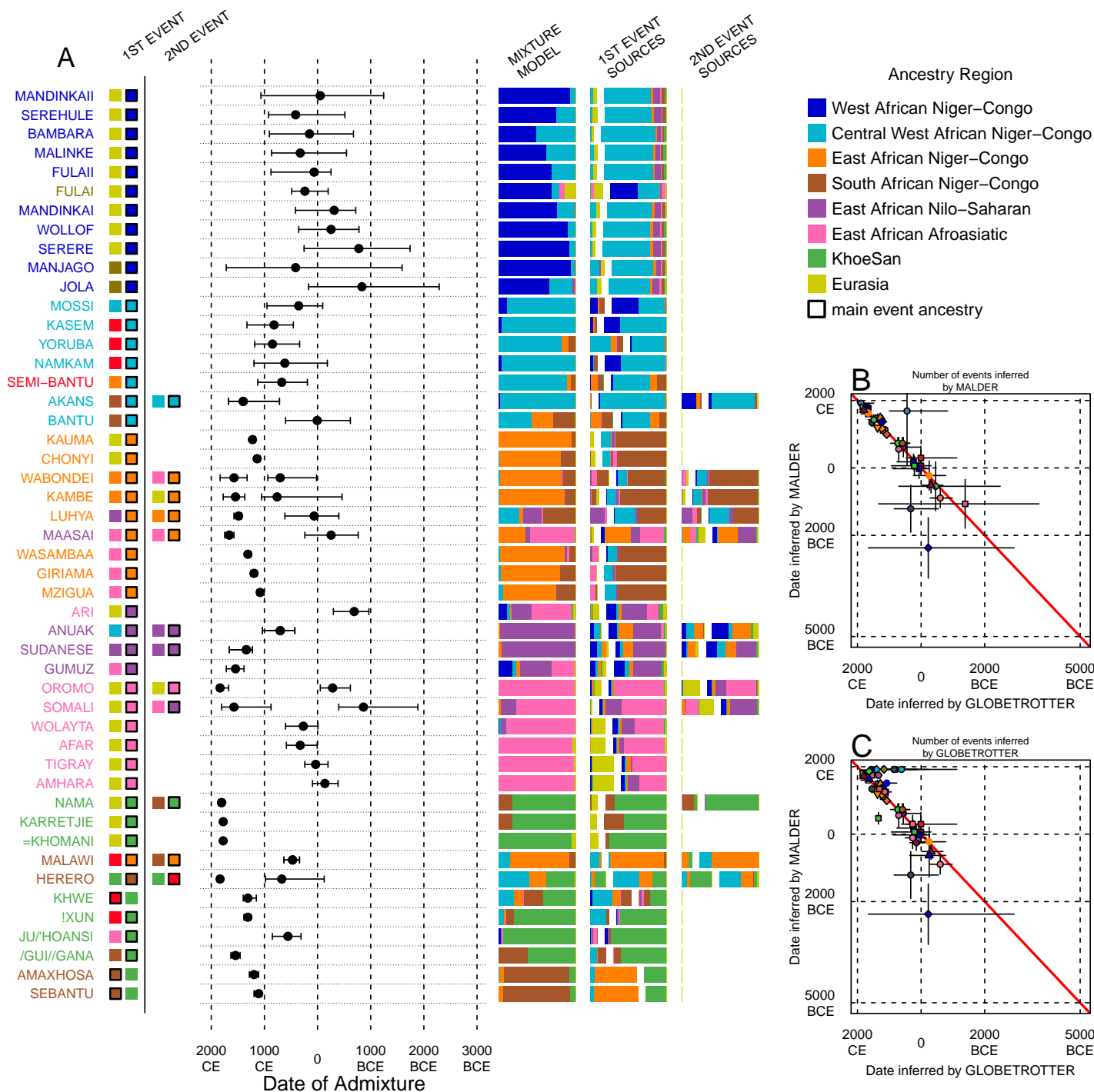


Figure 4. Inference of admixture in sub-Saharan African using GLOBETROTTER (A) For each group we show the ancestry region identity of the best matching source for the first and, if applicable, second events. Events involving sources that most closely match FULAI and SEMI-BANTU are highlighted by golden and red colours, respectively. Second events can be either multiway, in which case there is a single date estimate, or two-date in which case the 2ND EVENT refers to the earlier event. The point estimate of the admixture date is shown as a black point, with 95% CI shown with lines. MIXTURE MODEL: We infer the ancestry composition of each African group by fitting its copying vector as a mixture of all other population copying vectors. The coefficients of this regression sum to 1 and are coloured by broad ethno-linguistic origin. 1ST EVENT and 2ND EVENT SOURCES shows the ancestry breakdown of the admixture sources inferred by GLOBETROTTER, coloured by ancestry region as in the key top right. (B) and (C) Comparisons of dates inferred by MALDER and GLOBETROTTER. Because the two methods sometimes inferred different numbers of events, in (B) we show the comparison for based on the inferred number of events in the MALDER analysis, and in (C) for the number of events inferred by GLOBETROTTER. Point symbols refer to populations and are as in Figure 1 and source data can be found in Figure 4-Source Data 1.

nomadic pastoralist group found across West Africa, were sampled in The Gambia, at the very western edge of their current range, and have previously reported genetic affinities with Niger-Congo speaking, Sudanic, Saharan, and Eurasian populations [Tishkoff et al., 2009; Henn et al., 2012], consistent with the results of our mixture model analysis (Figure 4A). Admixture in the Fulani differs from other populations from this region, with sources containing greater amounts of Eurasian and Afroasiatic ancestry, but appears to have occurred during roughly the same period (c. 0CE; Figure 5).

The Fulani represent the best-matching surrogate to the minor source of recent admixture in the Jola and Manjago, which we interpret as resulting not from specific admixture from them into these groups, but because the mix of African and Eurasian ancestries in contemporary Fulani is the best proxy for the minor sources of admixture in this region. With the exception of the Fulani themselves, the major admixture source in groups across this region is a similar mixture of African ancestries that most closely matches contemporary Gambian and Malian surrogates (Jola, Serere, Serehule, and Malinke), suggesting ancestry from a common West African group within the last 3,000 years. The Ghana Empire flourished in West Africa between 300 and 1200CE, and is one of the earliest recorded African states [Roberts, 2007]. Whilst its origins are uncertain, it is clear that trade in gold, salt, and slaves across the Sahara, perhaps from as early as the Roman Period, as well as evolving agricultural technologies, were the driving forces behind its development [Oliver and Fagan, 1975; Roberts, 2007]. The observation of Eurasian admixture in our analysis is at least consistent with this moment in African history, and suggests that ancestry in groups from across this region of West Africa is the result of interactions through North Africa that were catalysed by trade across the Sahara.

We infer direct admixture from Eurasian sources in two populations from Kenya, where specifically South Asian populations (GIH, KHV) are the most closely matched surrogates to the minor sources of admixture (Figure 5). Interestingly, the Chonyi (1138CE: 1080-1182CE) and Kauma (1225CE: 1167-1254CE) are located on the so-called Swahili Coast, a region where Medieval trade across the Indian Ocean is historically documented [Allen, 1993]. In the Kambe, the third group from coastal Kenya, we infer two events, the more recent one involving local groups, and the earlier event involving a European-like source (GBR, 761CE: 461BCE-1053CE). In Tanzanian groups from the same ancestry region, we infer admixture during the same period, this time involving minor admixture sources with different, Afroasiatic ancestry: in the Giriama (1196CE: 1138-1254CE), Wasambaa (1312CE: 1254-1341CE), and Mzigua (1080: 1007-1138CE). Although the proportions of admixture from these sources differ, when we consider the major sources of admixture in East African Niger-Congo speakers, they are again similar, containing a mix of mainly local Southern Niger-Congo (Malawi), Central West African, Afroasiatic, and Nilo-Saharan ancestries.

In the Afroasiatic speaking populations of East Africa we infer admixture involving sources containing mostly Eurasian ancestry, which most closely matches the Tuscans (TSI, Figure 4). This ancestry appears to have entered the Horn of Africa in three distinct waves (Figure 5). We infer admixture involving Eurasian sources in the Afar (326CE: 7-587CE), Wolayta (268CE: 8BCE-602CE), Tigray (36CE: 196BCE-240CE), and Ari (689BCE:965-297BCE). There are no Middle Eastern groups in our analysis, and this latter group of events may represent previously observed migrations from the Arabian peninsular from the same time [Pagani et al., 2012; Hodgson et al., 2014a]. Considering Afroasiatic and Nilo-Saharan speakers

separately, the ancestry of the major sources of admixture of the former are predominantly local (purple), indicative of less historical interaction with Niger-Congo speakers due to their previously reported Middle Eastern ancestry [Pagani et al., 2012]. In Nilo-Saharan speaking groups, the Sudanese (1341CE: 1225-1660), Gumuz (1544CE: 1384-1718), Anuak (703: 427-1037CE), and Maasai (1646CE: 1584-1743CE), we infer greater proportions of West (blue) and East (orange) African Niger-Congo speaking surrogates in the major sources of admixture, indicating both that the Eurasian admixture occurred into groups with mixed Niger-Congo and Nilo-Saharan / Afroasiatic ancestry, and a clear recent link with Central and West African groups.

In two Khoesan speaking groups from South Africa, the ≠Khomani and Karretjie, we infer very recent direct admixture involving Eurasian groups most similar to Northern European populations, with dates aligning to European colonial period settlement in Southern Africa (c. 5 generations or 225 years ago; Figure 5) [Hellenthal et al., 2014]. Taken together, and in addition the MALDER analysis above, these observations suggest that gene flow back into Africa from Eurasia has been common around the edges of the continent, has been sustained over the last 3,000 years, and can often be attributed to specific and different historical time periods.

Population movements within Africa and the Bantu expansion

Admixture events involving sources that best match populations from within Africa tend to involve local groups. Even so, there are several long range admixture events of note. In the Ju/'hoansi, a San group from Namibia, we infer admixture involving a source that closely matches a local southern African Khoesan group, the Karretjie, and an East African Afroasiatic, specifically Somali, source at 558CE (311-851CE). We infer further events with minor sources most similar to present day Afroasiatic speakers in the Maasai (1660CE: 1573-1747CE and 254BCE: 764-239BCE) and, as mentioned, all four Tanzanian populations (Wabondei, Wasambaa, Giriama, and Mzigua). In contrast, the recent event inferred in the Luhya (1486: 1428-1573CE) involves a Nilo-Saharan-like minor source. The recent dates of these events imply not only that Eastern Niger-Congo speaking groups have been interacting with nearby Nilo-Saharan and Afroasiatic speakers – who themselves have ancestry from more northerly regions – after the putative arrival of Bantu-speaking groups to Eastern Africa, but also that the major sources of admixture contained both Central West and a majority of Southern Niger-Congo ancestry.

With the exception of Austronesian in Madagascar, African languages can be broadly classified into four major macro-families: Afroasiatic, Nilo-Saharan, Niger-Congo, and Khoesan [Blench, 2006]. Most of the sampled groups in this study, and indeed most sub-Saharan Africans, speak a language belonging to the Niger Congo linguistic phylum [Greenberg, 1972; Nurse and Philippson, 2003]. A sub-branch of this group are the so-called “Bantu” languages – a group of approximately 500 very closely related languages – that are of particular interest because they are spoken by the vast majority of Africans south of the line between Southern Nigeria/Cameroon and Somalia [Pakendorf et al., 2011]. Given the their high similarity and broad geographic range, it is likely that Bantu languages spread across Africa quickly. Whether this cultural expansion was accompanied by people is an active research question, but an increasing number of molecular studies, mostly using uni-parental genetic markers, indicate that the expansion of languages

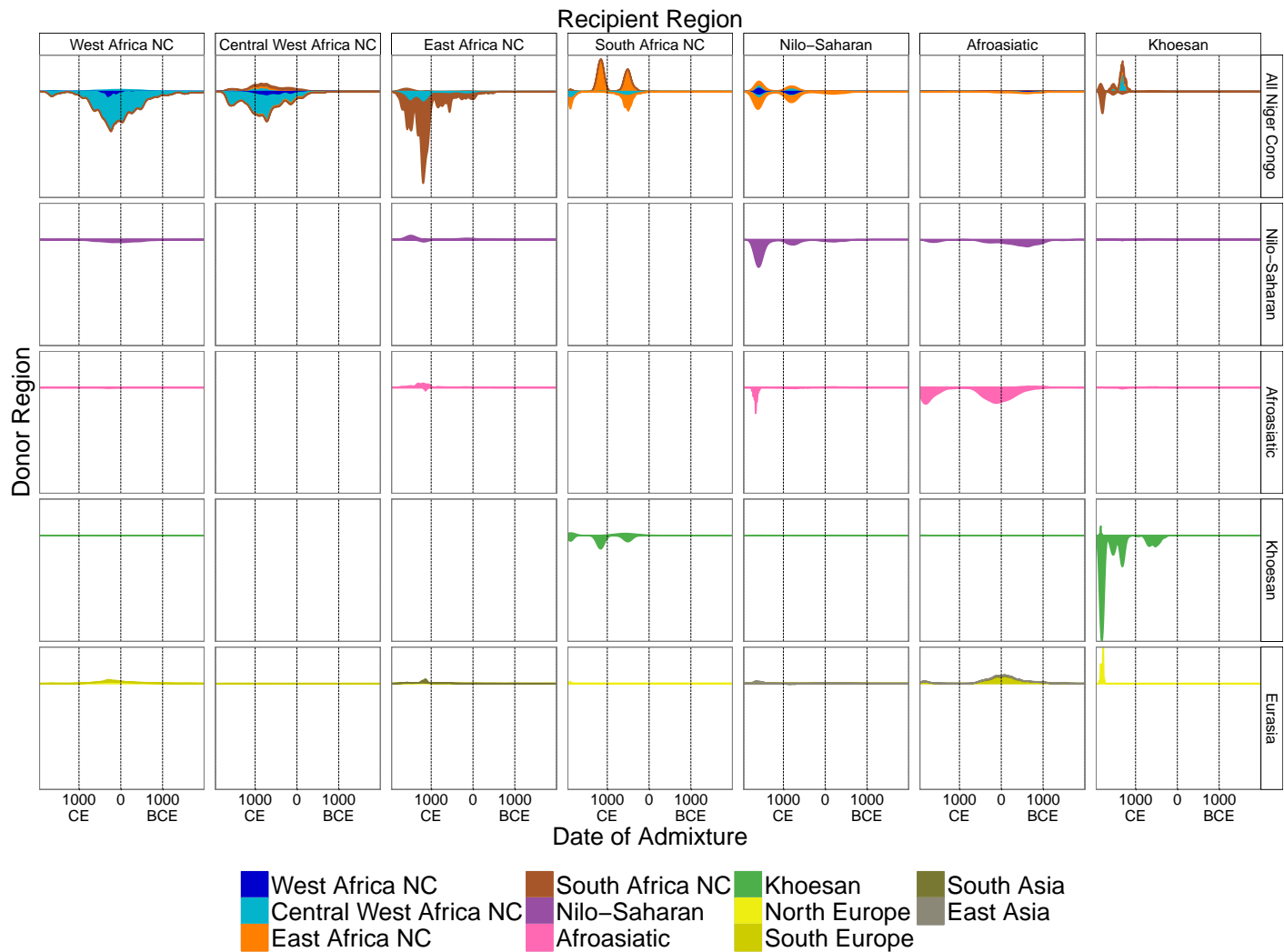


Figure 5. A timeline of recent admixture in sub-Saharan Africa. For all events involving recipient groups from each ancestry region (columns) we combine all date bootstrap estimates generated by GLOBETROTTER and show the densities of these dates separately for the minor (above line) and major (below line) sources of admixture. Dates are additionally stratified by the ancestry region of the surrogate populations (rows), with all dates involving Niger Congo speaking regions combined together (All Niger Congo). Within each panel, the densities are coloured by the geographic (country) origin of the surrogates and in proportion to the components of admixture involved in the admixture event. The integrals of the densities are proportional to the admixture proportions of the events contributing to them.

was accompanied by the diffusion of people [Beleza et al., 2005; Berniell-Lee et al., 2009; Pakendorf et al., 2011; de Filippo et al., 2012; Ansari Pour et al., 2013; Li et al., 2014; González-Santos et al., 2015]. Bantu languages can themselves be divided into three major groups: northwestern, which are spoken by groups near to the proto-Bantu heartland of Nigeria / Cameroon; western Bantu languages, spoken by groups situated down the west coast of Africa; and eastern, which are spoken across East and Central Africa [Li et al., 2014]. One particular debate concerns whether eastern languages are a primary branch that split off before the western groups began to spread south (the early-split hypothesis) or whether this occurred

after the migrants had begun to spread south (the late-split hypothesis) [Pakendorf et al., 2011]. Recent linguistic [Holden, 2002; Currie et al., 2013; Grollemund et al., 2015], and genetic analyses [Li et al., 2014] support the latter.

Whilst the current dataset does not cover all of Africa, and importantly contains no hunter-gather groups outside of southern Africa, we explored whether our admixture approach could be used to gain insight into the Bantu expansion. Specifically, we wanted to see whether the dates of admixture and composition of admixture sources were consistent with either of the two major models of the Bantu expansion. The major sources of admixture in East African Niger-Congo speakers tend to have both Central West and Southern Niger-Congo ancestry, although it is predominantly the latter (Figure 4). If Eastern Niger-Congo speakers derived all of their ancestry directly from Central West Africa (i.e. the early-split hypothesis) then we would not observe any ancestry from more southerly groups, but this is not what we see. That all East African Niger-Congo speakers that we sampled have admixture ancestry from a Southern group (Malawi), suggests that this group is more closely related to their Bantu ancestors than Central West Africans on their own. To explore this further, we performed additional GLOBETROTTER analyses where we restricted the surrogate populations used to infer admixture. For example, we re-ran the admixture inference process without allowing any groups from within the same ancestry region to contribute to the admixture event. As in the main analysis, when we disallow any East African Niger-Congo speakers from forming the admixture, Southern African Niger-Congo speakers, specifically Malawi, contribute most to the major source of admixture. In Southern African Niger-Congo groups, restricting admixture surrogates to those outside of their region leads to East African groups contributing most to the sources of admixture. The closer affinity between East and Southern Niger-Congo Bantu components is consistent with a common origin of these groups after the split from the Western Bantu. Moreover, this observation provides additional support for the hypothesis based on linguistic phylogenetics [Currie et al., 2013; Grollemund et al., 2015] that the main Bantu migration moved south and then east around the Congo rainforest. We performed a further restriction, disallowing any Southern or Eastern Niger-Congo speaking groups from being involved in admixture across both the South and East African Niger-Congo regions, which led to the vast majority of the ancestry of the major sources of admixture originating from Central West Africa, and almost exclusively from Cameroon populations (Figure 4-figure supplement 1 and Figure 4-figure supplement 2).

We inferred recent admixture involving major sources most similar to East or Southern African Niger-Congo speakers in the SEBantu (1109:1051-1196CE) and AmaXhosa (1196CE: 1109-1283CE) from South Africa. We infer two admixture dates in a Bantu-speaking group from Namibia, the Herero (1834CE: 1805-1892CE and 674CE: 124BCE-979CE), and a single date in the Khoesan-speaking Khwe (1312; 1152-1399CE), involving sources that more clearly contain ancestry from the Semi-Bantu. (Semi-Bantu and Bantu refer to two contemporary ethnic groups that currently reside in Cameroon.) In a third south-west African group, the !Xun from Angola, we infer admixture from a similar Cameroon-like source at around the same time as the Khwe (1312CE: 1254-1385CE). Assuming that Cameroon populations are the best proxy for Bantu ancestry in our dataset, these results suggest either a separate, more recent, arrival for Niger-Congo (Bantu) ancestry in south-west compared to south-east Africa, with the former coming recently directly down the west coast of Africa, specifically from Cameroon, and the latter deriving their

ancestry from earlier interactions via an eastern route [de Filippo et al., 2012; Li et al., 2014], or that there was limited detectable gene-flow in south-east Africa from the Bantus that eventually moved into south-west Africa.

Interestingly, in individuals from Malawi we infer a multi-way event with an older date (471: 340-631CE) involving a minor source which mostly contains ancestry from Cameroon, an event that is also seen in the Herero from Namibia. This Bantu admixture appears to have preceded that in some of the other South Africans by a few hundred years, suggesting either input from both western and eastern migrating Bantu speakers after the two branches split c.1,500 years ago, or that present day Malawi harbour ancestry from groups intermediate between the two branches [de Filippo et al., 2012]. Within this group we also see an admixture source with a significant proportion of non-Bantu (green) ancestry (2nd event, minor source), ancestry which we do not observe in the mixture model analysis, but which is also evident in other south-east African Niger-Congo speakers, the AmaXhosa and SEBantu, indicating that gene-flow must have occurred between the expanding Bantus and the resident hunter-gatherer groups [Marks et al., 2014]. The complex ancestry composition of the sources of these events, together with the dates of admixture, suggest large-scale interactions between groups during the Bantu expansion, and also hint that the expansion may have generated multiple waves of movement south. These events have resulted in recent haplotype sharing between groups across sub-Saharan Africa.

A haplotype-based model of gene flow in sub-Saharan Africa

Our haplotype-based analyses support a complex and dynamic picture of recent historical gene flow in Africa. We next attempted to summarise these events by producing a demographic model of African genetic history (Figure 6). Using genetics to infer historical demography will always depend somewhat on the available samples and population genetics methods used to infer population relationships. The model we present here is therefore unlikely to recapitulate the full history of the continent and will be refined in the future with additional data and methodology. We again caution that we do not have an exhaustive sample of African populations, and in particular lack significant representation from extant hunter-gatherer groups outside of southern Africa. Nevertheless, there are signals in the data that our haplotype-based approach allows us to pick up, and it is therefore possible to highlight the key gene-flow events and sources of coancestry in Africa, which we visualise in Figure 6:

- (1) **Colonial Era European admixture in the Khoesan.** In two southern African Khoesan groups we see very recent admixture involving northern European ancestry which likely resulted from Colonial Era movements from the UK, Germany, and the Netherlands into South Africa [Thompson, 2001].
- (2) **The recent arrival of the Western Bantu expansion in southern Africa.** Central West African, and in particular ancestry from Cameroon (red ancestry in Figure 6A), is seen in Southern African Niger-Congo and Khoesan speaking groups, the Herero, Khwe and !Xun, indicating that the gradual diffusion of Bantu ancestry reached the south of the continent only within the last 750 years. Central West African ancestry in Malawi appears to have appeared prior to this event.

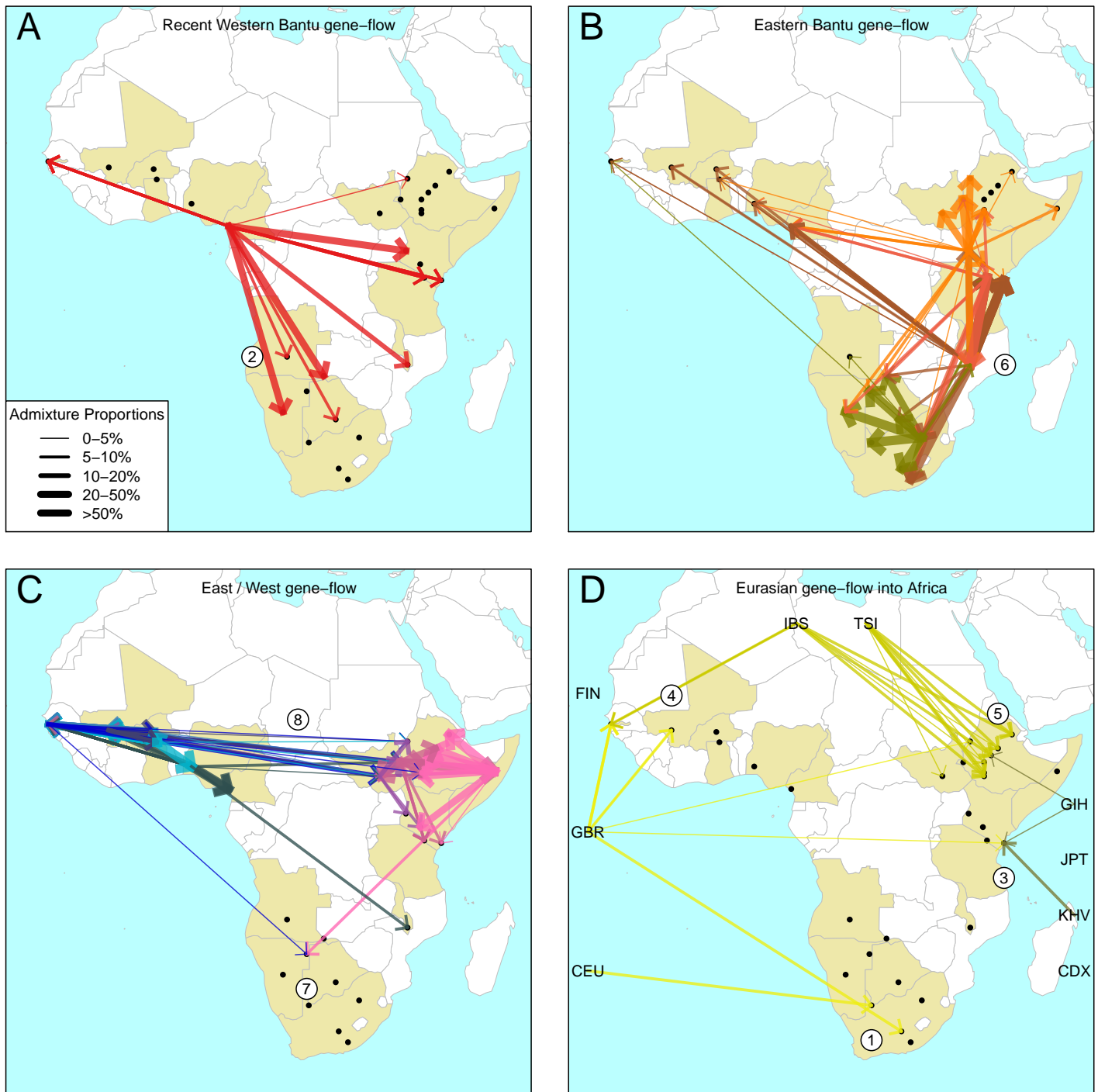


Figure 6. The geography of recent gene-flow in Africa. We summarise gene-flow events in Africa using the results of the GLOBE-TROTTER analysis. For each ethnic group, we inferred the composition of the admixture sources, and link recipient population to surrogates using arrows, the size of which is proportional to the amount it contributes to the admixture event. We separately plot (A) all events involving admixture source components from the Bantu and Semi-Bantu ethnic groups in Cameroon; (B) all events involving admixture sources from East and South African Niger-Congo speaking groups; (C) events involving admixture sources from West African Niger-Congo and East African Nilo-Saharan / Afroasiatic groups; (D) all events involving components from Eurasia. In (D) arrows are linked to the labelled 1KGP Eurasian groups. Arrows are coloured by country of origin, as in Figure 5. Numbers 1-8 in circles represent the events highlighted in section *A haplotype-based model of gene flow in sub-Saharan Africa*. An alternative version of this plot, stratified by date, is shown in Figure 6-figure supplement 1.

- (3) **Medieval contact between Asia and the East African Swahili Coast.** Specific Asian gene-flow is observed into two coastal Kenyan groups, the Kauma and Chonyi, which represents a distinct route of Eurasian, in this case Asian, ancestry into Africa, perhaps as a result of Medieval trade networks between Asia and the Swahili Coast around 1200CE.
- (4) **Gene-flow across the Sahara.** Over the last 3,000 years, admixture involving sources containing northern European ancestry is seen on the Western periphery of Africa, in The Gambia and Mali. This ancestry in West Africa is likely to be the result of more gradual diffusion of DNA across the Sahara from northern Africa and across the Iberian peninsular, and not via the Middle East, as in the latter scenario we would expect to see Spanish (IBS) and Italian (TSI) in the admixture sources. We do see limited southern European ancestry in West Africa (Figs. 5 and 6D) in the Fulani, suggesting that some Eurasian ancestry may also have entered West Africa via North East Africa [Henn et al., 2012].
- (5) **Several waves of Mediterranean / Middle Eastern ancestry into north-east Africa.** We observe southern European gene flow into East African Afroasiatic speakers over a more prolonged time period over the last 3,000 years, with a major wave 2,000 years ago (Figs. 5 and 6D). We do not have Middle-Eastern groups in our analysis, so the observed Italian ancestry in the minor sources of admixture – the Tuscans are the closest Eurasian group to the Middle East – is consistent with previous results using the same samples [Pagani et al., 2012; Hodgson et al., 2014a], indicating this region as a major route for the back migration of Eurasian DNA into sub-Saharan Africa [Pagani et al., 2012; Pickrell et al., 2014].
- (6) **The late split of the Eastern Bantus.** The major source of admixture in East Africa Niger-Congo speakers is consistently a mixture of Central West Africa and Southern Niger-Congo speaking groups, in particular Malawi. This result best fits a model where Bantu speakers initially spread south along the western side of the Congo rainforest before splitting off eastwards, and interacting with local groups in central south Africa – for which Malawi is our best proxy – and then moving further north-east and south (Figure 6B).
- (7) **Pre-Bantu pastoralist movements from East to South Africa.** In the Ju/'hoansi we infer an admixture event involving an East African Afroasiatic source. This event precedes the arrival of Bantu-speaking groups in southern Africa, and is consistent with several recent results linking Eastern and Southern Africa and the limited spread of cattle pastoralism prior to the Bantu expansion (Figs. 5 and 6C) [Pickrell et al., 2014; Ranciaro et al., 2014; Macholdt et al., 2015; Barham and Mitchell, 2008].
- (8) **Ancestral connections between West Africans and the Sudan.** Concentrating on older events, we observe old “Sudanese” (Nilotic) components in very small proportions in the Gambia (Figure 4-figure supplement 1 and Figure 5) which may represent ancient expansion relationships between East and West Africa. When we infer admixture in West and Central West African groups without allowing any West Africans to contribute to the inference, we observe a clear signal of Nilo-Saharan

ancestry in these groups, consistent with bidirectional movements across the Sahel [Tishkoff et al., 2009] and coancestry with (unsampled) Nilo-Saharan groups in Central West Africa. Indeed, if we look again at the PCA in Figure 1C, we observe that the Nilo-Saharan speakers are between West and East African Niger-Congo speaking individuals on PC3, an affinity which is supported by the presence of West African components in non-Niger-Congo speaking East Africans (Fig 6C).

- (9) **Ancient Eurasian gene-flow back into Africa and shared hunter-gatherer ancestry.** The MALDER analysis and f_3 statistics show the general presence of ancient Eurasian and/or Khoesan ancestry across much of sub-Saharan Africa. We tentatively interpret these results as being consistent with recent research suggesting very old (>10 kya) migrations back into Africa from Eurasia [Hodgson et al., 2014a], with the ubiquitous hunter-gatherer ancestry across the continent possibly related to the inhabitant populations present across Africa prior to these more recent movements. Future research involving ancient DNA from multiple African populations will help to further characterise these observations.

Discussion

Here we present an in-depth analysis of the genetic history of sub-Saharan Africa in order to characterise its impact on present day diversity. We show that gene-flow has taken place over a variety of different time scales which suggests that, rather than being static, populations have been sharing DNA, particularly over the last 3,000 years. An unanswered question in African history is how contemporary populations relate to those present in Africa before the transition to pastoralism that began some 2,000 years ago. Whilst the f_3 and MALDER analyses show evidence for deep Eurasian and some hunter-gatherer ancestry across Africa, our GLOBETROTTER analysis provides a greater precision on the admixture sources and a timeline of events and their impact on groups in our analysis (Figure 6). The transition from foraging to pastoralism and agriculture in Africa was likely to be complex, with its impact on existing populations varying substantially. However, the similarity of language and domesticates in different parts of Africa implies that this transition spread; there are few cereals or domesticated animals that are unique to particular parts of Africa. Whilst herding has likely been going on for several thousand years in some form in Africa, 2,000 years ago much of Africa still remained the domain of hunter-gatherers [Barham and Mitchell, 2008]. Our analysis provides an attempt at timing the spread of Niger-Congo speakers east from Central West Africa. Although we do not have representative forager (hunter-gatherer) groups from all parts of sub-Saharan Africa, we observe that in addition to their local region, many of the sampled groups share ancestry with Central West African groups, which is likely the result of genes spreading with the Bantu agriculturalists as their farming technology spread across Africa.

After OCE we begin to see admixture events shared between East and West Africa, with predominantly West to East direction, which appear to be most extensive around 1,000 years ago. During this time we see evidence of admixture from West, and West Central Africa into southern Khoesan speaking groups as well as down the Eastern side of the continent. Below, we outline some of the important technical challenges in using genetic data to interpret historical events exposed by our analyses. Nonetheless,

the study presented here shows that patterns of haplotype sharing in sub-Saharan African are largely determined by historical gene-flow events involving groups with ancestry from across and outside of the continent.

Interpreting haplotype similarity as historical admixture

Analyses that rely on correlations in allele frequencies (such as those performed here in the *Widespread evidence for admixture* section) provided initial evidence that the presence of Eurasian DNA across sub-Saharan Africa is the result of gene flow back into the continent within the last 10,000 years [Gurdasani et al., 2014; Pickrell et al., 2014; Hodgson et al., 2014a]. In addition, some groups have ancient (over 5 kya) shared ancestry with hunter-gather groups (Figure 3) [Gurdasani et al., 2014]. Whilst the weighted admixture LD decay curves between pairs of populations suggests that this admixture involved particular groups, for two main reasons, the interpretation of such events is difficult. Firstly, because our dataset includes closely related groups, in many populations, multiple pairs of reference populations generate significant admixture curves. Although the MALDER analysis helps us to identify which of these groups is closest to the true admixing source, it is not always possible to identify a single best matching reference, implying that sub-Saharan African groups share some ancestry with many different extant groups. So, as confirmed with the GLOBETROTTER analysis, it is unlikely that any single contemporary population adequately represents the true admixing source. On this basis of these analyses alone, it is not possible to characterise the composition of admixture sources.

Secondly, when we infer ancient events with ALDER, such as in the Mossi from Burkina Faso, where we estimate admixture around 5,000 years ago between a Eurasian (GBR) and a Khoesan speaking group (/Gui //Gana), we know that modern haplotypes are likely to only be an approximation of ancestral diversity [Pickrell and Reich, 2014]. Even the Ju/'hoansi, a San group from southern Africa traditionally thought to have undergone limited recent admixture, has experienced gene flow from non-Khoesan groups within this timeframe (Figure 4) [Pickrell et al., 2012, 2014]. Our result hints that the Mossi share deep ancestry with Eurasian and Khoesan groups, but any description of the historical event leading to this observation is potentially biased by the discontinuity between extant populations and those present in Africa in the past. In fact, this is one motivation for grouping populations into ancestry regions and defining admixture sources in this way. In this example we can (as we do) refer to Eurasian ancestry in general moving back into Africa, rather than British DNA in particular. Using contemporary populations as proxies for ancient groups is not the perfect approach, but it is the best that we have in parts of the world from which we do not have (for the moment) DNA from significant numbers of ancient human individuals, at sufficient quality, with which to calibrate temporal changes in population genetics.

An alternative approach is to characterise admixture as occurring between sources that have mixed ancestry themselves, to account for the fact that whilst groups in the past were not intact in the same way as they are now, DNA from such groups is nevertheless likely to be present in extant groups today. Haplotype-based methods allow for this type of analysis [Hellenthal et al., 2014; Leslie et al., 2015], and have the additional benefit of potentially identifying hidden structure and relationships (Figure 2). Whilst this approach still requires one to label groups by their present-day geographic or ethno-linguistic identity,

it at least allows for the construction of admixture sources containing ancestry from multiple different contemporary groups, which is likely to be closer to the truth. It also provides a framework for describing historical admixture in terms of gene flow networks, which taps into the increasingly supported notion that the genetic variation present in the vast majority of modern human groups is the result of mixture events in the past [Pickrell and Reich, 2014; Hellenthal et al., 2014]. One drawback of such an approach is that it is not always possible to assign a specific historical population label to mixed admixture sources, but it nevertheless accommodates the need to express admixture sources as containing multiple different ancestries.

Admixture in the context of infectious disease studies

Whilst inference of historical events are interesting in their own right, they also provide an important background for conducting large-scale genetic epidemiology in sub-Saharan Africa. Two study designs are of particular relevance: genetics association studies and inference of historical natural selection. Principally, the haplotype analysis described here suggests that while recent gene-flow events mean that absolute allele frequency differences across the sub-Sahara are small relative to between-Eurasian groups, there remains substantial haplotype diversity. Genetic studies often rely on assumptions about the similarity of haplotypes between two groups, either as part of testing for association, or in order to impute genetic variation that is not assayed directly. Extensive recent gene-flow can help in this regard by reducing the differences between groups. However, the clear signals of admixture suggest that novel and divergent haplotypes are likely to be pocketed with subtle variation occurring both ethno-geographically, due to differences in ancestry which pre-dates recent movements, and along the genome, either by chance or by the differential effects of natural selection.

When new haplotypes are introduced into a population their fate will be partly determined by the selective advantage they confer, as well as the action of genetic drift. These selective forces could include response to changes in infectious diseases, climate, and the cultural environment [Coop et al., 2009; Fumagalli et al., 2015]. Additionally, important causes of morbidity and ill-health in Africa are due to highly polymorphic parasites like *Plasmodium falciparum*, where moving into new environments might lead to exposure to new strains. An implication of widespread gene-flow is that it can provide a route for potentially beneficial novel mutations to enter populations allowing them to adapt to such change. Recent examples of this phenomena include the presence of altitude adaptations in Tibetans from admixture with Sherpa [Jeong et al., 2014] and Denisovans [Huerta-Sánchez et al., 2014]; higher than expected frequencies of the Duffy-null mutation in populations from Madagascar as a result of admixture with African Bantu speaking groups [Hodgson et al., 2014b]; and the observation of shared haplotype(s) in humans and Neanderthals immunity genes belonging to the Toll-like Receptor (*TLR*) gene cluster as a result of archaic admixture [Deschamps et al., 2016; Dannemann et al., 2016].

In the context of malaria, the spread of the Duffy-null allele, an ancient mutation which arose at least 30,000 years ago [Hamblin and Di Rienzo, 2000; Hamblin et al., 2002] and which confers resistance to *P. vivax* malaria, throughout Africa is only possible through contact and gene flow between populations right across the sub-Sahara. Conversely, the same sickle-cell causing mutation appears to have recently

occurred five times independently in Africa, causing multiple distinct haplotypes to be observed [Hedrick, 2011]. These mutations are young, within the order of 250-1750 years old [Curat et al., 2002; Modiano et al., 2008], so will have had limited opportunity to have been moved around by the gene-flow events that we describe.

We observed that some of this gene-flow may have coincided with the spread of Bantu languages and the concomitant introduction of new agricultural practices, which can jointly introduce both new genes and novel selective pressures into existing populations, such as specific agricultural products or unknown communicable disease. As an example, our analysis also sheds light on the spread of pastoralism into southern Africa prior to the Bantu expansion [Ehret, 1967; Guldemann, 2008; Henn et al., 2008]. Several recent resequencing studies have described the the distribution of mutations in and around the lactase gene (LCT) in African populations which are associated with the ability to digest lactase in adult life, lactase persistence (LP) [Ranciaro et al., 2014; Macholdt et al., 2015; Breton et al., 2014]. The C-14010 LP-associated variant, thought to have originated in East African Niger-Congo speaking groups [Tishkoff et al., 2007], was observed in the the San and !Xhosa, with the haplotype in the latter matching that of east African Bantu speaking populations [Ranciaro et al., 2014]. The same LCT variant was also observed at appreciable frequencies in the Nama [Breton et al., 2014]. We infer an admixture event in the Ju/'hoansi involving an Afroasiatic source, which at 732CE (616-993CE), is earlier than most of the Bantu admixture we observe in other southern African populations. In the Nama we infer a very recent admixture event (c. 1800CE) involving a source containing Niger-Congo (Malawi) ancestry. Whilst the event in the Nama is unlikely to have delivered the East African mutation to this group, the event in the Ju/'hoansi is consistent with pre-Bantu interactions between east African and southern African groups. Conversely, in the AmaXhosa the Bantu source of admixture derives most of its ancestry from East Africa, so the introduction the C-14010 mutation into this part of Africa may well have been the result of the Bantu expansion.

These examples demonstrate the utility of detailed inference exploiting the complex information that is captured by large-scale genome-wide studies of genetic diversity. Africa has an exciting opportunity to be part of bringing the genetic revolution to bear in understanding and treating both infectious and non-communicable disease. Further exploration and interpretation the the rich genetic diversity of the continent will help in this endeavour.

Materials and Methods

1 Overview of the dataset

The dataset comprises a mixture of 2,504 previously published individuals from Africa and elsewhere (see below) plus novel genotypes on 1,712 sampled by the Malaria Genomic Epidemiology Network (MalariaGEN) Figure 1-Source Data 1. The MalariaGEN samples were a subset of those collected at 8 locations in Africa as part of a consortial project on genetic resistance to severe malaria: details of the study sites and investigators involved are described elsewhere [Malaria Genomic Epidemiology Network, 2014]. Samples were genotyped on the Illumina Omni 2.5M chip in order to perform a multicentre genome-wide association study (GWAS) of severe malaria: initial GWAS findings from The Gambia, Kenya and Malawi have already been reported [Band et al., 2013; Malaria Genomic Epidemiology Network, 2015] and a manuscript describing findings at all 8 locations is in preparation

The MalariaGEN samples used in the present analysis were selected to be representative of the main ethnic groups present at each of the 8 African study sites. We screened the samples collected at each study site (typically >1000 individuals) to select individuals whose reported parental ethnicity matched their own ethnicity. This process identified 23 ethnic groups for which we had samples for approximately 50 unrelated individuals or more. For ethnic groups with more than 50 samples available, we performed a cluster analysis on cohort-wide principle components, generated as part of the GWAS, with the R statistical programming language [R Development Core Team, 2011] using the *MClust* package [Fraley et al., 2012], choosing individuals from the cluster containing the largest number of individuals, to avoid any accidental inclusion of outlying individuals and to ensure that the 50 individuals chosen were, when possible, relatively genetically homogeneous. We note that in several ethnic groups (Malawi, the Kambe from Kenya, and the Mandinka and Fula from The Gambia; Fig. 1-figure supplement 2) PCA of the genotype data showed a large amount of population structure. In these cases we chose two sub-groups of individuals from a given ethnic group, selected to represent the diversity of ancestry depicted by the PCs. In several other cases, following GWAS quality control (see below), genotype data for fewer than 50 control individuals was available, and in these cases we chose as many individuals as possible, regardless of the PC-based clustering or case/control status.

We additionally included further individuals from each of four Gambian ethnic groups: the Fula, Mandinka, Jola, Wolof. The genotype data from these individuals were included as the same individuals are also being sequenced as part of the **Gambian Genome Variation Project**¹. These subsets included ~30 trios from each ethnic group, information on which was used in phasing (see below). Full genome data from these individuals will be made available in the future.

1.1 Quality Control

Detailed quality control (QC) for the MalariaGEN dataset was performed for a genome-wide association study of severe malaria in Africa and is outlined in detail elsewhere [Malaria Genomic Epidemiology

¹The Gambian Genome Variation Project will sequence a number of full genomes from four Gambian ethnic groups as a basis for improving imputation for future West African specific GWAS

Network, 2015]. Briefly, genotype calls were formed by taking a consensus across three different calling algorithms (Illuminus, Gencall in Illumina’s BeadStudio software, and GenoSNP) [Band et al., 2013] and were aligned to the forward strand. Using the data from each country separately, SNPs with a minor allele frequency of <1% and missingness <5% were excluded, and additional QC to account for batch effects and SNPs not in Hardy-Weinberg equilibrium was also performed.

1.2 Combining the MalariaGEN populations with additional populations

The post-QC MalariaGEN data was combined with published data typed on the same Illumina 2.5M Omni chip from 21 populations typed for the 1000 Genomes Project (1KGP)², including and accounting for duos and trios, and with publicly available data from individuals from several populations from Southern Africa (Figure 1-Source Data 1) [Consortium, 2012; Schlebusch et al., 2012]. We merged samples to the forward strand, removing any ambiguous SNPs (A to T or C to G). Merging was checked by plotting allele frequencies between populations from both datasets, which should be generally correlated (data not shown). To describe population structure across sub-Saharan Africa we combined the above dataset with further publicly available samples typed on different Illumina (Omni 1M) chips, containing individuals from southern Africa [Petersen et al., 2013] and the Horn of Africa (Somalia/Ethiopia/Sudan)[Pagani et al., 2012] to generate a final dataset containing 4,216 individuals typed on 328,176 high quality common SNPs. To obtain the final set of analysis individuals we performed additional sample QC after phasing and removed American 1KG populations (Figure 1-Source Data 1).

1.3 Phasing

We used SHAPEITv2 [Delaneau et al., 2012] to generate haplotypically phased chromosomes for each individual. SHAPEITv2 conditions the underlying hidden Markov model (HMM) from Li and Stephens [2003] on all available haplotypes to quickly estimate haplotypic phase from genotype data. We split our dataset by chromosome and phased all individuals simultaneously, and used the most likely pairs of haplotypes (using the *-output-max* option) for each individual for downstream applications. We performed 30 iterations of the MCMC and used default values for all other parameters. As mentioned, we used known pedigree relationships to improve the phasing, using family data from both the 1KG and the Gambia Sequencing Projects.

1.4 Removing non-founders and cryptically related individuals

Our dataset included individuals who were known to be closely related (1KG duos and trios; Gambia Sequencing Project trios) and, because we took multiple samples from some population groups, there was also the potential to include cryptically related individuals. After phasing we therefore performed an additional step where we first removed all non-founders from the analysis and then identified individuals with high identity by descent (IBD), which is a measure of relatedness. Using an LD pruned set of SNPs generated by recursively removing SNPs with an $R^2 > 0.2$ using a 50kb sliding window, we calculated the

²data downloaded on 16th October 2013 from ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20120131_omni_genotypes_and_intensities/

proportion of loci that are IBD for each pair of individuals in the dataset using the R package *SNPRelate* [Zheng et al., 2012] and estimated kinship using the pi-hat statistic (the proportion of loci that are identical for both alleles (IBD=2) plus 0.5* the proportion of loci where one allele matches (IBD=1); i.e. $PI_{HAT} = P(IBD=2) + 0.5 * P(IBD=1)$). For any pair of individuals where $IBD > 0.2$, we randomly removed one of the individuals. 327 individuals were removed during this step.

1.5 1KG American populations and Native American Ancestry in 1KG Peruvians

With the exception of Peru, post-phasing we dropped all 1KG American populations from the analysis (97 ASW, 102 ACB, 107 CLM, 100 MXL and 111 PUR). We used a subset of the 107 Peruvian individuals that showed a large amount of putative Native American ancestry, with little apparent admixture from non-Amerindians (data not shown). Although Amerindians are not central to this study, and it is unlikely that there has been any recurrent admixture from the New World into Africa, we nevertheless generated a subset of 16 Peruvians to represent Amerindian admixture components in downstream analyses. When this subset was used, we refer to the population as PELII. The removal of these 606 American individuals left a final analysis dataset comprising 3,283 individuals from 60 different population groups (Figure 1-Source Data 1).

2 Analysis of population structure in sub-Saharan Africa

2.1 Principal Components Analysis

We performed Principal Components Analysis (PCA) using the *SNPRelate* package in R. We removed SNPs in LD by recursively removing SNPs with an $R^2 > 0.2$, using a 50kb sliding window, resulting in a subset of 162,322 SNPs.

2.2 Painting chromosomes with CHROMOPAINTER

We used fineSTRUCTURE [Lawson et al., 2012] to identify finescale population structure and to identify high level relationships between ethnic groups. The initial step of a fineSTRUCTURE analysis involves “painting” haplotypically phased chromosomes sequentially using an updated implementation of a model initially introduced by Li and Stephens [2003] and which is exploited by the CHROMOPAINTER package [Lawson et al., 2012]. The Li and Stephens copying model explicitly relates linkage disequilibrium to the underlying recombination process and CHROMOPAINTER uses an approximate method to reconstruct each “recipient” individual’s haplotypic genome as a series of recombination “chunks” from a set of sample “donor” individuals. The aim of this approach is to identify, at each SNP as we move along the genome, the closest relative genome among the members of the donor sample. Because of recombination, the identity of the closest relative will change depending on the admixture history between individual genomes. Even distantly related populations share some genetic ancestry since most human genetic variation is shared [The International HapMap Consortium, 2010; Ralph and Coop, 2013], but the amount of shared ancestry can differ widely. We use the term “painting” here to refer to the application of a different label to each of the donors, such that – conceptually – each donor is represented by a different

colour. Donors may be coloured individually, or in groups based on *a priori* defined labels, such as the geographic population that they come from. By recovering the changing identity of the closest ancestor along chromosomes we can understand the varying contributions of different donor groups to a given population, and by understanding the distribution of these chunks we can begin to uncover the historical relationships between groups.

2.3 Using painted chromosomes with fineSTRUCTURE

We used CHROMOPAINTER with 10 Expectation-Maximisation (E-M) steps to jointly estimate the program’s parameters N_e and θ , repeating this separately for chromosomes 1, 4, 10, and 15 and weight-averaging (using centimorgan sizes) the N_e and θ from the final E-M step across the four chromosomes. We performed E-M on 5 individuals from every population in the analysis and used a weighted average of the values across all pops to arrive at final values of 190.82 for N_e and 0.00045 for θ . We ran each chromosome from each population separately and combined the output to generate a final coancestry matrix to be used for fineSTRUCTURE.

As the focus of our analysis is population structure within Africa, we used a “continental force file” to combine all non-African individuals into single populations. The processing time of the algorithm is directly related to the number of individuals included in the analysis, so reducing the number of individuals speeds the analysis up. Furthermore, fineSTRUCTURE initially uses a prior that assumes that all individuals are equally distant from each other, which in the case of worldwide populations is likely to be untrue: African populations are likely to be more closely related to each other than to non-Africa populations, for example. The result is that not all of the substructure is identified in one run.

We therefore combined all individuals from each of the non-Africa 1KG populations into “continents”, which has the effect of combining all of the copying vectors from the individuals within them to look like (re-weighted) normal individuals but cannot be split and do not contribute to parameter inference, and can thus be considered as copying vectors that contain the average of the individuals within them. They are then included in the algorithm at minimal extra computational cost and exist primarily to provide chunks to (and from) the remaining groups. We combined all individuals from a labelled population (e.g. all IBS individuals were now contained in a “continent” grouping called IBS), with the exception of the three Chinese population CHB, CHD, and CHS where we combined all individuals into a single CHN continent.

2.4 Using fineSTRUCTURE to inform population groupings

The fineSTRUCTURE analysis identified 154 clusters of individuals, grouped on the basis of copying-vector similarity (Fig. 1-figure supplement 3). Some ethnic groups, such as Yoruba, Mossi, Jola and Ju/’hoansi form clusters containing only individuals from their own ethnic groups. In other cases, individuals from several different ethnic groups are shared across different clusters. We see this particularly in groups from The Gambia and Kenya. These are the two countries where the most ethnic groups were sampled, seven and four, respectively, and this putative mixing could be a result of sampling several

different ethnic groups from a limited geographic area. The reason that we do not see the same extent of haplotype sharing in Cameroon, for example, may be because we only have two ethnic groups sampled from that area. The mixed nature of the clusters in this region could also be the results of real ancestral relationships: some of the ethnic groups in The Gambia, for example, may have history of choosing partners from outside their own ethnic group.

As mentioned above, we combined data from different sources, and in some groups (e.g. the Fula) we specifically chose different groups of individuals in an attempt to cover the broad spectrum of ancestry present in that group. We used the fineSTRUCTURE tree to visually group individuals based on their ancestry. Our aim here is to try to maximise the number of individuals that we can include within an ethnic group, without merging together individuals that are distant on the tree. We also decided *not* to use the fineSTRUCTURE clusters themselves as analytical groups because of difficulties with the interpretation of the history of such clusters. We were interested in identifying the major admixture events that have occurred in the history of different populations, and it is not clear what an analytical group that is defined as, for example, a mixture of Manjago, Mandinka and Serere individuals, would mean in our admixture analyses. In practice, this meant that we used the original geographic population labelled groups for all populations except in the Fula and Mandinka from The Gambia, where individuals fell into two distinct groups of clusters. Here we defined two clusters for each group, with the two groups suffixed with an “I” or “II” (Figure 1-figure supplement 3).

2.5 Defining ancestry regions

We used a combination of genetic and ethno-linguistic information (see Supplementary File 1 below) to define seven ancestry regions in sub-Saharan Africa. The ancestry regions are reported in Figure 1-Source Data 1 and closely match the high level groupings we observed in the fineSTRUCTURE tree, with the following exceptions:

1. East African Niger-Congo speakers

- (a) The two ethnic groups from Cameroon – Bantu and Semi-Bantu – were included in the Central West African Niger-Congo ancestry region despite clustering more closely with East African groups from Kenya and Tanzania in Figure 1-figure supplement 3. In a preliminary fineSTRUCTURE analysis based on the MalariaGEN and 1KG populations only, using c. 1 million SNPs, the Cameroon populations clustered with other Central West African groups, and not East Africans (data not shown).
- (b) Malawi was included in the South African Niger-Congo ancestry region, despite being an outlying cluster in a clade with East African Niger-Congo speaking groups. A preliminary fineSTRUCTURE analysis based on the MalariaGEN, 1KG and Schlebusch populations, clustered Malawi with the Herero and SEBantu speakers (data not shown).

2. Southern Africa

- (a) We treated Southern African individuals slightly differently: even though the fineSTRUCTURE analysis did not split them into two separate clades of Khoesan and Niger-Congo

speaking individuals, we nevertheless did. Schlebusch et al. [2012] showed that these populations were inter-related and admixed, two properties in the data we were hoping to uncover. The final ancestry region assignments are outlined in Figure 1-Source Data 1.

2.6 Estimating pairwise F_{ST}

We used *smartpca* in the EIGENSOFT [Patterson et al., 2006] package to estimate pairwise F_{ST} between all populations. This implementation uses the Hudson estimator recently recommended by Bhatia et al. [2013]. Results are shown in Figure 2, Figure 2-Source Data 1 and Figure 2-Source Data 2.

2.7 Comparing sets of copying vectors

We used Total Variation Distance (TVD) to compare copying vectors [Leslie et al., 2015; van Dorp et al., 2015]. As the copying vectors are discrete probability distributions over the same set of donors, TVD is a natural metric for quantifying the difference between them. For a given pair of groups A and B with copying vectors describing the copying from n donors, a and b , we can compute TVD with the following equation:

$$TVD = 0.5 \times \sum_{i=1}^n (|a_i - b_i|)$$

3 Analysis of admixture in sub-Saharan African populations

We used a combination of approaches to explore admixture across Africa. Initially, we employed commonly used methods that utilise correlations in allele frequencies to infer historical relationships between populations. To understand ancient relationships between African groups we used the f_3 statistic [Reich et al., 2009] to look for shared drift components between a test population and two reference groups. We next used ALDER [Loh et al., 2013] and MALDER [Pickrell et al., 2014] – an updated implementation of ALDER that attempts to identify multiple admixture events – to identify admixture events through explicit modelling of admixture LD by generating weighted LD curves. The weightings of these curves are based on allele frequency differences, at varying genetic distances, between a test population and two putative admixing groups.

To identify more recent events we used two methods which aim to more fully model the mixed ancestry in a population by utilising the distribution and length of shared tracts of ancestry as identified with the CHROMOPAINTER algorithm [Lawson et al., 2012; Hellenthal et al., 2014]. We outline the details of this analysis below, but note here that, because this approach is based on the comparison and analysis of painted chromosomes, it offers a different perspective from approaches based on comparisons of allele frequencies.

3.1 Inferring admixture with the f_3 statistic and ALDER

We computed the f_3 statistic, introduced by Reich et al. [2009], as implemented in the TREEMIX package [Pickrell and Pritchard, 2012]. These tests are a 3-population generalization of F_{ST} , equal to the inner

product of the frequency differences between a group X and two other groups, A and B. The statistic, commonly denoted $f_3(X:A,B)$ is proportional to the correlated genetic drift between A and X and A and B. If X is related in a simple way to the common ancestor with A and B, we expect this quantity to be positive. Significantly negative values of f_3 suggest that X has arisen as a mixture of A and B, which is thus an unambiguous signal of mixture. Standard errors are computed using a block jackknife procedure in blocks of 500 SNPs (Supplementary Table 1).

This analysis shows that sub-Saharan African populations tend to have either a hunter-gatherer or Eurasian source contributing to their most significant f_3 statistics (Supplementary Table 1). That is, there is clear evidence of admixture involving both Eurasian and hunter-gatherer groups across sub-Saharan Africa. Whilst we do not infer admixture using this statistic in the Jola, Kasem, Sudanese, Gumuz, and Ju/'hoansi, in most other groups the most significant f_3 statistic includes either the Ju/'hoansi or a 1KGP European source (GBR, CEU, FIN, or TSI). Niger-Congo speaking groups from Central West and Southern Africa tend to show most significant statistics involving the Ju/'hoansi, where as West and East African and Southern Khoesan speaking groups tended to show most significant statistics involving European sources, consistent with an recent analysis on a similar (albeit smaller) set of African populations [Gurdasani et al., 2014].

We used ALDER [Patterson et al., 2012; Loh et al., 2013] to test for the presence of admixture LD in different populations. This approach works by generating weighted admixture curves for pairs of populations and tests for admixture. As noted in Loh et al. [2013] the use of f_3 statistics and weighted LD curves are somewhat complementary, and there are several reasons why f_3 statistics might pick up signals of admixture when ALDER does not. In particular, admixture identified using f_3 statistics but not by ALDER is potentially related to more ancient events because whilst shared drift signals will still be present, admixture LD will have been broken over (potentially) millennia of recombination.

As previously shown by Loh et al. [2013] and Pickrell et al. [2014], weighted LD curves can be used to identify the source of the gene flow by comparing curves computed using different reference populations. This is possible because theory predicts that the amplitude (i.e. the y-axis intercept) of these curves becomes larger as one uses reference populations that are closer to the true mixing populations. Loh et al. [2013] demonstrated that this theory holds even when using the admixed population itself as one of the reference populations. Pickrell et al. [2014] used this concept to identify west Eurasian ancestry in a number of East African and Khoesan speaking groups from southern Africa.

We thus initially ran ALDER in “one-reference” mode, where for each focal population, we generated curves involving itself with every other reference population in turn. We used the average amplitude of the curves generated in this way to identify the groups important in describing admixture in the history of the focal group. Figure 3-figure supplement 1 shows comparative plots to those by Pickrell et al. [2014] for a selection of African populations, including the Ju/'hoansi, who we also infer to have largest curve amplitudes with Eurasian groups, consistent with that previous analysis. We summarise the results of the analysis of curve amplitudes across all populations in Figure 3-figure supplement 2, which shows the rank of the curve amplitudes for each reference population across different sub-Saharan African groups. Across much of Africa, amplitudes from non-African (European) groups are amongst the highest, indicative of admixture from Eurasian-like sources.

Next, for each focal ethnic group in turn, we used ALDER to characterise admixture using all other ethnic groups as potential reference groups (i.e. in two-population mode). In effect, this approach compares every pair of reference groups, identifying those pairs that show evidence of shared admixture LD ($P < 0.05$ after multiple-hypothesis testing). As many of the groups are closely related, we often observed more than one pair of ethnic groups as displaying admixture in a given focal population, the results of which are highly correlated. In Supplementary Table 1, we show the evidence for admixture only for the pair of groups with the lowest P-value for each focal group. Dates for admixture events were generated using a generation time of 29 years [Fenner, 2005] and the following equation:

$$D = 1950 - (n + 1) * g$$

where D is the inferred date of admixture, n is the inferred number of generations since admixture, and g is the generation time in years.

3.2 Inferring multiple waves of admixture in African populations using weighted LD curves

We used MALDER [Pickrell et al., 2014], an implementation of ALDER designed to fit multiple exponentials to LD decay curves and therefore characterise multiple admixture events to allele frequency data. For each event we recorded (a) the curve, C , with the largest overall amplitude $C_{Pop1;Pop2}^{max}$, and (b) the curves which gave the largest amplitude where each of the two reference populations came from a different ancestry region, and for which a significant signal of admixture was inferred. To identify the source of an admixture event we compared curves involving populations from the same ancestry region as the two populations involved in generating $C_{Pop1;Pop2}^{max}$. For example, in the Jola, the population pair that gave C^{max} were the Ju/'hoansi and GBR. Substituting these populations for their ancestry regions we get $C_{Khoesan;Eurasia}^{max}$. To understand whether this event represents a specific admixture involving the Khoesan in the history of the Jola, we identified the amplitude of the curves from (b) of the form $C_{M;Eurasia}^{max}$, where M represents a population from any ancestry region other than Eurasia that gave a significant MALDER curve. We generated a Z-score for this curve comparison using the following formula [Pickrell et al., 2014]:

$$Z = \frac{C_{Khoesan;Eurasia}^{max} - C_{Khoesan;M}^{max}}{\sqrt{se(C_{Khoesan;Eurasia}^{max})^2 + se(C_{Khoesan;M}^{max})^2}}$$

The purpose of this was to determine, for a given event, whether the sources of admixture could be represented by a single ancestry region, in which case the overall $C_{ancestry1;ancestry2}^{max}$ will be significantly greater than curves involving other regions, or whether populations from multiple ancestry regions can generate admixture curves with similar amplitudes, in which case there will be a number of ancestry regions that best represent the admixing source. We combined all values of M where the Z-score computed from the above test gave a value of < 2 , and define the sources of admixture in this way.

To identify the major source of admixture, we performed a similar test. We determined the regional identity of the two populations used to generate C^{max} . In the example above, these are Khoesan and Eurasia. Separately for each region, we identify the curve, C , with the maximum amplitude where either of the two reference populations was from the Khoesan region, $C_{Khoesan}^{max}$ as well as the curve where neither of the reference populations was Khoesan, $C_{notKhoesan}^{max}$. We compute a Z-score as follows:

$$Z = \frac{C_{Khoesan}^{max} - C_{notKhoesan}^{max}}{\sqrt{se(C_{Khoesan}^{max})^2 + se(C_{notKhoesan}^{max})^2}}$$

This test generates two Z-scores, in this example, one for the Khoesan/not-Khoesan comparison, and one for the Eurasia/not-Eurasia comparison. We assign the main ancestry of an event to be the region(s) that generate $Z > 2$. If neither region generates a Z-score > 2 , then we do not assign a major ancestry to the event.

3.3 Comparisons of MALDER dating using the HAPMAP worldwide and African-specific recombination maps

Recombination maps inferred from different populations are correlated on a broad scale, but differ in the fine-scale characterisation of recombination rates [Hinch et al., 2011]. Having observed significantly older dates in many West and Central West African groups, we investigated the effect of recombination map choice by recomputing MALDER results with an African specific genetic map which was inferred through patterns of LD from the HAPMAP Yoruba (YRI) sample (Fig. 3C).

We next re-inferred admixture parameters with MALDER using all populations with the African (YRI) and additionally with a European (CEU) map [Hinch et al., 2011]. We show comparison of the dates inferred with these different maps in the main paper, and here we shows the equivalent figures to Figure 3 for events inferred using the African (Fig. 3-figure supplement 6) and European (Fig. 3-figure supplement 7) maps.

3.4 Analysis of the minimum genetic distance over which to start curve fitting when using ALDER/MALDER

A key consideration when using weighted LD to infer admixture parameters is the minimum genetic distance over which to begin computing admixture curves. Short-range LD correlations between two reference populations and a target may not only be the result of admixture, but may also be due to demography unrelated to admixture, such as shared recent bottlenecks between the target population and one of the references, or from an extended period of low population size [Loh et al., 2013]. Indeed, the authors of the ALDER algorithm specifically incorporate checks into the default ALDER analysis pipeline that define the threshold at which a test population shares short-range LD with either of the two reference populations. Subsequent curve analyses then ignore data from pairs of SNPs at smaller distances than this correlation threshold [Loh et al., 2013].

The authors nevertheless provide the option of over-riding this LD correlation threshold, allowing the user to define the minimum genetic distance over which the algorithm will begin to compute curves and therefore infer admixture. So there are (at least) two different approaches that can be used to infer admixture using weighted LD. The first is to infer the minimum distance to start building admixture curves from the data (the default), and the second is to assume that any short-range correlations that we observe in the data result from true admixture, and prescribe a minimum distance over which to infer admixture.

3.5 All African populations share correlated LD at short genetic distances

We tested these two approaches by inferring admixture using MALDER/ALDER using a minimum distance defined by the data on the one hand, and a prescribed minimum distance of 0.5cM on the other. This value is commonly used in MALDER analyses, for example by the African Genome Variation Project [Gurdasani et al., 2014]. For each of the 48 African populations as a target, we used ALDER to infer the genetic distance over which LD correlations are shared with every other population as a reference. In Figure 3-figure supplement 3 we show the distribution of these values across all targets for each reference population.

Across all African populations we observe LD correlations with other African populations at genetic distances $> 0.5\text{cM}$, with median values ranging between 0.7cM when GUMUZ is used as a reference to 1.4cM when FULAI is used as a reference. In fact, when we further explore the range of these values across each region separately (Fig. 3-figure supplement 3), we note that, as expected, these distances are greater between more closely related groups.

When we plot the inverse of these distributions, that is, the range of these values for each target population separately, we see that the median minimum distance that all sub-Saharan African populations is always greater than 0.5cM (Fig. 3-figure supplement 4). Taken together, these results suggest that all African populations share some LD over short genetic distances, that may be the result of shared demography or admixture. In order to reduce the confounding effect of demography, with the exception of the next section, all ALDER/MALDER analyses presented in the paper were performed after accounting for this short range shared LD.

3.6 Results of MALDER analysis using a fixed minimum genetic distance of 0.5cM

In order to compare our MALDER analysis to previously published studies, we show the results of the MALDER analysis where we fixed the minimum genetic distance to 0.5cM (Figure 3-figure supplement 5). The main differences between this analysis and that presented in the main part of the paper are:

1. Ancient ($>5\text{ky}$) admixture in Central West African populations where the main analysis found no signal of admixture
2. A second ancient admixture in Malawi c.10ky
3. More of the events appear to involve Eurasian and Khoesan groups mixing.

3.7 Chromosome painting for mixture model and GLOBETROTTER

For the mixture model and GLOBETROTTER analyses, we generated painted samples where we disallowed closely related groups from being painting donors. In practice, this meant removing all populations from the same ancestry region as a given population from the painting analysis. The exception to this are populations from the Nilo-Saharan and Afroasiatic ancestry regions. In these groups, no population from either ancestry region was used as painting donors. We refer to this as the “non-local” painting analysis.

3.8 Modelling populations as mixtures of each other using linear regression

Copying vector summaries generated from painted chromosomes describe how populations relate to one another in terms of the relative time to a common shared ancestor, subsequent recent admixture, and population-specific drift [Hellenthal et al., 2014; Leslie et al., 2015]. For the following analysis, we used the GLOBETROTTER package to generate the mixing coefficients used in Figure 4.

Given a number of potential admixing donor populations, a key step in assessing the extent of admixture in a given population k is to identify which of these donors is relevant; that is, we want to identify the set T^* containing all populations $l \neq k \in [1, \dots, K]$ believed to be involved in any admixture generative to population k . Using copying vectors from the non-local painting analysis, we generate an initial estimate of the mixing coefficients that describe the copying vector of population k by fitting f^k as a mixture of f^l where $l \neq k \in [1, \dots, K]$. The purpose of this step is to assess the evidence for putative admixture in our populations, as described by Hellenthal et al. [2014] and Leslie et al. [2015]. In practice, we remove the self-copying (drift) element from these vectors, i.e. we set $f_i^k = 0$, and rescale each population's copying vector such that $\sum_{i=1}^K f_i^l = 1.0$ for all $l = k \in [1, \dots, K]$.

We assume a standard linear model form for the relationship between f^k and terms f^l for $l \neq k \in [1, \dots, K]$:

$$f^k = \sum_{l \neq k}^K \beta_l^k f^l + \epsilon$$

where ϵ is a vector of errors which we seek to choose the β terms to minimise using non negative least squares regression with the R “nnls” package. Here, β_l^k is the coefficient for f^l under the mixture model, and we estimate the β_l^k s under the constraints that all $\beta_l^k \geq 0$ and $\sum_{l \neq k}^K \beta_l^k = 1.0$. We refer to the estimated coefficient for the l^{th} population as $\hat{\beta}_l^k$; to avoid over-fitting we exclude all populations for which $\hat{\beta}_l^k < 0.001$ and rescale so that $\sum_{l \neq k}^K \hat{\beta}_l^k = 1.0$. T^* is the set of all populations whose $\hat{\beta}_l^k > 0.001$.

The $\hat{\beta}_l^k$ s represent the mixing coefficients that describe a recipient population's DNA as a linear combination of the set T^* donor populations. This process identifies donor populations whose copying vectors match the copying vector of the recipient, as inferred by the painting algorithm.

3.9 Overview of GLOBETROTTER analysis pipeline

In the current setting we are interesting in identifying the general historical relationships between the different African and non-African groups in our dataset. We used GLOBETROTTER [Hellenthal et al., 2014] to characterise patterns of ancestral gene flow and admixture. Individuals tend to share longer stretches of DNA with more closely related individuals, so we used a focused approach where we disallowed copying from local populations.

GLOBETROTTER was originally described by Hellenthal et al. [2014] and a detailed description of the algorithm and the extensive validation of the method is presented in that paper. Here we run over the general framework as used in the current study, with the key difference between our approach and the default use of the algorithm being that we do not allow any groups from within the same ancestry region as a target group to be donors in the painting analysis. Throughout we use GLOBETROTTERv2.

For a given test population k :

1. We define the set of populations present in the same broad ancestry region as k as m , with the caveats outlined below. Using CHROMOPAINTER, we generate painting samples using a reduced set of donors T^* , that included only those populations not present in the same region of Africa as k , i.e. $T^* = l \neq m \in [1, \dots, K]$. The effect of this is to generate mosaic painted chromosomes whose ancestral chunks do not come from closely related individuals or groups, which can mask more subtle signals of admixture. For each group in turn, prior to this final painting, we ran 10 iterations of CHROMOPAINTER's EM algorithm to infer the population-specific prior copying probabilities (using the *-ip* flag), and use these for the final sampled paintings.
2. For each population, l in T^* , we generate a copying-vector, f^l , allowing all individuals from l to copy from every individual in T^* ; i.e. we paint every population l with the same set of restricted donors as k in (1). For each recipient and surrogate group in turn, we sum the chunklengths donated by all individuals within all of our final donor groups (i.e. all 59 groups: including the recipient's own group) and average across all recipients to generate a single 59 element copying vector for each recipient group.
3. To account for noise due to haplotype sharing among groups, we perform a non-negative-least-squares regression (mixture model; outlined above) that takes the copying vector of the recipient group as the response and the copying vectors for each donor group as the predictors. We take the coefficients of this regression, which are restricted to be ≥ 0 and to sum to 1 across donors, as our initial estimates of mixing coefficients describing the genetic make-up of the recipient population as a mixture of other sampled groups.
4. Within and between every pairing of 10 painting samples generated for each haploid of a recipient individual, we consider every pair of chunks (i.e. contiguous segments of DNA copied from a single donor haploid) separated by genetic distance g . For every two donor populations, we tabulate the number of chunk pairs where the two chunks come from the two populations. This is done in a manner to account for phasing switch errors, a common source of error when inferring haplotypes.
5. An appropriate weighting and rescaling of the curves calculated in step 4 gives us the observed coancestry curves illustrating the decay in ancestry linkage disequilibrium versus genetic distance. There is one such curve for each pair of donor populations.
6. We find the maximum likelihood estimate (MLE) of rate parameter λ of an exponential distribution fit to all coancestry curves simultaneously. Specifically, we perform a set of linear regressions that takes each curve in turn as a response and the exponential distribution with parameter λ as a predictor, finding the λ that minimizes the mean-squared residuals of these regressions. This value of λ is our estimated date of admixture. We take the coefficients from each regression. (In the case of 2 dates, we fit two independent exponential distributions with separate rate parameters to all curves simultaneously and take the MLEs of these two rate parameters as our estimates of the two

respective admixture dates. We hence get two sets of coefficients, with each set representing the coefficients for one of the two exponential distributions.)

7. We perform an eigen decomposition of a matrix of values formed using the coefficients inferred in step 6. (In the case of 2 dates, we perform an eigen decomposition of each of the two matrices of coefficients, one for each inferred date.)
8. We use the eigen decomposition from step 7 and the copying vectors to infer both the proportion of admixture α and the mixing coefficients that describe each of the admixing source groups as a linear combination of the donor populations. (In the case of 2 dates, we perform separate fits on each of the two eigen decompositions described in step 7 to describe each admixture event separately.)
9. We re-estimate the mixing coefficients of step 3 to be $\hat{\alpha}$ times the inferred mixing coefficients of the first source plus $1 - \hat{\alpha}$ times the inferred mixing coefficients of the second source.
10. We repeat steps 5-9 for five iterations.
11. We repeat steps 4-5 using a new set of coancestry curves that should eliminate any putative signal of admixture (by taking into account the background distribution of chunks, the so-called **null procedure**), normalize our previous curves using these new ones, and repeat steps 6-10 to re-estimate dates using these normalized curves. We generate 101 date estimates via bootstrapping and assess the proportion of inferred dates that are $= 1$ or ≥ 400 , setting this proportion as our empirical p -value for showing *any* evidence of admixture.
12. Using values calculated in the final iteration of step 10, we classify the admixture event into one of five categories as: (A) 'no admixture', (B) 'uncertain', (C) 'one date', (D) 'multiple dates' and (E) 'one date, multiway'.

3.10 Inferring admixture with GLOBETROTTER

We use the painting samples from (1) and the copying-vectors from (2) detailed in Section 3.9 to implement GLOBETROTTER, characterising admixture in group k ; the intuition being that any admixture observed is likely to be representative of gene-flow from across larger geographic scales.

We report the results of this analysis in Figure 4-Source Data 1 as well as Figures 4, 5, and 6. We generate date estimates by simultaneously fitting an exponential curve to the coancestry curves output by GLOBETROTTER and generate confidence intervals based on 100 bootstrap replicates of the GLOBETROTTER procedure, each time bootstrapping across chromosomes. Because it is unlikely that the true admixing group is present in our set of donor groups, GLOBETROTTER infers the sources of admixture as mixtures of donor groups, which are in some sense equivalent to the β coefficients described above, but are inferred using the additional information present in the coancestry curves. We infer the composition of the admixing sources by using the β s output by GLOBETROTTER from the two (or more) sources of admixture to arrive at an understanding of the genetic basis of the the admixing source groups. These contrasts show us the contribution of each population – which we sum together into regions – to the admixture event and thus provide further intuition into historical gene flow.

3.11 Defining GLOBETROTTER admixture events

The GLOBETROTTER algorithm provides multiple metrics as evidence that admixture has taken place which are combined to arrive at an understanding of the nature of the observed admixture event. In particular, as the authors suggest, to generate an admixture P value, we ran GLOBETROTTER’s “NULL” procedure, which estimates admixture parameters accounting for unusual patterns of LD, and then inferred 100 date bootstraps using this inference, identifying the proportion of inferred dates(s) that are ≤ 1 or ≥ 400 .

Although the algorithm provides a “best-guess” for observed admixture event, we performed the following post-GLOBETROTTER filtering to arrive at our final characterisation of events. We outline the full GLOBETROTTER output in Figure 4-Source Data 2.

1. **Southern African populations** In all Southern African groups we present the results of the GLOBETROTTER runs where results are standardised by using the “NULL” individual see Hellenthal et al. [2014] for further details. We also note that in both the AmaXhosa and SEBantu GLOBETROTTER found evidence for two admixture events but on running the date bootstrap inference process, in both populations the most recent date confidence interval contained 1 generation, suggesting that the dating is not reliable. Inspection of the coancestry curves in this case showed that evidence for a single date of admixture.
2. **East Africa Afroasiatic speaking populations** In all Afroasiatic groups we present the results of the GLOBETROTTER runs where results are standardised by using a “NULL” individual see Hellenthal et al. [2014] for further details.
3. **West African Niger-Congo speaking populations** In all West African Niger-Congo speaking groups, with the exception of the Jola, GLOBETROTTER found evidence for two dates of admixture. In all cases the most recent event was young (1-10 generations) and the date bootstrap confidence interval often contained very small values. Inspection of the coancestry curves showed a sharp decrease at short genetic distances – consistent with the old inferred event – but there was little evidence of a more recent event based on these curves. In groups from this region we therefore show inference of a single date, which we take to be the older of the two dates inferred by GLOBETROTTER.

In all other cases we used the result output by GLOBETROTTER using the default approach.

3.12 Comparison of weighted LD curve dates with GLOBETROTTER dates

Noting that there were differences between the dates inferred by the two dating methods we employed, we compared the dates generated by ALDER/MALDER with those inferred from GLOBETROTTER. Figure 3B shows a comparison of dates using the MALDER event inference; that is, for each population, we used the MALDER inference (either one or two dates) and used the corresponding GLOBETROTTER date inference (either one or two dates) irrespective of whether GLOBETROTTER’s inferred event was different to that of MALDER. Figure 3C is the opposite: we use GLOBETROTTER’s event inference to

define whether we select one or two dates, and then use MALDER’s two date inferences if two dates are inferred, or ALDER’s inference if MALDER infers two dates and GLOBETROTTER infers one. Each point represents a comparison of dates for a single ethnic group, with the symbol and colour defining the ethnic group as in previous plots. In general there appears to be an inflation in dates inferred from West African groups (blues) with MALDER/ALDER compared to GLOBETROTTER.

3.13 Analysis of admixture using sets of restricted surrogates

Recall that for GLOBETROTTER analyses two painting steps are required. One needs to (a) paint target individuals with a set of “donor” individuals to generate mosaic painted chromosomes, and (b) paint all potential “surrogate” groups with the same set of painting donors, such that we then describe admixture in the target individuals with this particular set of surrogate groups. One major benefit of GLOBETROTTER is its ability to represent admixing source groups as mixtures of surrogates.

To infer admixture in sub-Saharan African groups, we painted individuals with a set of painting donors which did not include other individuals from the same ancestry region. We refer to these as “non-local” donors. We then painted individuals from *all* ethnic groups with the same reduced set of painting donors, allowing us to include all ethnic groups as surrogates in the inference process. This allows us to infer sources of admixture for which the components can potentially come from any ethnic group. For example, when inferring admixture in the Jola from the West African Niger-Congo ancestry region, we painted all Jola individuals without including individuals from any other ethnic group from the West African Niger-Congo region. We then painted all ethnic groups with the same set of non-West African Niger-Congo donors, allowing us to model admixture using information from all ethnic groups which have been painted with the same set of painting donors.

The purpose of using this approach was to “mask” the effects of including many closely-related individuals in the painting process: including such individuals in a painting analysis causes the resultant mosaic chromosomes to be dominated by chunklengths from these local groups. As mentioned, this allows us to concentrate on identifying the ancestral processes that have occurred at broader spatio-temporal scales. In effect, removing closely-related individuals from the painting step means that we identify the changing identity of the non-local donor along chromosomes and use these identities to infer historical admixture processes.

3.14 Removing non-local surrogates

In the main analysis we inferred admixture in each of the 48 target sub-Saharan African ethnic groups using all other 47 sub-Saharan African and 12 Eurasian groups as surrogates. We were interested in seeing how the admixture inference changed as we removed surrogate groups from the analysis. Masking surrogates like this provides further insight into the historical relationships between groups. By removing non-local surrogates, we can infer admixture parameters and characterize admixture sources as mixtures of this reduced set of surrogates. Given that, by definition, local groups are more closely related to the target of interest, this approach effectively asks who, outside of the targets region is next best at describing the sources of admixture.

We performed several “restricted surrogate” analyses, for different sets of targets, where we infer admixture using sub-sets of surrogates. One aim of this analysis was to track the spread of Niger-Congo ancestry in the four Niger-Congo ancestry regions. For example, in the full analysis, the major sources of admixture in East African groups tended to be dominated by Southern African Niger-Congo (specifically Malawi) components. If we remove South African Niger-Congo groups from the admixture inference, how is the admixture source now composed?

We performed the following restricted surrogate analyses:

1. **No local region:** for all 48 African groups, we re-ran GLOBETROTTER without allowing any surrogates from the same ancestry region.
2. **No local, east or south:** for groups from the East African and South African Niger-Congo ancestry regions, we disallowed groups from both East and South African Niger-Congo regions from being surrogates. In effect, this asks where in West/Central African is their Niger-Congo ancestry likely to come from.
3. **No local or west:** For West and Central West African groups, we disallowed both West and Central West African Niger-Congo groups from being admixture surrogates. In effect, this asks where in East/South African their ancestry comes from.
4. **No local or Malawi:** As previously noted (Materials and Methods Section 2.5), Malawi was included in the South Africa Niger-Congo ancestry region. There is some evidence, for example from the fineSTRUCTURE analysis, that Malawi is closely related to the East African groups. We therefore wanted to assess whether the inference of a large amount of South African Niger-Congo ancestry in the major sources of admixture in East African Niger-Congo groups was a function of the genetic proximity of Malawi to East Africa. We removed East African Niger-Congo and Malawi as surrogates, and re-inferred admixture parameters.

3.15 Summary of inferred gene-flow from West to East and South Africa

In Figure 4-figure supplement 1 we show the changing composition of GLOBETROTTER’s inferred sources of admixture. Figure 4-figure supplement 2 shows these same results with Niger-Congo surrogates only coloured, and with Malawi and Cameroon groups additionally emphasized. This analysis highlights several key insights:

1. **West and Central West Africa Niger-Congo.** Removing local surrogates from the West African Niger-Congo ancestry region does not substantially change the admixture source inference as admixture in groups from this region already involved Central West African donors. Sources of admixture in Central West African are replaced by West African Niger-Congo components in the non-local analysis. In both regions, admixture source components contain limited ancestry from East and South Africa regions. Interestingly, when we remove all West and Central West African surrogates, we observe similar major sources of admixture across groups from this region, which contains a significant amount of Nilo-Saharan / Afroasiatic ancestry (specifically Sudanese). We

do not have Nilo-Saharan populations from Central and West Africa, although these exist. This observation may represent gene-flow between Nilo-Saharans populations across sub-Saharan Africa.

2. **East African Niger-Congo.** The major sources of admixture inferred in groups from this region tend to contain a majority of Southern African (Malawi) ancestry (Fig. 4-figure supplement 2) whether or not local surrogates are removed from the analysis. When we remove East African Niger-Congo groups and Malawi from the analysis, other South African Niger-Congo speaking groups (SEBantu, Herero, Amaxhosa) are involved in the admixture source. Subsequent removal of these groups leads to an admixture source that is almost exclusively made up of groups from Cameroon, suggesting that Niger-Congo ancestry in East and Southern Africa is most closely related to populations from this part of Central West Africa. Intriguingly, the remainder of major admixture source is made up of a small amount of Khoesan ancestry. Whilst we do not have autochthonous groups from East Africa, this component likely represents ancestry from groups that were present in the area before the major West to East population migrations.
3. **East African Nilo-Saharan and Afroasiatic.** In the full analysis the major source of admixture tends to contain local (purple) groups (Fig. 4-figure supplement 1). When we remove local surrogates from the analysis we see an increase in West and East African donors in the major sources of admixture.
4. **South African Niger-Congo.** We similarly observe mostly East African Niger-Congo ancestry in South African Niger-Congo speakers whether or not we remove non-local surrogates, although there is also a significant Malawi component in the SEBantu and Amaxhosa in the full analysis, which is replaced by East African Niger-Congo groups in the non-local analysis. As in East African Niger-Congo speakers, when we disallow any East or South African Niger-Congo surrogates, the major sources of admixture are almost exclusively from Cameroon donors.
5. **South African Khoesan.** Several Khoesan groups show evidence of specific admixture events involving Central West African donors, in the Khwe, /Gui//Gana, and Xun. When we remove Khoesan groups from being surrogates, interestingly we see that the major sources of admixture now tend to contain ancestry from Southern African Niger-Congo populations, which specifically are not Malawi-like (Fig. 4-figure supplement 2). We also observe a small amount of West African (dark blue) ancestry in these admixture sources, and very little Central West African ancestry (and none from Cameroon). This suggests that the non-Khoesan ancestry in groups from this region is unlikely to be associated with the migrations that spread the Cameroon (Bantu) ancestry into East and South Africa. We also observe a small amount of East African Nilo-Saharan / Afroasiatic ancestry in the major sources of admixture in the non-local analysis.

A key insight from these analyses is that a large amount of the ancestry in East and Southern Africa Niger-Congo speaking groups comes from populations currently present in Central West Africa. The opposite is not true for Western and Central West African groups, who share DNA amongst themselves, but when local surrogates are removed, they instead choose mainly non-Niger-Congo speaking groups as admixture sources.

4 Plotting date densities

For each admixture event we split the admixture sources into their constituent components (i.e. we used the β coefficients inferred by GLOBETROTTER) at the appropriate admixture proportions. For a given event, these components sum to 1. We multiplied these components by 100 to estimate the percentage of ancestry from a given event that originates from each donor group. We then assigned each of the components the set of date bootstraps associated with the event. For example, in the Kauma we infer an admixture event with an admixture proportion α of 6% involving a minor source containing the following coefficients: Massai 0.02 Afar 0.15 GBR 0.26 GIH 0.57. We multiply each of these coefficients by α to obtain a final proportion that each group gives to the admixture event: Massai 0.5 Afar 1 GBR 1.5 GIH 3. We assign all the inferred date bootstraps for the Kauma to each of the populations in these proportions. In this example, GIH has twice the density of GBR. We then additionally sum components across the same country to finally arrive at the density plots in Figures 5.

5 Gene-flow maps

We generated maps with the *rworldmaps* package in R. To generate arrows, we combined the inferred ancestral components (i.e. 1ST and 2ND EVENT SOURCES in Figure 3) for each population and estimated the proportion of a group's ancestry coming from each component, summed across all surrogates from a particular country. For example, if an admixture contained source contains components from both the Jola and Wolof (both from The Gambia), then these components were added together. As such, the arrows point from the country of component origin to the country of the recipient. We then plot only those arrows which relate to events pertaining to the different broad gene-flow events. For each map, we used plot arrows for any event involving the following:

- (A) **Recent Western Bantu gene-flow**: any admixture source which has a component from either of the two Cameroon ethnic groups, Bantu and Semi-Bantu.
- (B) **Eastern Bantu gene-flow**: any admixture source which has a component from Kenya, Tanzania, Malawi, or South Africa (Niger-Congo speakers).
- (C) **East / West gene-flow**: any admixture event which has a component from Gambia, Burkina Faso, Ghana, Mali, Nigeria, Ethiopia, Sudan or Somalia.
- (D) **Eurasian gene-flow into Africa**: any admixture event which has a component from any Eurasian population.

An alternative map stratified by time window, rather than admixture component is shown in Figure 6-figure supplement 1.

6 Analysis and plotting code

Code used for analyses and plotting will be made available at <https://github.com/georgebusby/popgen>

Acknowledgements

We thank all the MalariaGEN study sites that contributed samples to this analysis: a list of researchers involved at each study site can be found at <https://www.malariagen.net/projects/host/consortium-members>.

MalariaGEN is funded by the Wellcome Trust (WT077383/Z/05/Z, 090770/Z/09/Z) and the Bill and Melinda Gates Foundation through the Foundation for the National Institutes of Health (566). Genotyping was performed at the Wellcome Trust Sanger Institute, partly funded by its core award from the Wellcome Trust (098051/Z/05/Z). This research was also supported by Centre grants from the Wellcome Trust (090532/Z/09/Z) and the Medical Research Council (G0600718). C.C.A.S. was supported by a Wellcome Trust Career Development Fellowship (097364/Z/11/Z).

The Malaria Research and Training Center–Bandiagara Malaria Project (MRTC-BMP) in Mali group is supported by an Interagency Committee on Disability Research (ICDR) grant from the National Institute of Allergy and Infectious Diseases/US National Institutes of Health (NIAID/NIH) to the University of Maryland and the University of Bamako (USTTB) and by the Mali-NIAID/NIH International Centers for Excellence in Research (ICER) at USTTB. The Kenya Medical Research Institute (KEMRI)–Wellcome Trust Programme is funded through core support from the Wellcome Trust. This paper is published with the permission of the director of KEMRI. C.M.N. is supported through a strategic award to the KEMRI–Wellcome Trust Programme from the Wellcome Trust (084538). The Joint Malaria Programme, Kilimanjaro Christian Medical Centre in Tanzania received funding from a UK MRC grant (G9901439).

We thank Clare Bycroft, Lucy van Dorp, and Cristian Capelli for critically evaluating the manuscript and Francesco Montinaro for insightful discussions on interpretation of MALDER analyses. Genotype data for the MalariaGen samples included in this paper will be made available at the European Nucleotide Archive (accession number TBC).

References

- Allen, J. D. V. (1993). *Swahili Origins: Swahili Culture & the Shungwaya Phenomenon*. James Currey Publishers.
- Ansari Pour, N., Plaster, C. A., and Bradman, N. (2013). Evidence from Y-chromosome analysis for a late exclusively eastern expansion of the Bantu-speaking people. *European Journal of Human Genetics*, 21(4):423–429.
- Band, G., Le, Q. S., Jostins, L., Pirinen, M., Kivinen, K., Jallow, M., Sisay-Joof, F., Bojang, K., Pinder, M., Sirugo, G., Conway, D. J., Nyirongo, V., Kachala, D., Molyneux, M., Taylor, T., Ndila, C., Peshu, N., Marsh, K., Williams, T. N., Alcock, D., Andrews, R., Edkins, S., Gray, E., Hubbart, C., Jeffreys, A., Rowlands, K., Schuldt, K., Clark, T. G., Small, K. S., Teo, Y. Y., Kwiatkowski, D. P., Rockett, K. A., Barrett, J. C., Spencer, C. C. A., and Malaria Genomic Epidemiological Network ¶ (2013). Imputation-Based Meta-Analysis of Severe Malaria in Three African Populations. *PLoS Genet*, 9(5):e1003509.
- Barham, L. and Mitchell, P. (2008). *The First Africans: African Archaeology from the Earliest Toolmakers to Most Recent Foragers*. Cambridge University Press, Cambridge ; New York, 1 edition edition.
- Behar, D. M., Yunusbayev, B., Metspalu, M., Metspalu, E., Rosset, S., Parik, J., Rootsi, S., Chaubey, G., Kutuev, I., Yudkovsky, G., and others (2010). The genome-wide structure of the Jewish people. *Nature*, 466(7303):238–242.
- Beleza, S., Gusmão, L., Amorim, A., Carracedo, A., and Salas, A. (2005). The genetic legacy of western Bantu migrations. *Human Genetics*, 117(4):366–375.
- Berniell-Lee, G., Calafell, F., Bosch, E., Heyer, E., Sica, L., Mouguiama-Daouda, P., van der Veen, L., Hombert, J.-M., Quintana-Murci, L., and Comas, D. (2009). Genetic and Demographic Implications of the Bantu Expansion: Insights from Human Paternal Lineages. *Molecular Biology and Evolution*, 26(7):1581–1589.
- Bhatia, G., Patterson, N., Sankararaman, S., and Price, A. L. (2013). Estimating and interpreting FST: The impact of rare variants. *Genome Research*, 23(9):1514–1521.
- Blench, R. (2006). *Archaeology-Language-and-the-African-Past*.
- Breton, G., Schlebusch, C. M., Lombard, M., Sjödin, P., Soodyall, H., and Jakobsson, M. (2014). Lactase Persistence Alleles Reveal Partial East African Ancestry of Southern African Khoe Pastoralists. *Current Biology*, 24(8):852–858.
- Bryc, K., Auton, A., Nelson, M. R., Oksenberg, J. R., Hauser, S. L., Williams, S., Froment, A., Bodo, J.-M., Wambebe, C., Tishkoff, S. A., and Bustamante, C. D. (2010). Genome-wide patterns of population structure and admixture in West Africans and African Americans. *Proceedings of the National Academy of Sciences*, 107(2):786–791.

- Busby, G. B. J., Hellenthal, G., Montinaro, F., Tofanelli, S., Bulayeva, K., Rudan, I., Zemunik, T., Hayward, C., Toncheva, D., Karachanak-Yankova, S., Nesheva, D., Anagnostou, P., Cali, F., Brisighelli, F., Romano, V., Lefranc, G., Buresi, C., Ben Chibani, J., Haj-Khelil, A., Denden, S., Ploski, R., Krajewski, P., Hervig, T., Moen, T., Herrera, R. J., Wilson, J. F., Myers, S., and Capelli, C. (2015). The Role of Recent Admixture in Forming the Contemporary West Eurasian Genomic Landscape. *Current Biology*, 25(19):2518–2526.
- Campbell, M. C. and Tishkoff, S. A. (2008). African Genetic Diversity: Implications for Human Demographic History, Modern Human Origins, and Complex Disease Mapping. *Annual Review of Genomics and Human Genetics*, 9(1):403–433.
- Consortium, T. . G. P. (2012). An integrated map of genetic variation from 1,092 human genomes. *Nature*, 491(7422):56–65.
- Coop, G., Pickrell, J., Novembre, J., Kudaravalli, S., Li, J., Absher, D., Myers, R., Cavalli-Sforza, L., Feldman, M., and Pritchard, J. (2009). The role of geography in human adaptation. *PLoS Genetics*, 5(6).
- Currat, M., Trabuchet, G., Rees, D., Perrin, P., Harding, R. M., Clegg, J. B., Langaney, A., and Excoffier, L. (2002). Molecular Analysis of the β -Globin Gene Cluster in the Niokholo Mandenka Population Reveals a Recent Origin of the β S Senegal Mutation. *The American Journal of Human Genetics*, 70(1):207–223.
- Currie, T. E., Meade, A., Guillon, M., and Mace, R. (2013). Cultural phylogeography of the Bantu Languages of sub-Saharan Africa. *Proceedings of the Royal Society of London B: Biological Sciences*, 280(1762):20130695.
- Dannemann, M., Andrés, A. M., and Kelso, J. (2016). Introgression of Neandertal- and Denisovan-like Haplotypes Contributes to Adaptive Variation in Human Toll-like Receptors. *The American Journal of Human Genetics*, 98(1):22–33.
- de Filippo, C., Bostoen, K., Stoneking, M., and Pakendorf, B. (2012). Bringing together linguistic and genetic evidence to test the Bantu expansion. *Proceedings of the Royal Society B: Biological Sciences*, 279(1741):3256–3263.
- Delaneau, O., Marchini, J., and Zagury, J.-F. (2012). A linear complexity phasing method for thousands of genomes. *Nature Methods*, 9(2):179–181.
- Deschamps, M., Laval, G., Fagny, M., Itan, Y., Abel, L., Casanova, J.-L., Patin, E., and Quintana-Murci, L. (2016). Genomic Signatures of Selective Pressures and Introgression from Archaic Hominins at Human Innate Immunity Genes. *The American Journal of Human Genetics*, 98(1):5–21.
- Diamond, J. and Bellwood, P. (2003). Farmers and their languages: The first expansions. *Science*, 300(5619):597–603.

- Ehret, C. (1967). Cattle-Keeping and Milking in Eastern and Southern African History: The Linguistic Evidence. *The Journal of African History*, 8(1):1–17.
- Fenner, J. (2005). Cross-cultural estimation of the human generation interval for use in genetics-based population divergence studies. *American Journal of Physical Anthropology*, 128(2):415–423.
- Fraley, C., Raftery, A. E., Murphy, T. B., and Scrucca, L. (2012). *mclust Version 4 for R: Normal Mixture Modeling for Model-Based Clustering, Classification, and Density Estimation*. Number 597. Department of Statistics, University of Washington.
- Fumagalli, M., Moltke, I., Grarup, N., Racimo, F., Bjerregaard, P., Jørgensen, M. E., Korneliussen, T. S., Gerbault, P., Skotte, L., Linneberg, A., Christensen, C., Brandslund, I., Jørgensen, T., Huerta-Sánchez, E., Schmidt, E. B., Pedersen, O., Hansen, T., Albrechtsen, A., and Nielsen, R. (2015). Greenlandic Inuit show genetic signatures of diet and climate adaptation. *Science*, 349(6254):1343–1347.
- González-Santos, M., Montinaro, F., Oosthuizen, O., Oosthuizen, E., Busby, G. B. J., Anagnostou, P., Destro-Bisol, G., Pascali, V., and Capelli, C. (2015). Genome-Wide SNP Analysis of Southern African Populations Provides New Insights into the Dispersal of Bantu-Speaking Groups. *Genome Biology and Evolution*, 7(9):2560–2568.
- Greenberg, J. H. (1972). Linguistic evidence regarding Bantu origins. *The Journal of African History*, 13(02):189–216.
- Grollemund, R., Branford, S., Bostoen, K., Meade, A., Venditti, C., and Pagel, M. (2015). Bantu expansion shows that habitat alters the route and pace of human dispersals. *Proceedings of the National Academy of Sciences*, 112(43):13296–13301.
- Guldemann, T. (2008). A linguist’s view : Khoe-Kwadi speakers as the earliest food-producers of southern Africa.
- Guldemann, T. and Fehn, A.-M., editors (2014). *Beyond 'Khoisan': Historical relations in the Kalahari Basin*, volume 330 of *Current Issues in Linguistic Theory*. John Benjamins Publishing Company, Amsterdam.
- Gurdasani, D., Carstensen, T., Tekola-Ayele, F., Pagani, L., Tachmazidou, I., Hatzikotoulas, K., Karthikeyan, S., Iles, L., Pollard, M. O., Choudhury, A., Ritchie, G. R. S., Xue, Y., Asimit, J., Nsubuga, R. N., Young, E. H., Pomilla, C., Kivinen, K., Rockett, K., Kamali, A., Doumatey, A. P., Asiki, G., Seeley, J., Sisay-Joof, F., Jallow, M., Tollman, S., Mekonnen, E., Ekong, R., Oljira, T., Bradman, N., Bojang, K., Ramsay, M., Adeyemo, A., Bekele, E., Motala, A., Norris, S. A., Pirie, F., Kaleebu, P., Kwiatkowski, D., Tyler-Smith, C., Rotimi, C., Zeggini, E., and Sandhu, M. S. (2014). The African Genome Variation Project shapes medical genetics in Africa. *Nature*, advance online publication.
- H3Africa Consortium (2014). Research capacity. Enabling the genomic revolution in Africa. *Science (New York, N.Y.)*, 344(6190):1346–1348.

- Hamblin, M. T. and Di Rienzo, A. (2000). Detection of the signature of natural selection in humans: evidence from the Duffy blood group locus. *American Journal of Human Genetics*, 66(5):1669–1679.
- Hamblin, M. T., Thompson, E. E., and Di Rienzo, A. (2002). Complex signatures of natural selection at the duffy blood group locus. *American Journal of Human Genetics*, 70(2):369–383.
- Hedrick, P. W. (2011). Population genetics of malaria resistance in humans. *Heredity*, 107(4):283–304.
- Hellenthal, G., Busby, G. B. J., Band, G., Wilson, J. F., Capelli, C., Falush, D., and Myers, S. (2014). A Genetic Atlas of Human Admixture History. *Science*, 343(6172):747–751.
- Henn, B. M., Botigué, L. R., Gravel, S., Wang, W., Brisbin, A., Byrnes, J. K., Fadhlouzi-Zid, K., Zalloua, P. A., Moreno-Estrada, A., Bertranpetit, J., Bustamante, C. D., and Comas, D. (2012). Genomic Ancestry of North Africans Supports Back-to-Africa Migrations. *PLoS Genet*, 8(1):e1002397.
- Henn, B. M., Gignoux, C., Lin, A. A., Oefner, P. J., Shen, P., Scozzari, R., Cruciani, F., Tishkoff, S. A., Mountain, J. L., and Underhill, P. A. (2008). Y-chromosomal evidence of a pastoralist migration through Tanzania to southern Africa. *Proceedings of the National Academy of Sciences*, 105(31):10693.
- Hinch, A. G., Tandon, A., Patterson, N., Song, Y., Rohland, N., Palmer, C. D., Chen, G. K., Wang, K., Buxbaum, S. G., Akyzbekova, E. L., Aldrich, M. C., Ambrosone, C. B., Amos, C., Bandera, E. V., Berndt, S. I., Bernstein, L., Blot, W. J., Bock, C. H., Boerwinkle, E., Cai, Q., Caporaso, N., Casey, G., Adrienne Cupples, L., Deming, S. L., Ryan Diver, W., Divers, J., Fornage, M., Gillanders, E. M., Glessner, J., Harris, C. C., Hu, J. J., Ingles, S. A., Isaacs, W., John, E. M., Linda Kao, W. H., Keating, B., Kittles, R. A., Kolonel, L. N., Larkin, E., Le Marchand, L., McNeill, L. H., Millikan, R. C., Murphy, Musani, S., Neslund-Dudas, C., Nyante, S., Papanicolaou, G. J., Press, M. F., Psaty, B. M., Reiner, A. P., Rich, S. S., Rodriguez-Gil, J. L., Rotter, J. I., Rybicki, B. A., Schwartz, A. G., Signorello, L. B., Spitz, M., Strom, S. S., Thun, M. J., Tucker, M. A., Wang, Z., Wiencke, J. K., Witte, J. S., Wrensch, M., Wu, X., Yamamura, Y., Zanetti, K. A., Zheng, W., Ziegler, R. G., Zhu, X., Redline, S., Hirschhorn, J. N., Henderson, B. E., Taylor Jr, H. A., Price, A. L., Hakonarson, H., Chanock, S. J., Haiman, C. A., Wilson, J. G., Reich, D., and Myers, S. R. (2011). The landscape of recombination in African Americans. *Nature*, 476(7359):170–175.
- Hodgson, J. A., Mulligan, C. J., Al-Meer, A., and Raaum, R. L. (2014a). Early Back-to-Africa Migration into the Horn of Africa. *PLoS Genet*, 10(6):e1004393.
- Hodgson, J. A., Pickrell, J. K., Pearson, L. N., Quillen, E. E., Prista, A., Rocha, J., Soodyall, H., Shriver, M. D., and Perry, G. H. (2014b). Natural selection for the Duffy-null allele in the recently admixed people of Madagascar. *Proceedings of the Royal Society B: Biological Sciences*, 281(1789):20140930.
- Holden, C. J. (2002). Bantu language trees reflect the spread of farming across sub-Saharan Africa: a maximum-parsimony analysis. *Proceedings of the Royal Society of London B: Biological Sciences*, 269(1493):793–799.

- Hudson, R., Slatkin, M., and Maddison, W. (1992). Estimation of levels of gene flow from DNA sequence data. *Genetics*, 132(2):583–589.
- Huerta-Sánchez, E., Jin, X., Asan, Bianba, Z., Peter, B. M., Vinckenbosch, N., Liang, Y., Yi, X., He, M., Somel, M., Ni, P., Wang, B., Ou, X., Huasang, Luosang, J., Cuo, Z. X. P., Li, K., Gao, G., Yin, Y., Wang, W., Zhang, X., Xu, X., Yang, H., Li, Y., Wang, J., Wang, J., and Nielsen, R. (2014). Altitude adaptation in Tibetans caused by introgression of Denisovan-like DNA. *Nature*, 512(7513):194–197.
- Jeong, C., Alkorta-Aranburu, G., Basnyat, B., Neupane, M., Witonsky, D. B., Pritchard, J. K., Beall, C. M., and Di Rienzo, A. (2014). Admixture facilitates genetic adaptations to high altitude in Tibet. *Nature Communications*, 5.
- Lawson, D. J., Hellenthal, G., Myers, S., and Falush, D. (2012). Inference of Population Structure using Dense Haplotype Data. *PLoS Genetics*, 8(1):e1002453.
- Leslie, S., Winney, B., Hellenthal, G., Davison, D., Boumertit, A., Day, T., Hutnik, K., Royrvik, E. C., Cunliffe, B., Wellcome Trust Case Control Consortium 2, International Multiple Sclerosis Genetics Consortium, Lawson, D. J., Falush, D., Freeman, C., Pirinen, M., Myers, S., Robinson, M., Donnelly, P., and Bodmer, W. (2015). The fine-scale genetic structure of the British population. *Nature*, 519(7543):309–314.
- Li, N. and Stephens, M. (2003). Modeling linkage disequilibrium and identifying recombination hotspots using single-nucleotide polymorphism data. *Genetics*, 165(4):2213–2233.
- Li, S., Schlebusch, C., and Jakobsson, M. (2014). Genetic variation reveals large-scale population expansion and migration during the expansion of Bantu-speaking peoples. *Proceedings of the Royal Society B: Biological Sciences*, 281(1793):20141448.
- Llorente, M. G., Jones, E. R., Eriksson, A., Siska, V., Arthur, K. W., Arthur, J. W., Curtis, M. C., Stock, J. T., Coltorti, M., Pieruccini, P., Stretton, S., Brock, F., Higham, T., Park, Y., Hofreiter, M., Bradley, D. G., Bhak, J., Pinhasi, R., and Manica, A. (2015). Ancient Ethiopian genome reveals extensive Eurasian admixture throughout the African continent. *Science*, page aad2879.
- Loh, P.-R., Lipson, M., Patterson, N., Moorjani, P., Pickrell, J. K., Reich, D., and Berger, B. (2013). Inferring Admixture Histories of Human Populations Using Linkage Disequilibrium. *Genetics*, 193(4):1233–1254.
- Macholdt, E., Slatkin, M., Pakendorf, B., and Stoneking, M. (2015). New insights into the history of the C-14010 lactase persistence variant in Eastern and Southern Africa. *American Journal of Physical Anthropology*, 156(4):661–664.
- Malaria Genomic Epidemiology Network (2008). A global network for investigating the genomic epidemiology of malaria. *Nature*, 456(7223):732–737.
- Malaria Genomic Epidemiology Network (2014). Reappraisal of known malaria resistance loci in a large multicenter study. *Nature Genetics*, 46(11):1197–1204.

- Malaria Genomic Epidemiology Network (2015). A novel locus of resistance to severe malaria in a region of ancient balancing selection. *Nature*, 526(7572):253–257.
- Marks, S. J., Montinaro, F., Levy, H., Brisighelli, F., Ferri, G., Bertoncini, S., Batini, C., Busby, G. B., Arthur, C., Mitchell, P., Stewart, B. A., Oosthuizen, O., Oosthuizen, E., D’Amato, M. E., Davison, S., Pascali, V., and Capelli, C. (2014). Static and moving frontiers: the genetic landscape of Southern African Bantu-speaking populations. *Molecular Biology and Evolution*, page msu263.
- Mitchell, P. (2002). *The archaeology of Southern Africa*. Cambridge University Press, Cambridge, UK.
- Modiano, D., Bancone, G., Ciminelli, B. M., Pompei, F., Blot, I., Simporé, J., and Modiano, G. (2008). Haemoglobin S and haemoglobin C: ‘quick but costly’ versus ‘slow but gratis’ genetic adaptations to *Plasmodium falciparum* malaria. *Human Molecular Genetics*, 17(6):789–799.
- Montinaro, F., Busby, G. B. J., Pascali, V. L., Myers, S., Hellenthal, G., and Capelli, C. (2015). Unravelling the hidden ancestry of American admixed populations. *Nature Communications*, 6.
- Need, A. C. and Goldstein, D. B. (2009). Next generation disparities in human genomics: concerns and remedies. *Trends in Genetics*, 25(11):489–494.
- Novembre, J., Johnson, T., Bryc, K., Kutalik, Z., Boyko, A., Auton, A., Indap, A., King, K., Bergmann, S., Nelson, M., Stephens, M., and Bustamante, C. (2008). Genes mirror geography within Europe. *Nature*, 456(7218):98–101.
- Nurse, D. and Philippson, G. (2003). *The Bantu Languages*. Number 4 in Routledge Language Family Series. Routledge, London, UK.
- Oliver, R. A. and Fagan, B. M. (1975). *Africa in the Iron Age: C.500 BC-1400 AD*. Cambridge University Press.
- Pagani, L., Kivisild, T., Tarekegn, A., Ekong, R., Plaster, C., Gallego Romero, I., Ayub, Q., Mehdi, S. Q., Thomas, M. G., Luiselli, D., Bekele, E., Bradman, N., Balding, D. J., and Tyler-Smith, C. (2012). Ethiopian Genetic Diversity Reveals Linguistic Stratification and Complex Influences on the Ethiopian Gene Pool. *The American Journal of Human Genetics*, 91(1):83–96.
- Pakendorf, B., Bostoen, K., and de Filippo, C. (2011). Molecular Perspectives on the Bantu Expansion: A Synthesis. *Language Dynamics and Change*, 1(1):50–88.
- Patterson, N., Moorjani, P., Luo, Y., Mallick, S., Rohland, N., Zhan, Y., Genschoreck, T., Webster, T., and Reich, D. (2012). Ancient Admixture in Human History. *Genetics*, 192(3):1065–1093.
- Patterson, N., Price, A. L., and Reich, D. (2006). Population Structure and Eigenanalysis. *PLoS Genetics*, 2(12):e190.
- Petersen, D. C., Libiger, O., Tindall, E. A., Hardie, R.-A., Hannick, L. I., Glashoff, R. H., Mukerji, M., Fernandez, P., Haacke, W., Schork, N. J., Hayes, V. M., and Indian Genome Variation Consortium (2013). Complex Patterns of Genomic Admixture within Southern Africa. *PLoS Genet*, 9(3):e1003309.

- Pickrell, J. K., Patterson, N., Barbieri, C., Berthold, F., Gerlach, L., Güldemann, T., Kure, B., Mpoloka, S. W., Nakagawa, H., Naumann, C., Lipson, M., Loh, P.-R., Lachance, J., Mountain, J., Bustamante, C. D., Berger, B., Tishkoff, S. A., Henn, B. M., Stoneking, M., Reich, D., and Pakendorf, B. (2012). The genetic prehistory of southern Africa. *Nature Communications*, 3:1143.
- Pickrell, J. K., Patterson, N., Loh, P.-R., Lipson, M., Berger, B., Stoneking, M., Pakendorf, B., and Reich, D. (2014). Ancient west eurasian ancestry in southern and eastern africa. *Proceedings of the National Academy of Sciences of the United States of America*, 111(7):2632–2637.
- Pickrell, J. K. and Pritchard, J. K. (2012). Inference of Population Splits and Mixtures from Genome-Wide Allele Frequency Data. *PLoS Genetics*, 8(11):e1002967.
- Pickrell, J. K. and Reich, D. (2014). Toward a new history and geography of human genes informed by ancient DNA. *Trends in Genetics*, 30(9):377–389.
- R Development Core Team (2011). R: a language and environment for statistical computing. Vienna (Austria): R Foundation for Statistical Computing.
- Ralph, P. and Coop, G. (2013). The Geography of Recent Genetic Ancestry across Europe. *PLoS Biol*, 11(5):e1001555.
- Ranciaro, A., Campbell, M. C., Hirbo, J. B., Ko, W.-Y., Froment, A., Anagnostou, P., Kotze, M. J., Ibrahim, M., Nyambo, T., Omar, S. A., and Tishkoff, S. A. (2014). Genetic Origins of Lactase Persistence and the Spread of Pastoralism in Africa. *The American Journal of Human Genetics*, 94(4):496–510.
- Reich, D., Thangaraj, K., Patterson, N., Price, A. L., and Singh, L. (2009). Reconstructing Indian population history. *Nature*, 461(7263):489–494.
- Roberts, J. (2007). *The New Penguin History of the World*. Penguin Books, London, UK, 5th edition.
- Schlebusch, C. M., Skoglund, P., Sjödin, P., Gattepaille, L. M., Hernandez, D., Jay, F., Li, S., Jongh, M. D., Singleton, A., Blum, M. G. B., Soodyall, H., and Jakobsson, M. (2012). Genomic Variation in Seven Khoe-San Groups Reveals Adaptation and Complex African History. *Science*, 338(6105):374–379.
- Smith, A. B. (2005). *African Herders: Emergence of Pastoral Traditions*. Rowman Altamira.
- The International HapMap Consortium (2010). Integrating common and rare genetic variation in diverse human populations. *Nature*, 467(7311):52–58.
- Thompson, L. (2001). *A history of South Africa*. Yale University Press, USA, 3rd edition.
- Tishkoff, S. A., Reed, F. A., Friedlaender, F. R., Ehret, C., Ranciaro, A., Froment, A., Hirbo, J. B., Awomoyi, A. A., Bodo, J.-M., Doumbo, O., Ibrahim, M., Juma, A. T., Kotze, M. J., Lema, G., Moore, J. H., Mortensen, H., Nyambo, T. B., Omar, S. A., Powell, K., Pretorius, G. S., Smith, M. W., Thera,

- M. A., Wambebe, C., Weber, J. L., and Williams, S. M. (2009). The Genetic Structure and History of Africans and African Americans. *Science*, 324(5930):1035–1044.
- Tishkoff, S. A., Reed, F. A., Ranciaro, A., Voight, B. F., Babbitt, C. C., Silverman, J. S., Powell, K., Mortensen, H. M., Hirbo, J. B., Osman, M., Ibrahim, M., Omar, S. A., Lema, G., Nyambo, T. B., Gori, J., Bumpstead, S., Pritchard, J. K., Wray, G. A., and Deloukas, P. (2007). Convergent adaptation of human lactase persistence in Africa and Europe. *Nature Genetics*, 39(1):31–40.
- van Dorp, L., Balding, D., Myers, S., Pagani, L., Tyler-Smith, C., Bekele, E., Tarekegn, A., Thomas, M. G., Bradman, N., and Hellenthal, G. (2015). Evidence for a Common Origin of Blacksmiths and Cultivators in the Ethiopian Ari within the Last 4500 Years: Lessons for Clustering-Based Inference. *PLoS Genet*, 11(8):e1005397.
- Welter, D., MacArthur, J., Morales, J., Burdett, T., Hall, P., Junkins, H., Klemm, A., Flicek, P., Manolio, T., Hindorff, L., and Parkinson, H. (2014). The NHGRI GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Research*, 42(Database issue):D1001–D1006.
- Zheng, X., Levine, D., Shen, J., Gogarten, S. M., Laurie, C., and Weir, B. S. (2012). A high-performance computing toolset for relatedness and principal component analysis of SNP data. *Bioinformatics*, 28(24):3326–3328.

Figures and Figure Supplements

Figure 1:	Sub-Saharan African genetic variation is shaped by ethno-linguistic and geographical similarity.	5
Figure 2:	Haplotypes capture more population structure than independent loci	7
Figure 3:	Inference of admixture in sub-Saharan Africa using MALDER . . .	11
Figure 4:	Inference of admixture in sub-Saharan African using GLOBETROTTER	13
Figure 5:	A timeline of recent admixture in sub-Saharan Africa.	16
Figure 6:	The geography of recent gene-flow in Africa	19
Figure 1-figure supplement 1:	Map of populations used in the analysis.	54
Figure 1-figure supplement 2:	An example of hierarchical clustering to chose two groups of similar individuals from the Fula based on a PCA of the Gambia.	55
Figure 1-figure supplement 3:	fineSTRUCTURE analysis of the full dataset	56
Figure 3-figure supplement 1:	Weighted LD amplitudes for a selection of 9 ethnic groups	57
Figure 3-figure supplement 2:	Comparison of weighted LD amplitude scores across all African ethnic groups	58
Figure 3-figure supplement 3:	Comparison of the minimum distance to begin computing admixture LD	59
Figure 3-figure supplement 4:	Comparison of the minimum distance to begin computing admixture LD split by region	60
Figure 3-figure supplement 5:	Results of the MALDER analysis computing weighted admixture decay curves from 0.5cM	61
Figure 3-figure supplement 6:	Results of MALDER for all populations using an African specific recombination map	62
Figure 3-figure supplement 7:	Results of MALDER for all populations using a European specific recombination map	63
Figure 4-figure supplement 1:	Admixture source inference by GLOBETROTTER after sequentially removing local surrogates from the analysis	64
Figure 4-figure supplement 2:	Admixture source inference by GLOBETROTTER after sequentially removing local surrogates from the analysis	65
Figure 6-figure supplement 1:	Gene-flow in Africa over the last 2,000 years.	66

Supplementary Files

Supplementary File 1 A note on ethnolinguistic groupings	67
--	----

Source Data Tables

Figure 1-Source Data 1 Overview of sampled populations	69
Figure 2-Source Data 1 Pairwise F_{ST} for African populations	70
Figure 2-Source Data 2 Pairwise F_{ST} for Eurasian populations	71
Figure 2-Source Data 3 Pairwise TVD for African populations	72
Figure 2-Source Data 4 Pairwise TVD for Eurasian populations	73
Figure 3-Source Data 1 The evidence for multiple waves of admixture in African populations using MALDER and the HAPMAP recombination map.	74
Figure 3-Source Data 2 The evidence for multiple waves of admixture in African populations using MALDER and the African recombination map.	75
Figure 3-Source Data 3 The evidence for multiple waves of admixture in African populations using MALDER and the European recombination map.	76
Figure 3-Source Data 4 The evidence for multiple waves of admixture in African populations using MALDER and the HAPMAP recombination map and a mindis of 0.5cM.	77
Figure 4-Source Data 1 Results of the main GLOBETROTTER analysis.	79
Figure 4-Source Data 2 Results of the main GLOBETROTTER analysis.	81
Supplementary Table 1 Evidence for admixture across the ethnic groups included in the study using f_3 tests and ALDER.	85

Figure Supplements

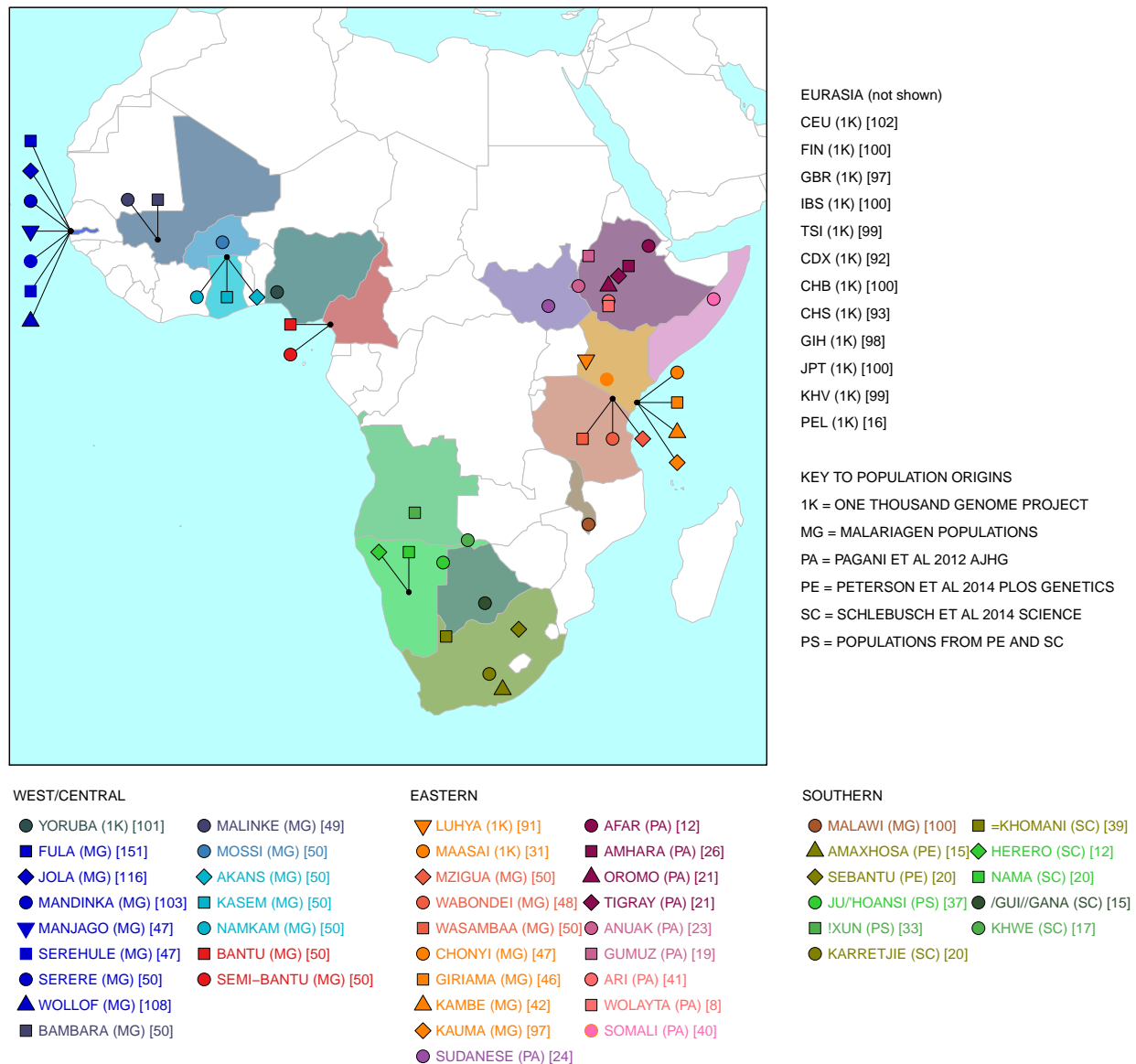


Figure 1-figure supplement 1. Map of populations used in the analysis. Population names are coloured by the country of origin; positions of the countries are shown on the map. Individual point labels, which are used throughout this paper, are shown for each population in the legend. Sample provenance is shown immediately after the population name in circular parentheses and final number of individuals is shown in square parentheses.

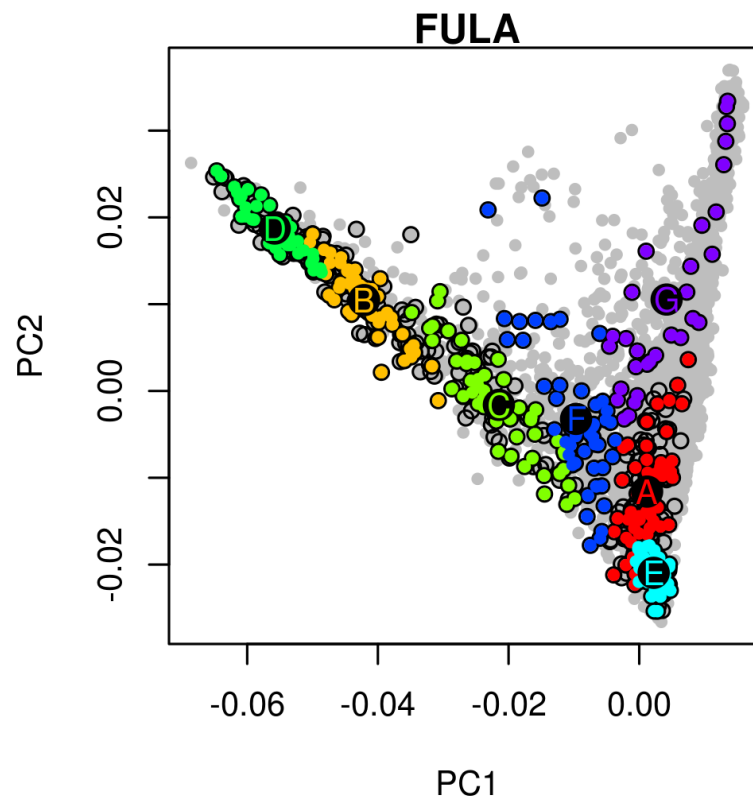


Figure 1-figure supplement 2. An example of hierarchical clustering to chose two groups of similar individuals from the Fula based on a PCA of the Gambia. Projected onto a PCA of Gambian genetic variation where each point represents an individual, all Fula individuals are coloured, with the colour depicting their cluster assignment, based on the *MClust* clustering algorithm. We chose individuals from the green (D) and light blue (E) clusters to maximise the representation of Fula genetic variation. Note that the majority of the individuals from the other 6 Gambian ethnic groups occur in the right arm of the PCA. An analogous process was preformed for all ethnic groups from the MalariaGEN dataset where more than 50 individuals were available.

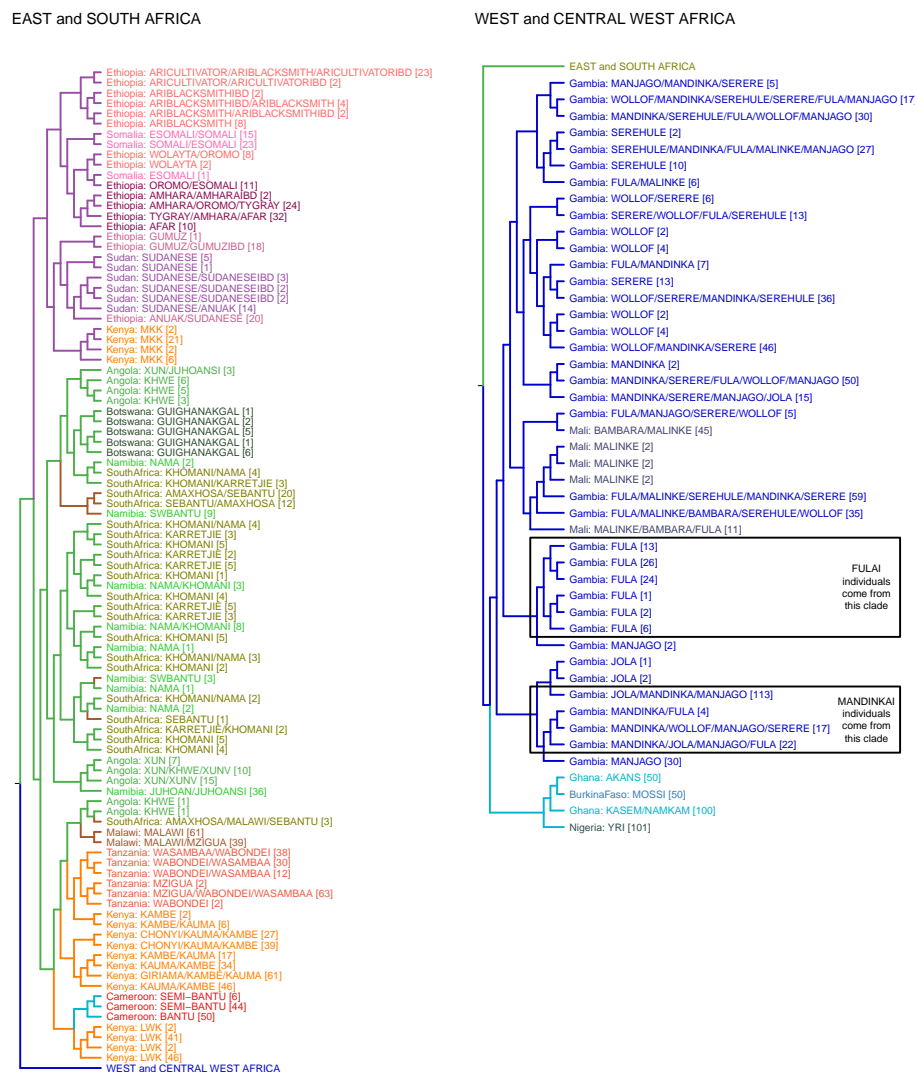


Figure 1-figure supplement 3. fineSTRUCTURE analysis of the full dataset. We show the tree output from a single run of the fineSTRUCTURE algorithm. To aid reading, the tree has been split in two, East and Southern African groups are on the left, West and Central West African groups are on the right. Leaves are labelled by the identity of the individuals within them, with the total number of individuals in the clusters shown in parentheses. Leaves are coloured by the country of origin (as in Fig. 1-figure supplement 1) and branches are coloured by the final ancestry region that the clusters were assigned to. Note that although Malawi and Cameroon individuals were located in a clade with mostly East African individuals, they were assigned to Southern and Central West African ancestry regions, respectively. Clades containing outlying individuals from the Fula and Mandinka are also shown.

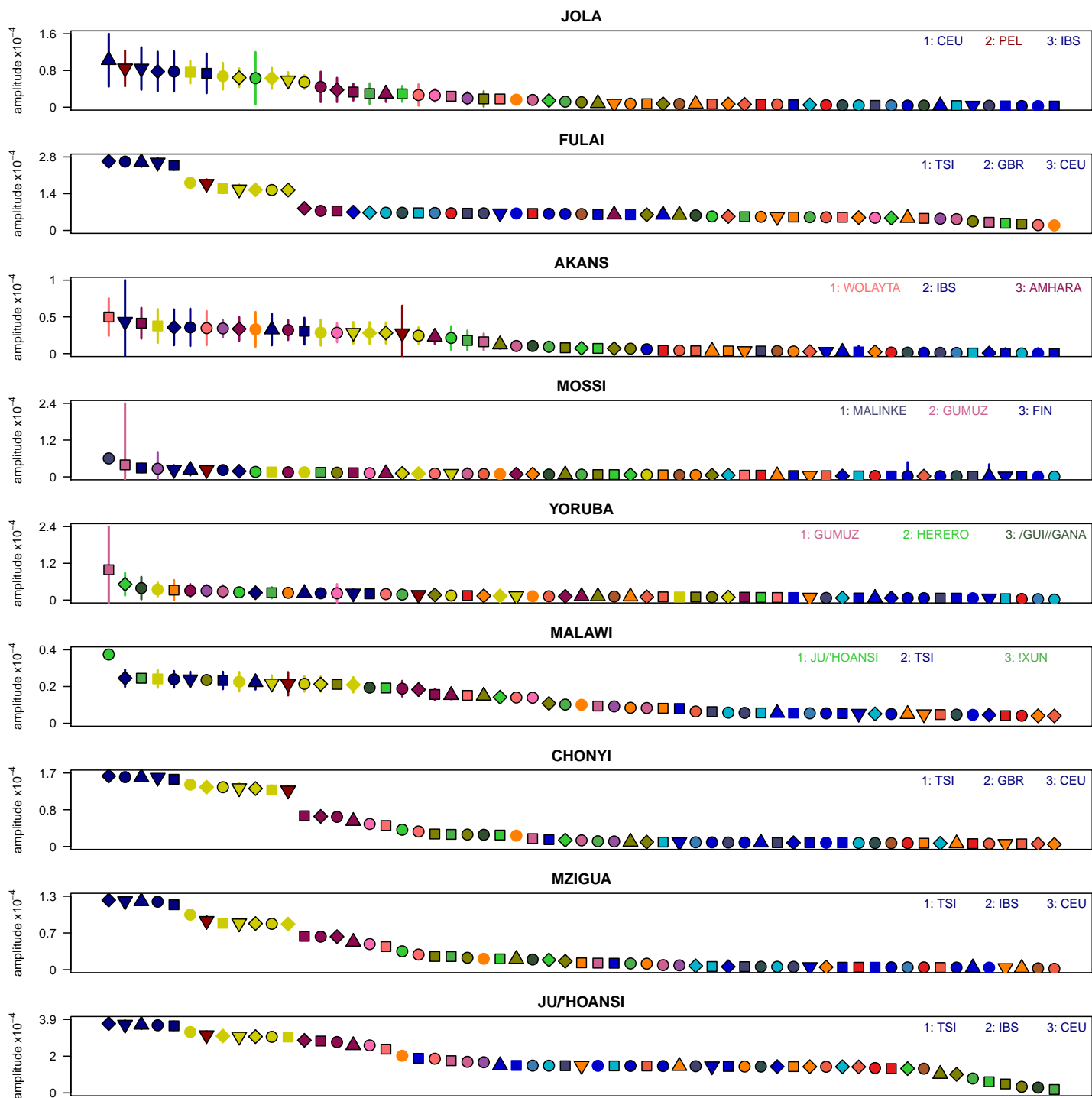


Figure 3-figure supplement 1. Weighted LD amplitudes for a selection of 9 ethnic groups. For a given test population we show the amplitude (± 1 s.e.) computed using a test population and every other population as the second reference. Plotted are the fitted amplitudes for each set of curves with the population used labelled beneath, with populations ordered by amplitude. A large number of population showed a similar profile to (A), that is with Eurasian populations showing the highest amplitudes. Other populations, e.g. Malawi, obtained the largest amplitudes from an African population.

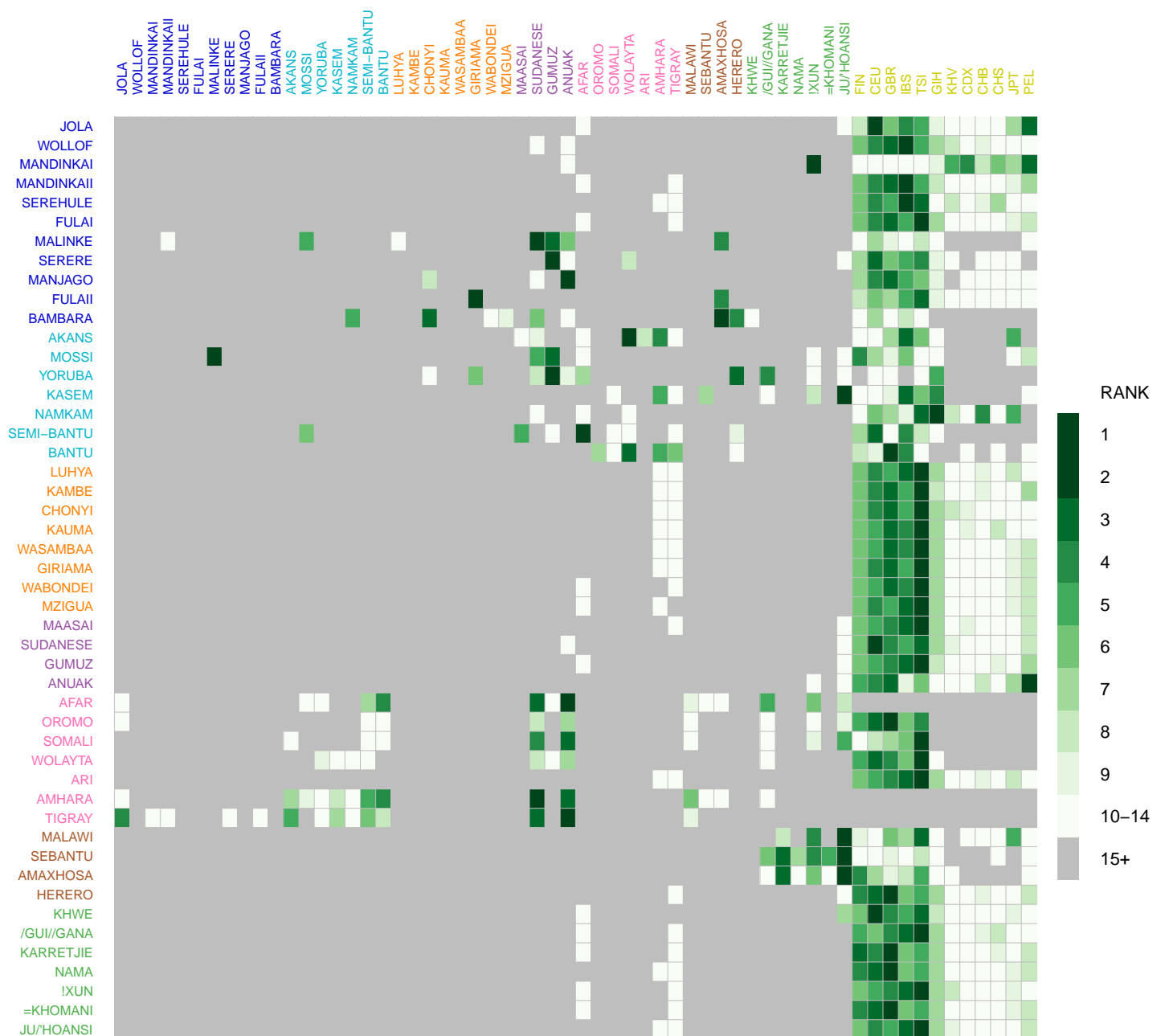


Figure 3-figure supplement 2. Comparison of weighted LD amplitude scores across all African ethnic groups. For a given test population we computed the ALDER amplitude (y-axis intercept) using the test population and every other population as the second reference. We then ranked the amplitudes across a given test population: populations who gave the top-ranked (i.e. largest) amplitude are in green, with those beneath a rank 15 shown in grey. This analysis shows that for many populations the reference populations giving the largest amplitudes (i.e. have the highest rank) are often non-African groups.

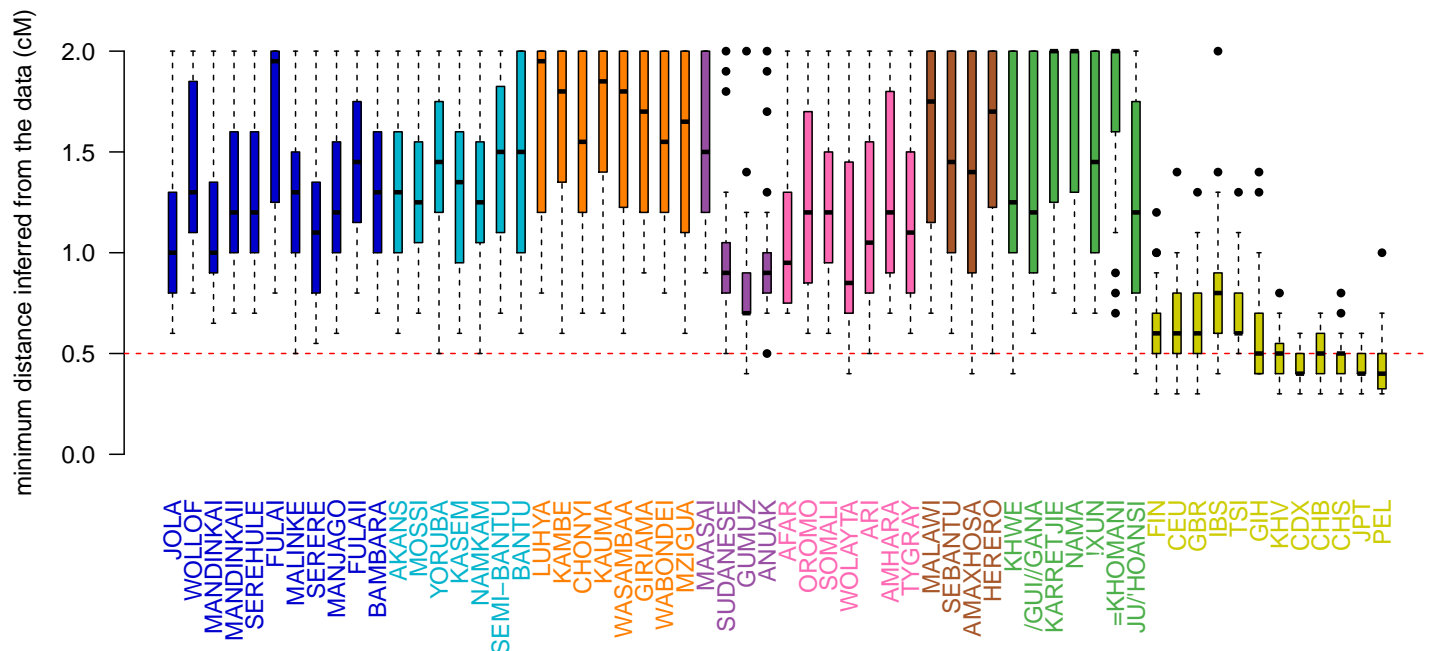


Figure 3-figure supplement 3. Comparison of the minimum distance to begin computing admixture LD. For each of the 48 African populations as a target, we used ALDER to compute the minimum distance over which short-range LD is shared with each of the 47 other African and 12 Eurasian reference populations. Here we show boxplots showing the distribution of minimum inferred genetic distances (y-axis) over which LD is shared for each of the reference populations separately (x-axis). We performed two analyses using weighted LD, one using these values of the minimum distance inferred from the data, and another where this distance was forced to be 0.5cM (dotted red line). The comparisons show that all African populations share LD correlations at distances $> 0.5\text{cM}$ with all other African populations. Note that ALDER computes LD correlations at distances $< 2\text{cM}$.

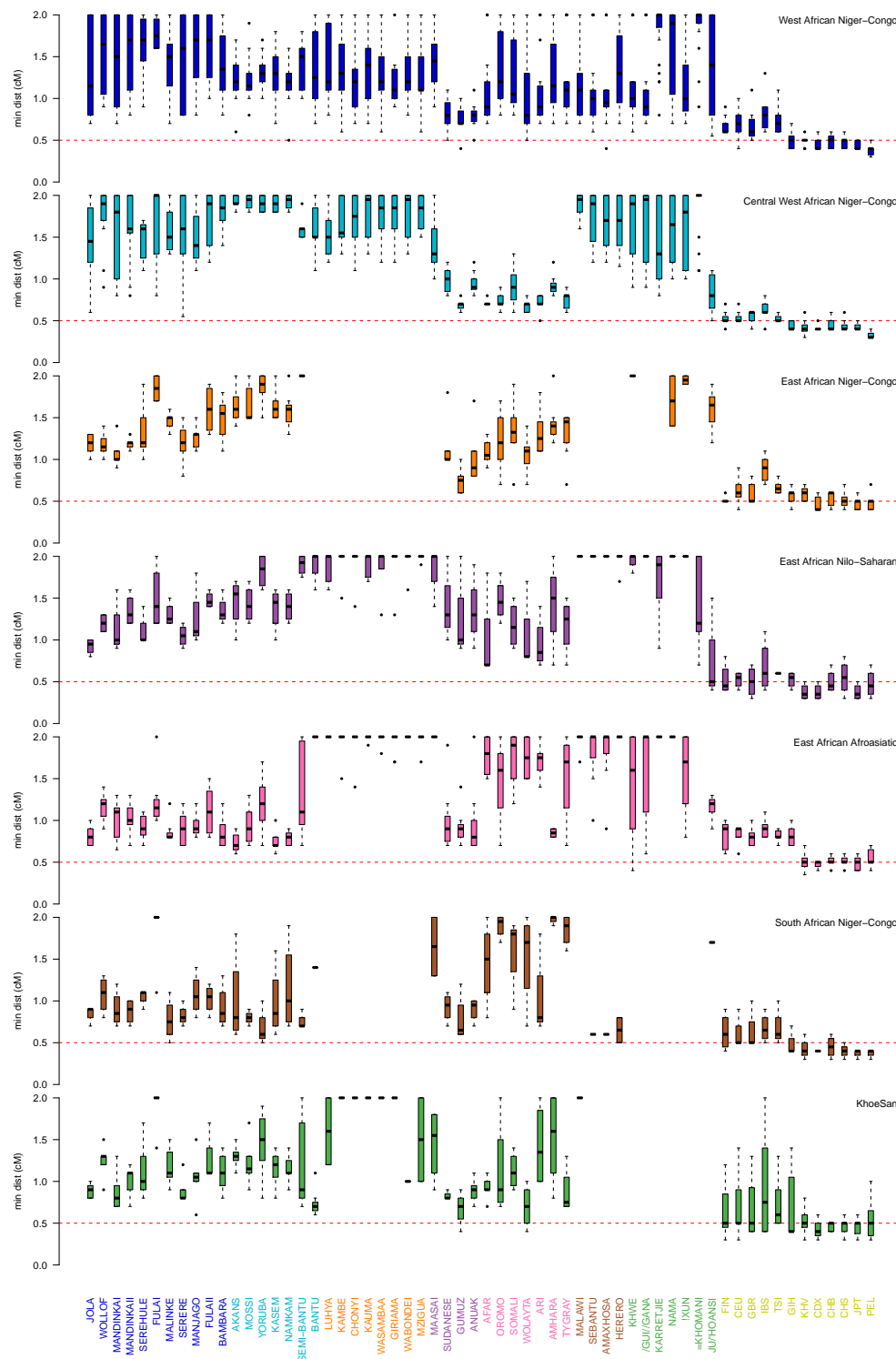


Figure 3-figure supplement 4. Comparison of the minimum distance to begin computing admixture LD split by region. As in Figure 3-figure supplement 3 except distances are stratified by region. The comparisons show that all African populations share LD correlations at distances $> 0.5\text{cM}$ with all other African populations. Note that ALDER computes LD correlations at distances $< 2\text{cM}$.

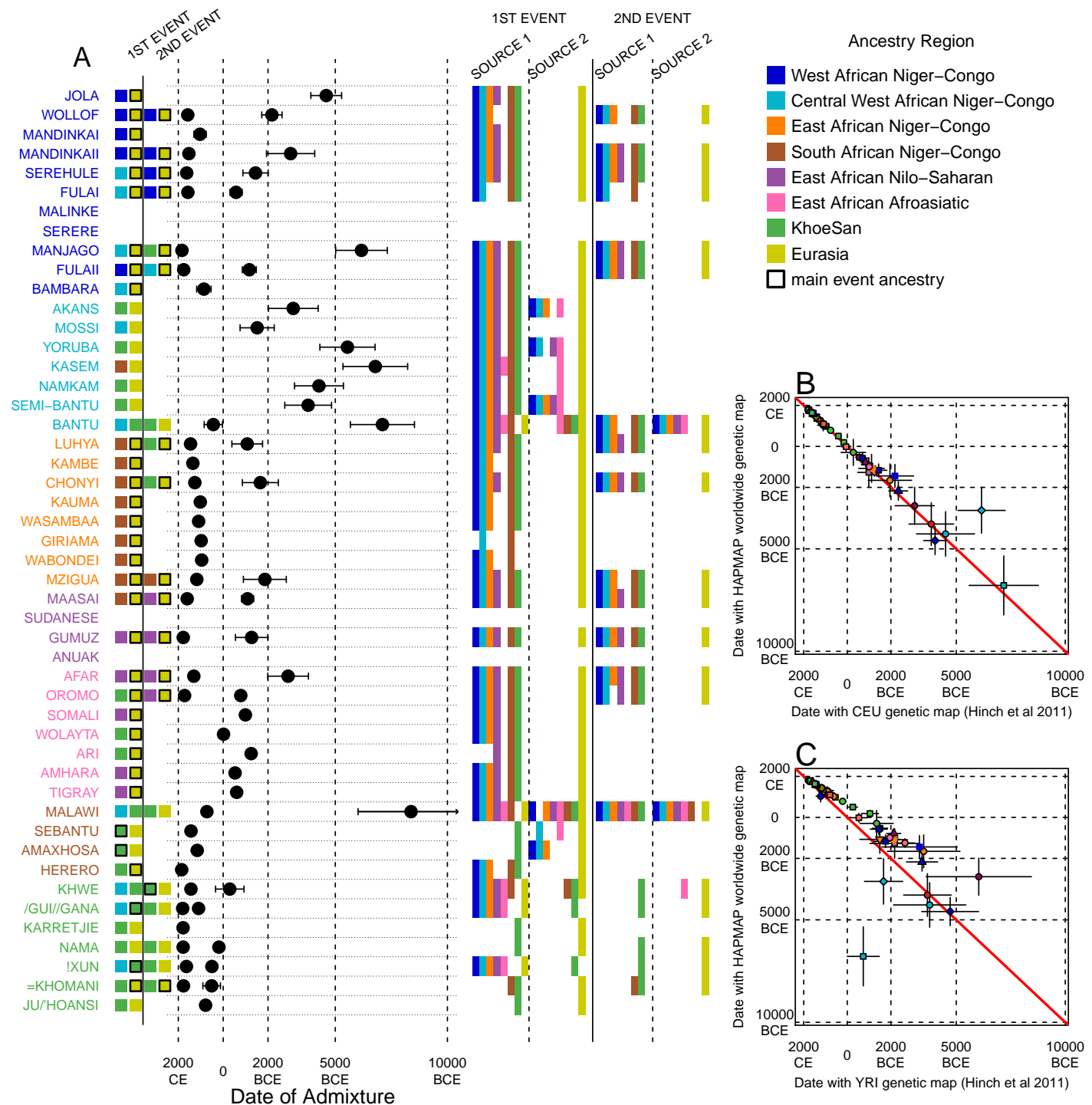
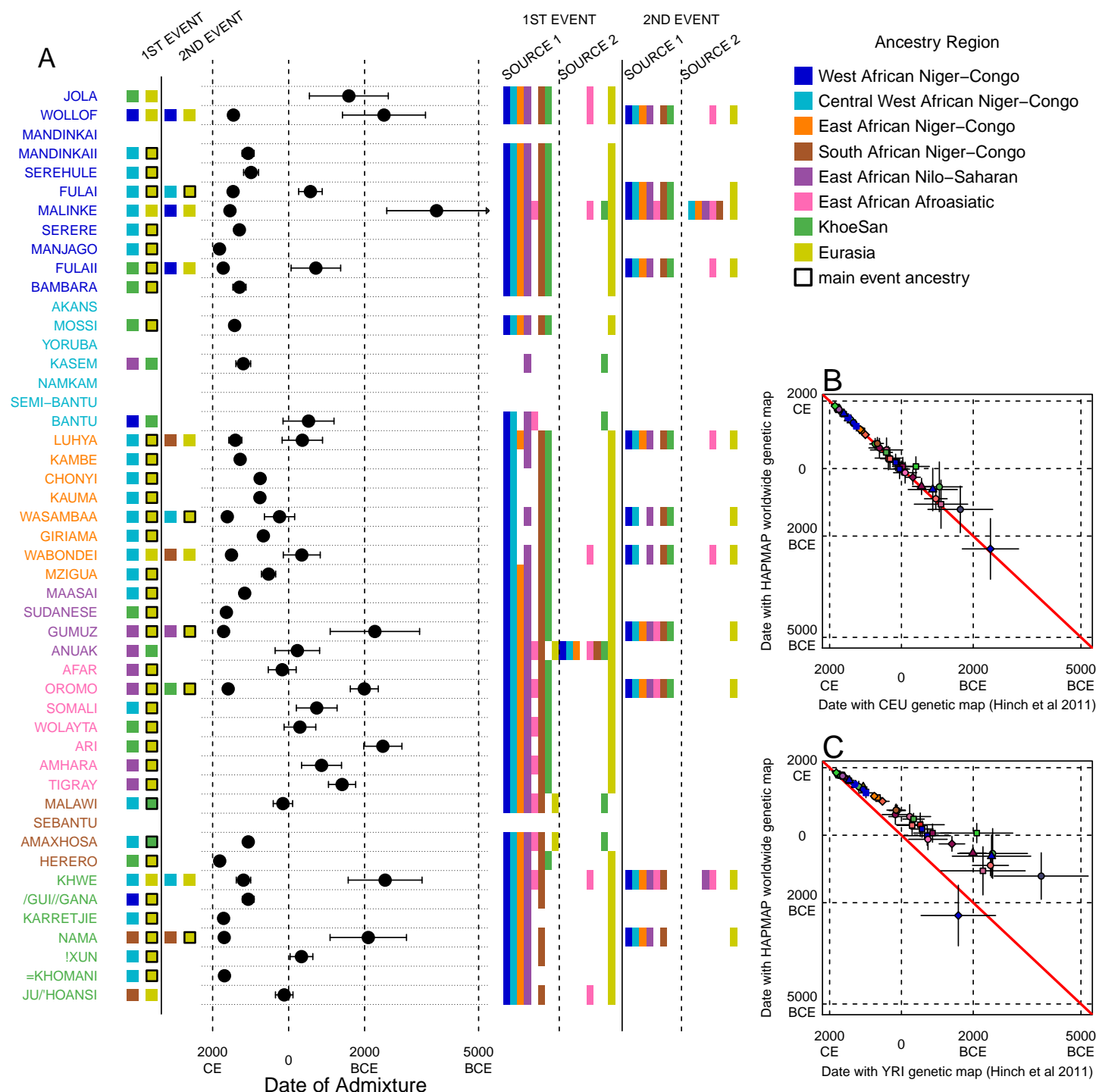


Figure 3-figure supplement 5. Results of the MALDER analysis computing weighted admixture decay curves from 0.5cM. As in the main analyses, the algorithm was run independently three times with the HAPMAP, YRI, and CEU genetic maps. The main results shown here are from the HAPMAP analysis. For each population, we show the ancestry region identity of the two populations involved in generating the MALDER curves with the greatest amplitudes (which are the closest to the true admixing sources amongst the reference populations) for at most two events. The sources generating the greatest amplitude are highlighted with a black box. Populations are ordered by ancestry of the admixture sources and dates estimates which are shown ± 1 s.e. (B) Comparison of dates of admixture ± 1 s.e. for MALDER dates inferred using the HAPMAP recombination map and a recombination map inferred from European (CEU) individuals from Hinch et al. [2011]. We only show comparisons for dates where the same number of events were inferred using both methods. Point symbols refer to populations and are as in Figure 1. (C) as (B) but comparing with an African (YRI) map.



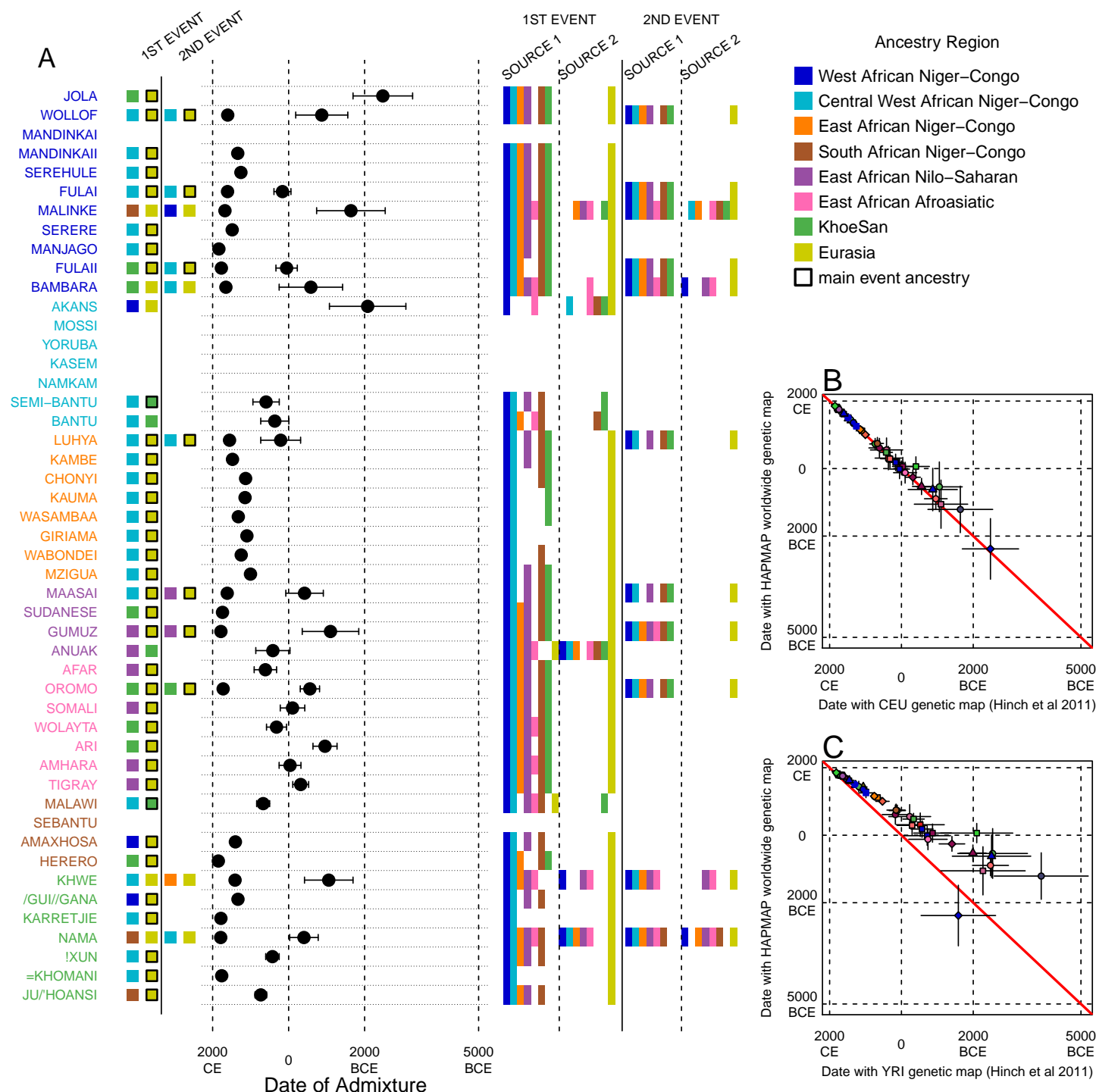


Figure 3-figure supplement 7. Results of MALDER for all populations using a European specific recombination map. We used MALDER to identify the evidence for multiple waves of admixture in each population. (A) For each population, we show the ancestry region identity of the two populations involved in generating the MALDER curves with the greatest amplitudes (which are the closest to the true admixing sources amongst the reference populations) for at most two events. The sources generating the greatest amplitude are highlighted with a black box. Populations are ordered by ancestry of the admixture sources and dates estimates which are shown ± 1 s.e. (B) Comparison of dates of admixture ± 1 s.e. for MALDER dates inferred using the HAPMAP recombination map and a recombination map inferred from European (CEU) individuals from Hinch et al. [2011]. We only show comparisons for dates where the same number of events were inferred using both methods. Point symbols refer to populations and are as in Figure 1. (C) as (B) but comparing with an African (YRI) map.

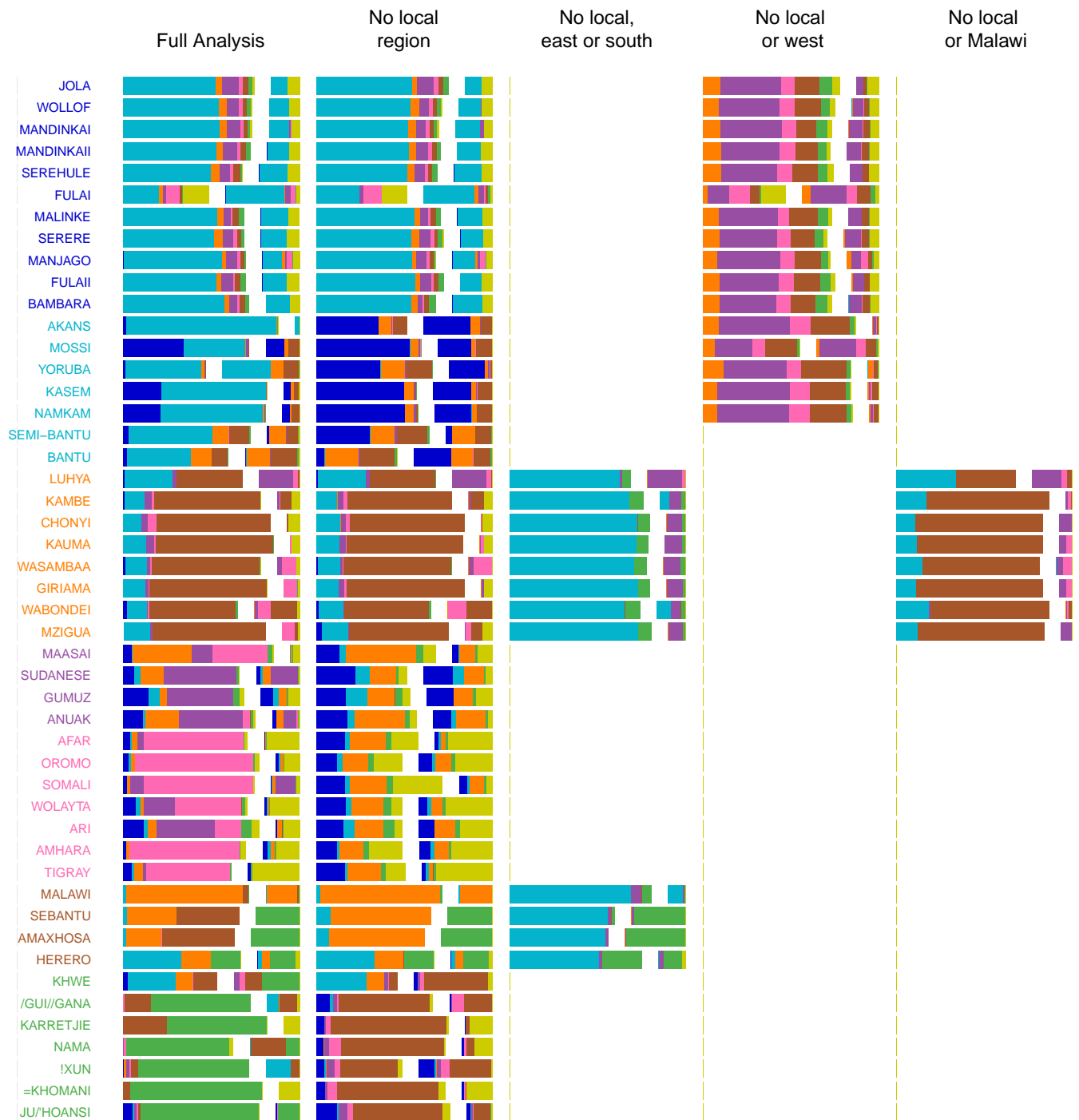


Figure 4-figure supplement 1. Admixture source inference by GLOBETROTTER after sequentially removing local surrogates from the analysis. In addition to the Full analysis, we show the inferred composition of admixture sources for different, restricted surrogate analyses. Components and y-axis labels are coloured by ancestry region. In each case we show admixture sources inferred by GLOBETROTTER for a single date of admixture.

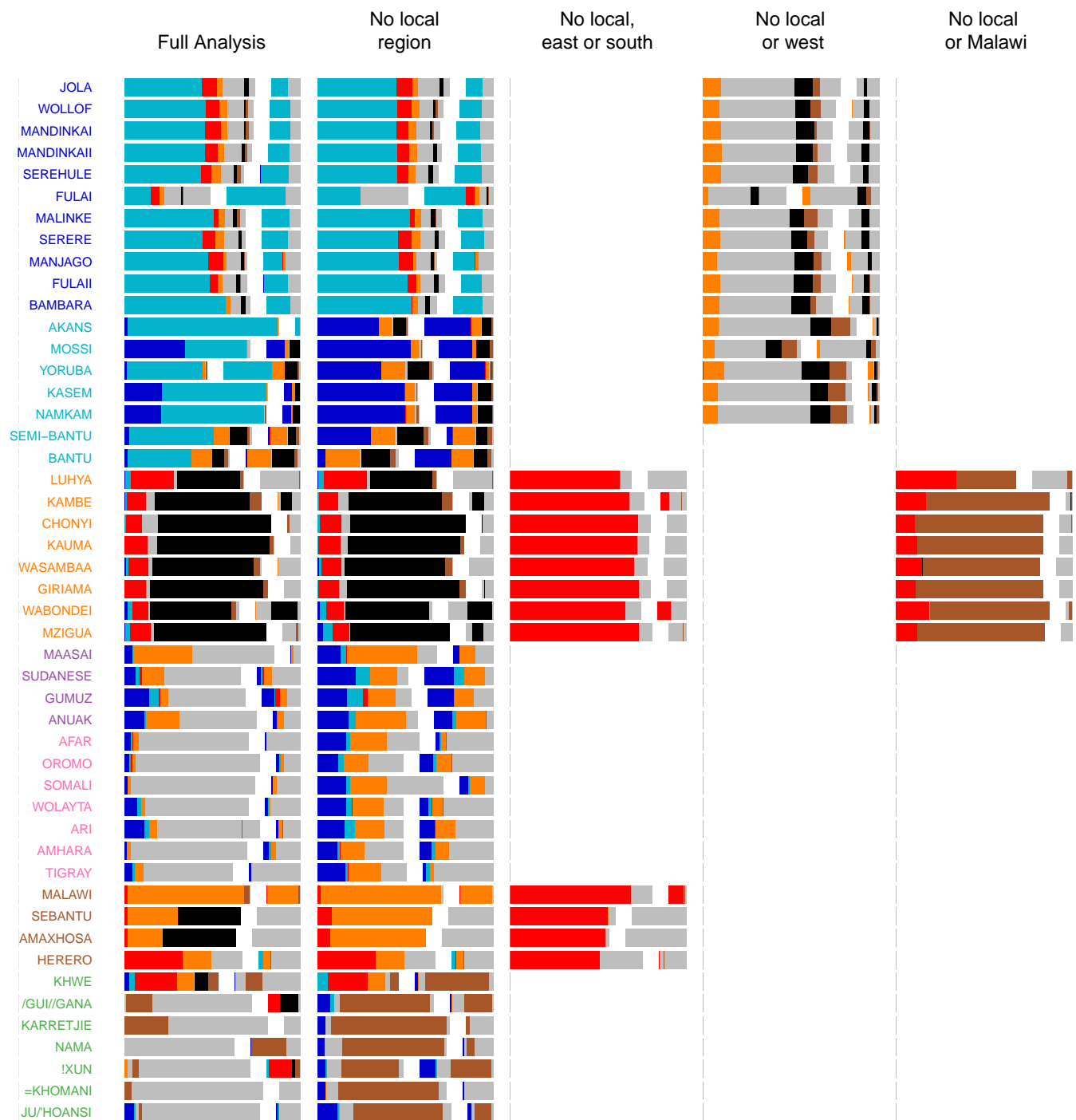


Figure 4-figure supplement 2. Admixture source inference by GLOBETROTTER after sequentially removing local surrogates from the analysis. The results are the same as Figure 4-figure supplement 1, but only Niger-Congo speaking groups are coloured. We highlight Malawi components in black, and Cameroon (Bantu and Semi-Bantu) in red.

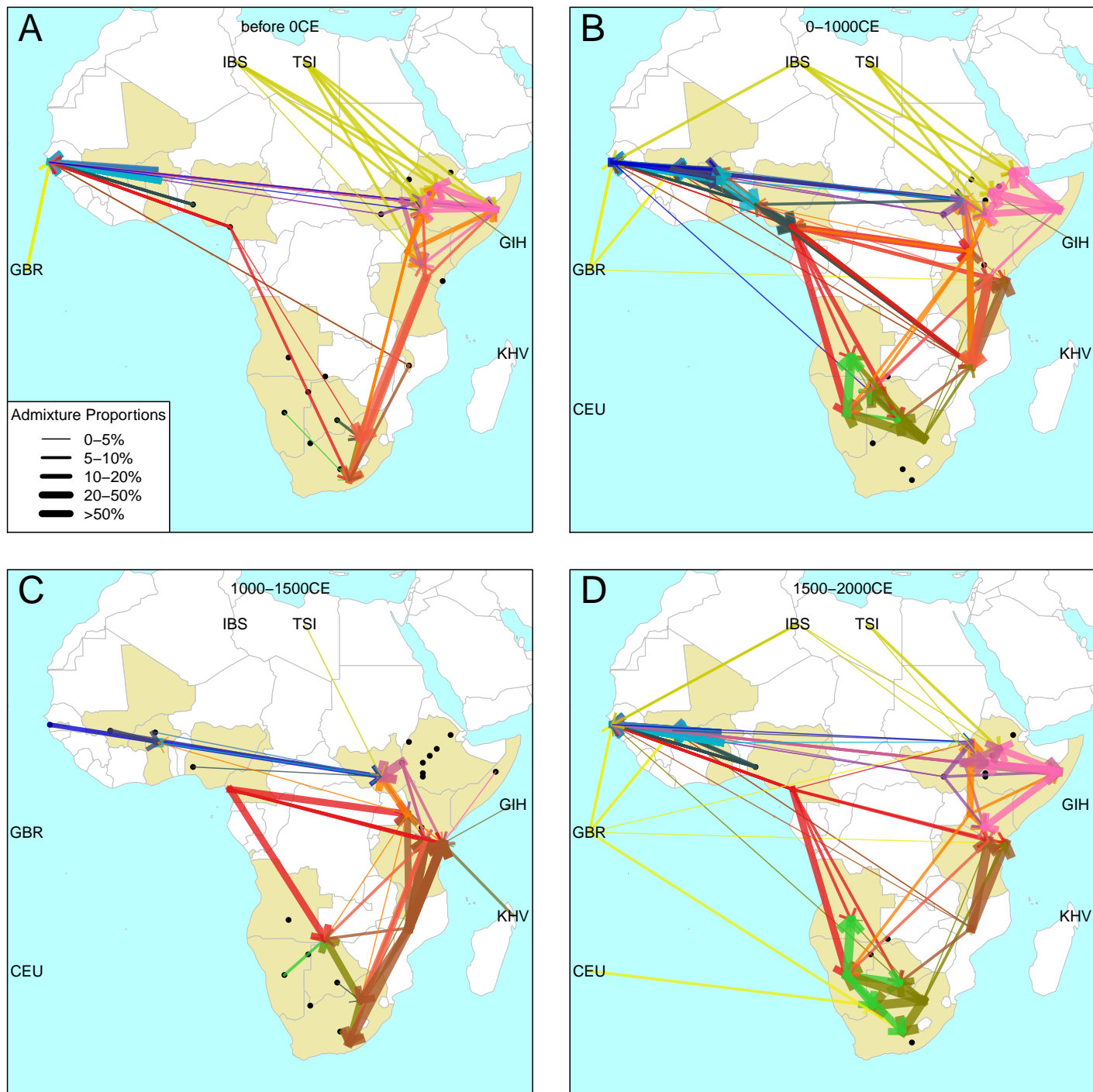


Figure 6-figure supplement 1. Gene-flow in Africa over the last 2,000 years. Using the results of the GLOBETROTTER analysis we show the connections between different groups in sub-Saharan Africa over time. For each population, we inferred the date of admixture and the composition of the admixing sources. We link each recipient population to its donor components using arrows, the size of which is proportional to the amount it contributes to the admixture event. Arrows are coloured by country of origin, as in Figure 4 in the main text.

Supplementary File 1 A note on ethnolinguistic groupings

The results of the population genetic analysis shows that population structure is largely the result of ethno-linguistic similarity, which itself is largely but not completely correlated with geographical proximity. These divisions are shown below and referred to in the text, together with the latest Ethnologue classification[‡] of the languages spoken, where possible.

1. 1st major Niger-Congo speaking group from West Africa: Gambian and Malian ethnic groups
 - Niger-Congo, **Mande** {Mandinka, Malinke, Bambara}
 - Niger-Congo, Atlantic-Congo, Atlantic, **Northern, Senegambian, Fula-Wolof** {Fula, Wollof}
 - Niger-Congo, Atlantic-Congo, Atlantic, **Northern, Senegambian, Serer** {Serere, Serehule}
 - Niger-Congo, Atlantic-Congo, Atlantic, **Northern, Bak** {Jola}
2. 2nd major Niger-Congo speaking group from West Africa: Ghana/BF/Nigerian ethnic groups:
 - Niger-Congo, Atlantic-Congo, Atlantic, **Volta-Congo, Kwa**{Akan, Yoruba}
 - Niger-Congo, Atlantic-Congo, Atlantic, **Volta-Congo, North, Gur** {Mossi, Kasem, Namkam?}
3. A Central and Eastern African Niger-Congo / “Bantoid” speaking group, split into two sub-divisions:
 - (a) North Western
 - Niger-Congo, Atlantic-Congo, Atlantic, Volta-Congo, **Benue-Congo, exNarrow-Bantu** {Cameroon: Bantu, Semi-Bantu?}
 - (b) Eastern
 - Niger-Congo, Atlantic-Congo, Atlantic, Volta-Congo, **Benue-Congo, Bantoid, Southern, Narrow-Bantu, Central, I** {Masaba-Luhya?}
 - Niger-Congo, Atlantic-Congo, Atlantic, Volta-Congo, **Benue-Congo, Bantoid, Southern, Narrow-Bantu, Central, E** {Mijikenda (Kenya)}
 - Niger-Congo, Atlantic-Congo, Atlantic, Volta-Congo, **Benue-Congo, Bantoid, Southern, Narrow-Bantu, Central, F-G** {Tanzania}
 - Niger-Congo, Atlantic-Congo, Atlantic, Volta-Congo, **Benue-Congo, Bantoid, Southern, Narrow-Bantu, Central, N** {Malawi (Chewa)}
4. Afroasiatic and Nilo-Saharan speakers from the Horn of Africa:
 - Afroasiatic, **Cushitic** {Afar, Somali, Oromo}
 - Afroasiatic, **Semitic** {Amhara, Tigrayan}

[‡]information accessed on 28th July 2014 at <https://www.ethnologue.com/>

- Afroasiatic, **Omotic** {Ari, Wolayta}
- Nilo-Saharan, **Komuz** {Gumuz}
- Nilo-Saharan, **Eastern Sudanic, Nilotic** {Maasai}
- Nilo-Saharan, **Eastern Sudanic, Nilotic** {Anuak}
- Nilo-Saharan {Sudanese}

5. Khoesan and Bantu speaking groups from Southern Africa

(a) Khoesan

- Khoesan, Southern Africa, **Northern** {Ju/'hoansi, !Xun}
- Khoesan, Southern Africa, **Central** {Nama}
- Khoesan, Southern Africa, **Central, Tshu-Khwe, Northwest** {/Gui//Gana, Khwe}
- Khoesan, Southern Africa, **Southern** {≠Khomani}
- uncertain {Karretijie}

(b) Southern Bantu speakers

- Niger-Congo, Atlantic-Congo, Atlantic, Volta-Congo, **Benue-Congo, Bantoid, Southern, Narrow-Bantu, Central, R** {Herero}
- Niger-Congo, Atlantic-Congo, Atlantic, Volta-Congo, **Benue-Congo, Bantoid, Southern, Narrow-Bantu, Central, S** {Amaxhosa, SEBantu}

Source Data

Figure 1-Source Data 1. Overview of sampled populations describing the continent, region, numbers of individuals used, and the source of any previously published datasets.

Ancestry Region	Country	Ethnic Group	Origin	N(phased)	IBD	REL	N(final)
Western Africa Niger-Congo	Gambia	JOLA	this study	147	0	31	116
Western Africa Niger-Congo	Gambia	MANJAGO	this study	47	0	0	47
Western Africa Niger-Congo	Gambia	SEREHULE	this study	47	0	0	47
Western Africa Niger-Congo	Gambia	SERERE	this study	50	0	0	50
Western Africa Niger-Congo	Gambia	WOLLOF	this study	135	0	27	108
Western Africa Niger-Congo	Gambia	MANDINKAI	this study	60	0	28	32
Western Africa Niger-Congo	Gambia	MANDINKAI	this study	71	0	0	71
Western Africa Niger-Congo	Gambia	FULAI	this study	97	0	25	72
Western Africa Niger-Congo	Gambia	FULAI	this study	79	0	0	79
Western Africa Niger-Congo	Mali	BAMBARA	this study	50	0	0	50
Western Africa Niger-Congo	Mali	MALINKE	this study	49	0	0	49
Central West Africa Niger-Congo	BurkinaFaso	MOSSI	this study	50	0	0	50
Central West Africa Niger-Congo	Ghana	AKANS	this study	50	0	0	50
Central West Africa Niger-Congo	Ghana	KASEM	this study	50	0	0	50
Central West Africa Niger-Congo	Ghana	NAMKAM	this study	50	0	0	50
Central West Africa Niger-Congo	Nigeria	YORUBA	1KG	161	0	60	101
East Africa Niger-Congo	Cameroon	BANTU	this study	50	0	0	50
East Africa Niger-Congo	Cameroon	SEMI-BANTU	this study	50	0	0	50
East Africa Niger-Congo	Kenya	LUHYA	1KG	100	9	0	91
East Africa Niger-Congo	Kenya	CHONYI	this study	47	0	0	47
East Africa Niger-Congo	Kenya	KAUMA	this study	97	0	0	97
East Africa Niger-Congo	Kenya	KAMBE	this study	42	0	0	42
East Africa Niger-Congo	Kenya	GIRIAMA	this study	46	0	0	46
East Africa Niger-Congo	Tanzania	MZIGUA	this study	50	0	0	50
East Africa Niger-Congo	Tanzania	WABONDEI	this study	48	0	0	48
East Africa Niger-Congo	Tanzania	WASAMBAA	this study	50	0	0	50
East Africa Nilo-Saharan	Kenya	MAASAI	1KG	31	0	0	31
East Africa Nilo-Saharan	Sudan	SUDANESE	Pagani 2012	24	0	0	24
East Africa Nilo-Saharan	Ethiopia	GUMUZ	Pagani 2012	19	0	0	19
East Africa Nilo-Saharan	Ethiopia	ANUAK	Pagani 2012	23	0	0	23
East Africa Nilo-Saharan	Somalia	SOMALI	Pagani 2012	40	0	0	40
East Africa Afroasiatic	Ethiopia	ARI	Pagani 2012	41	0	0	41
East Africa Afroasiatic	Ethiopia	AFAR	Pagani 2012	12	0	0	12
East Africa Afroasiatic	Ethiopia	TIGRAY	Pagani 2012	21	0	0	21
East Africa Afroasiatic	Ethiopia	AMHARA	Pagani 2012	26	0	0	26
East Africa Afroasiatic	Ethiopia	WOLAYTA	Pagani 2012	8	0	0	8
East Africa Afroasiatic	Ethiopia	OROMO	Pagani 2012	21	0	0	21
South Africa Niger-Congo	Malawi	MALAWI	this study	100	0	0	100
South Africa Niger-Congo	Namibia	HERERO	Schlebusch 2012	12	0	0	12
South Africa Niger-Congo	SouthAfrica	SEBANTU	Schlebusch 2012	20	0	0	20
South Africa Niger-Congo	SouthAfrica	AMAXHOSA	Peterson 2013	15	0	0	15
South Africa KhoeSan	SouthAfrica	KARRETJIE	Schlebusch 2012	20	0	0	20
South Africa KhoeSan	SouthAfrica	≠KHOMANI	Schlebusch 2012	39	0	0	39
South Africa KhoeSan	Namibia	NAMA	Schlebusch 2012	20	0	0	20
South Africa KhoeSan	Namibia	KHWE	Schlebusch 2012	17	0	0	17
South Africa KhoeSan	Angola	!XUN	Schlebusch 2012	33	0	0	33
South Africa KhoeSan	Botswana	/GUI //GANA	Peterson 2013	15	0	0	15
South Africa KhoeSan	Namibia	JU/'HOANSI	Schlebusch 2012	37	0	0	37
South Africa KhoeSan	Namibia	JU/'HOANSI	Peterson 2013	37	0	0	37
Europe	Finland	FIN	1KG	100	0	0	100
Europe	NorthEurope	CEU	1KG	104	0	2	102
Europe	Britain	GBR	1KG	101	1	3	97
Europe	Spain	IBS	1KG	150	0	50	100
Europe	Italy	TSI	1KG	100	1	0	99
Asia	India	GIH	1KG	100	2	0	98
Asia	Vietnam	KHV	1KG	121	1	21	99
Asia	China	CDX	1KG	100	8	0	92
Asia	China	CHB	1KG	101	0	1	100
Asia	China	CHS	1KG	150	7	50	93
Asia	Japan	JPT	1KG	100	0	0	100
Americas	Peru	PELII	1KG	105 [§]	0	0	16
				3699	29	298	3283

[§]89 admixed Peruvians and 517 individuals from five other 1KG American populations were removed from further analysis after phasing

Figure 2-Source Data 1. Pairwise F_{ST} for African populations. We used *smartpca* to compute F_{ST} for each pair of populations, upper left diagonal, together with standard errors computed using a block jackknife. F_{ST} has been multiplied by 1000

	JOLA	MANDINKAI	MANDINKAIH	SEREHULE	FULAI	MALINKE	SERERE	MANJAGO	FULAI	BANBARA	AKANS	MOSSI	YORUBA	KASEM	NANKAM	SEMI-BANTU	BANTU	LUHYA	KAMBE	CHONYI	KAUMA	WASAMBAA	GIRIAMA	WABONDEI	MZIGUA	MAASAI	WOLAYTA	SOMALI	ARI	SUDANESE	GUMUZ	OROMO	AMHARA	AFAR	ANUAK	TIGRAY	MALAWI	SEBANTU	MAAXHOSA	HERERO	KHWE	/GWI//GANA	KARRETJE	NAMA	IXUN	#KHOMANI	JU//HOANSI
	4	1	3	5	23	6	3	3	5	7	10	9	10	10	9	11	13	12	15	19	17	15	17	14	11	32	52	58	54	30	46	57	69	61	20	21	19	20	21	58	66	51	68	59	106		
WOLLUF	0.10	2	1	2	17	3	1	3	2	4	7	6	7	6	6	8	10	12	15	13	12	14	11	11	27	46	51	48	26	41	57	69	30	68	15	19	20	21	18	27	56	63	48	66	54	104	
MANDINKAI	0.10	0.10	0	1	2	18	2	1	2	4	7	6	7	6	8	10	12	15	13	12	14	11	11	27	46	51	48	26	41	57	69	30	68	15	19	20	21	18	27	56	63	48	66	54	104		
MANDINKAIH	0	0	0.10	1	1	18	2	1	2	4	7	6	7	6	8	10	12	15	13	12	14	11	11	28	48	53	50	27	42	51	61	64	26	63	12	16	17	18	26	62	48	65	56	64	104		
SEREHULE	0.10	0	0.10	0	17	1	1	2	3	1	5	4	6	5	7	8	11	14	12	11	13	10	10	26	45	48	45	24	41	50	59	61	25	62	11	15	16	17	25	55	61	47	64	55	102		
FULAI	0.20	0.20	0.20	0.20	20	19	18	20	17	20	23	22	24	23	23	24	25	24	23	22	26	23	24	10	26	45	48	45	24	41	50	59	61	25	62	11	15	16	17	25	55	61	47	64	55	102	
MALINKE	0.10	0.10	0.10	0.10	0.10	0.20	2	4	1	1	3	2	4	3	5	7	9	10	13	11	9	11	8	27	47	53	49	26	41	52	61	64	35	38	27	31	29	38	68	50	77	56	114	102	104		
SERERE	0.10	0	0.10	0	0.10	0.20	2	4	1	1	3	2	4	3	5	7	9	10	13	11	9	11	8	27	47	53	49	26	41	52	61	64	35	38	27	31	29	38	68	50	77	56	114	102	104		
MANJAGO	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0.2	4	3	6	5	8	10	12	14	12	11	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
BANBARA	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	3	6	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3	5	8	10	12	14	12	11	11	29	48	54	50	27	42	50	62	64	35	62	11	13	14	16	26	62	47	65	56	104	104			
KASSI	0.10	0.10	0.10	0.10	0.10	0.20	0.10	0.10	0.10	0	2	2	4	3																																	

Figure 2-Source Data 2. Pairwise F_{ST} for Eurasian populations. We used *smartpca* to compute F_{ST} for each pair of populations, upper left diagonal, together with standard errors computed using a block jackknife. F_{ST} has been multiplied by 1000

	FIN	CEU	GBR	IBS	TSI	GIH	KHV	CDX	CHB	CHS	JPT	PEL
JOLA	160	157	158	151	153	143	182	186	184	186	186	225
WOLLOF	151	148	149	142	144	135	175	178	177	178	178	217
MANDINKAI	154	151	152	145	147	138	177	180	179	180	180	219
MANDINKAI	153	150	151	144	146	137	176	179	178	180	180	218
SEREHULE	150	148	148	142	143	135	174	178	176	178	178	216
FULAI	112	108	109	102	103	102	147	151	149	151	151	188
MALINKE	154	152	152	146	148	138	177	181	180	181	181	220
SERERE	154	151	152	145	147	137	177	180	179	180	180	219
MANJAGO	154	152	152	146	147	138	178	181	180	181	181	220
FULAI	152	149	150	143	145	136	175	179	177	179	179	217
BAMBARA	155	152	153	147	148	139	178	181	180	181	182	220
AKANS	158	156	157	150	152	142	181	184	183	184	184	223
MOSSI	156	154	154	148	149	140	179	182	181	182	182	221
YORUBA	158	155	156	150	151	141	180	183	182	184	184	222
KASEM	157	155	155	149	150	141	180	183	182	183	183	222
NAMKAM	157	155	155	149	150	141	180	183	182	183	183	222
SEMI-BANTU	157	154	155	149	150	140	179	183	181	183	183	221
BANTU	158	155	156	150	151	141	180	183	182	183	183	222
LUHYA	147	144	145	138	140	130	170	174	172	174	174	213
KAMBE	142	140	140	134	135	127	167	170	169	171	171	209
CHONYI	149	146	147	141	142	133	173	176	175	176	176	215
KAUMA	146	143	144	138	139	130	170	174	172	174	174	213
WASAMBAA	142	139	140	133	135	126	167	170	169	170	170	209
GIRIAMA	148	145	146	140	141	132	172	175	174	175	176	214
WABONDEI	145	142	143	137	138	129	169	173	171	173	173	212
MZIGUA	148	145	146	140	141	132	172	175	174	175	175	214
MAASAI	104		101	94	95	92	137	141	139	141	141	179
WOLAYTA	77	73	73	67	67	70	120	124	122	123	123	159
SOMALI	78	73	74	67	67	71	122	126	124	126	126	162
ARI	115	112	112	106	106	103	147	150	149	150	150	189
SUDANESE	152	150	150	144	146	135	174	178	176	178	178	218
GUMUZ	147	144	145	139	140	130	170	173	172	173	174	213
OROMO	66	61	62	55	55	61	113	117	115	117	117	153
AMHARA	59	53	54	47	46	55	111	115	112	114	114	149
AFAR	64	58	59	52	52	60	114	118	116	117	117	153
ANUAK	150	148	149	143	144	134	173	176	175	176	177	216
TIGRAY	57	51	51	44	44	54	109	113	111	113	113	148
MALAWI	158	156	157	150	152	142	181	184	183	184	184	223
SEBANTU	158	156	157	150	152	142	181	184	183	184	184	223
AMAXHOSA	157	155	156	150	151	141	180	183	182	184	184	222
HERERO	145	143	143	138	139	131	172	175	174	175	176	213
KHWE	160	157	158	152	153	144	183	186	185	187	187	225
/GUI//GANA	190	187	188	182	183	173	212	215	214	215	216	250
KARRETJIE	165	162	163	157	159	150	190	193	192	193	194	232
NAMA	142	139	139	134	135	130	172	176	174	176	176	213
IXUN	193	191	192	186	187	178	216	220	218	220	220	250
≠KHOMANI	138	135	136	131	132	126	168	172	170	172	172	210
JU//HOANSI	226	224	224	219	220	210	249	250	250	250	250	250
FIN		6	7	10	12	36	101	106	102	104	104	127
CEU	0.10		0	2	4	35	109	113	110	112	112	134
GBR	0.10	0		3	4	35	109	113	110	112	112	134
IBS	0.10	0.10	0.10		2	36	108	113	110	112	111	136
TSI	0.10	0.10	0.10	0		34	108	113	110	112	112	137
GIH	0.30	0.30	0.30	0.30	0.30		72	77	75	77	76	114
KHV	0.90	1	1	1	1	0.60		2	7	4	14	111
CDX	0.90	1	1	1	1	0.70	0.10		9	6	17	115
CHB	0.90	1	1	1	1	0.70	0.10	0.10		1	7	107
CHS	0.90	1	1	1	1	0.70	0.10	0.10	0		9	110
JPT	0.90	1	1	1	1	0.70	0.20	0.20	0.10	0.10		107
PEL	1.20	1.30	1.30	1.30	1.30	1.10	1.20	1.20	1.20	1.20	1.20	

[illegible]

Figure 2-Source Data 4. Pairwise TVD for Eurasian populations. TVD has been multiplied by 1000

	FIN	CEU	GBR	IBS	TSI	GIH	KHV	CDX	CHB	CHS	JPT	PEL
JOLA	886	857	861	827	835	841	871	878	869	871	873	861
WOLLOF	859	827	832	797	805	812	846	856	844	846	850	833
MANDINKAI	869	838	842	808	816	822	854	863	852	854	858	843
MANDINKAI	854	821	826	790	798	805	842	852	840	842	846	828
SEREHULE	847	814	819	783	791	798	836	846	834	837	840	822
FULAI	860	829	834	800	808	814	844	853	842	844	847	834
MALINKE	851	817	822	786	793	801	840	850	838	840	844	826
SERERE	856	823	828	792	800	807	843	853	841	844	847	830
MANJAGO	854	822	827	792	800	807	841	851	839	842	845	829
FULAI	852	819	823	788	795	803	841	851	839	842	845	827
BAMBARA	854	820	825	789	797	805	843	853	841	843	847	829
AKANS	859	825	830	794	802	810	848	858	846	848	852	835
MOSSI	860	827	832	797	804	812	848	859	847	849	852	836
YORUBA	864	830	836	800	808	816	853	863	851	853	857	840
KASEM	868	836	841	806	813	821	855	864	853	855	859	843
NAMKAM	868	836	840	805	813	821	854	864	853	855	858	843
SEMI-BANTU	863	828	833	798	806	814	852	862	850	853	856	840
BANTU	869	833	839	807	813	820	857	867	856	858	861	846
LUHYA	855	814	820	785	791	799	847	858	845	848	852	831
KAMBE	867	837	842	812	817	820	855	866	854	856	860	843
CHONYI	891	865	869	848	850	850	876	884	875	877	880	867
KAUMA	889	862	866	845	846	847	873	882	872	874	877	864
WASAMBAA	855	815	820	786	791	799	847	857	845	847	851	830
GIRIAMA	880	851	856	828	832	834	867	876	865	867	870	855
WABONDEI	857	820	825	793	799	804	848	859	846	849	853	833
MZIGUA	861	825	830	799	804	810	852	863	850	853	856	837
MAASAI	790	741	748	702	705	729	782	793	780	782	786	765
WOLAYTA	664	608	616	567	572	597	716	738	705	720	720	633
SOMALI	706	659	666	616	621	646	736	759	728	740	741	682
ARI	813	769	776	731	734	754	805	815	803	805	809	789
SUDANESE	823	776	783	733	739	763	815	825	813	816	819	799
GUMUZ	798	750	757	708	713	738	790	800	788	790	794	774
OROMO	635	578	587	536	541	566	711	736	701	715	716	615
AMHARA	610	552	562	511	515	540	707	735	697	714	714	600
AFAR	630	577	584	534	540	564	713	740	703	719	719	615
ANUAK	816	768	775	726	731	756	808	818	806	808	812	792
TIGRAY	595	535	546	493	497	522	704	732	695	711	712	591
MALAWI	882	849	854	825	831	835	870	879	869	871	874	859
SEBANTU	877	845	850	821	827	831	865	874	864	866	869	854
AMAXHOSA	879	847	852	822	828	832	867	876	865	867	871	856
SWBANTU	825	790	795	765	772	776	813	822	811	813	817	801
KHWE	870	835	841	810	816	820	859	868	857	859	863	846
/GUI//GANA	874	841	846	816	822	826	862	871	860	863	866	851
KARRETIJE	781	747	753	723	729	733	769	783	767	770	773	757
NAMA	794	760	765	734	740	745	785	800	783	787	790	771
IXUN	884	851	856	829	834	837	871	879	869	871	874	860
≠KHOMANI	757	724	729	700	705	709	757	779	751	761	763	734
JU//HOANSI	894	866	870	856	856	857	881	890	880	882	885	874
FIN		341	354	378	378	523	733	754	725	740	738	620
CEU			50	107	136	452	704	732	695	711	711	573
GBR				354	50	126	463	709	735	700	716	579
IBS					378	107	126	87	431	702	693	709
TSI						378	136	153	87	416	692	698
GIH							523	452	463	431	416	483
KHV								733	704	709	702	617
CDX									754	732	735	146
CHB										725	695	217
CHS											740	137
JPT												34
PEL												

Figure 3-Source Data 1. The evidence for multiple waves of admixture in African populations using MALDER and the HAPMAP recombination map. For each event in each ethnic group we show the largest inferred amplitude and date of an admixture event involving two reference populations (Pop1 and Pop2). We additionally provide the ancestry region identity of the two main reference populations, together with Z scores for curve comparisons between this best curve and those containing populations from different ancestry regions. We use a cut-off of $Z < 2$ to decide whether sources from multiple ancestries best describe the admixture source.

Ethnic Group	amp	Date Date (CI)	Pop1	Pop1Anc	Pop1Anc-West Africa NC(Z)	Pop1Anc-Central West Africa NC(Z)	Pop1Anc-East Africa NC(Z)	Pop1Anc-Nilo-Saharan(Z)	Pop1Anc-Afroasiatic(Z)	Pop1Anc-South Africa NC(Z)	Pop1Anc-Khoesan(Z)	Pop1Anc-Eurasia(Z)	Pop1Anc-(Z)	Pop2	Pop2Anc-West Africa NC(Z)	Pop2Anc-Central West Africa NC(Z)	Pop2Anc-East Africa NC(Z)	Pop2Anc-Nilo-Saharan(Z)	Pop2Anc-Afroasiatic(Z)	Pop2Anc-South Africa NC(Z)	Pop2Anc-Khoesan(Z)	Pop2Anc-Eurasia(Z)	Pop2Anc-(Z)
JOLA	7.3e-05	148	31	JU//HOANSI	Khoesan	0.57 0.43 0.72 0.69 2.25 0.22 0.00	0.39 0.06 0.78 1.04 4.21 0.33 0.00	0.43 0.00 0.52 1.01 3.47 0.23 0.20	0.00 0.10 1.02 1.15 4.59 0.40 0.79	0.00 0.00 0.63 1.11 3.42 0.35 0.17	0.00 0.00 0.63 1.11 3.42 0.35 0.17	0.00	GBR	0.22	GBR	0.22	GBR	0.22	GBR	0.22	GBR	0.22	
WOLLOF	1e-04	88	22	AKANS	Central West Africa NC	-0.22 0.00 0.79 1.35 3.39 0.23 0.18	-0.22 0.00 0.79 1.35 3.39 0.23 0.18	-0.22 0.00 0.79 1.35 3.39 0.23 0.18	-0.22 0.00 0.79 1.35 3.39 0.23 0.18	-0.22 0.00 0.79 1.35 3.39 0.23 0.18	-0.22 0.00 0.79 1.35 3.39 0.23 0.18	0.00	GBR	0.00	GBR	0.00	GBR	0.00	GBR	0.00	GBR	0.00	
WOLLOF	2.3e-05	10	3	NAMKAM	Central West Africa NC	0.00 0.00 0.63 1.11 3.42 0.35 0.17	0.00 0.00 0.63 1.11 3.42 0.35 0.17	0.00 0.00 0.63 1.11 3.42 0.35 0.17	0.00 0.00 0.63 1.11 3.42 0.35 0.17	0.00 0.00 0.63 1.11 3.42 0.35 0.17	0.00 0.00 0.63 1.11 3.42 0.35 0.17	0.00	TSI	-0.22	TSI	-0.22	TSI	-0.22	TSI	-0.22	TSI	-0.22	
MANDINKAI	3.7e-05	19	3	JU//HOANSI	Khoesan	0.39 0.06 0.78 1.04 4.21 0.33 0.00	0.39 0.06 0.78 1.04 4.21 0.33 0.00	0.43 0.00 0.52 1.01 3.47 0.23 0.20	0.00 0.10 1.02 1.15 4.59 0.40 0.79	0.00 0.00 0.63 1.11 3.42 0.35 0.17	0.00 0.00 0.63 1.11 3.42 0.35 0.17	0.00	IBS	0.06	IBS	0.06	IBS	0.06	IBS	0.06	IBS	0.06	
SENHULE	6.3e-05	23	5	AKANS	Central West Africa NC	0.43 0.00 0.52 1.01 3.47 0.23 0.20	0.43 0.00 0.52 1.01 3.47 0.23 0.20	0.00 0.10 1.02 1.15 4.59 0.40 0.79	0.00 0.10 1.02 1.15 4.59 0.40 0.79	0.00 0.00 0.63 1.11 3.42 0.35 0.17	0.00 0.00 0.63 1.11 3.42 0.35 0.17	0.20	IBS	0.20	IBS	0.20	IBS	0.20	IBS	0.20	IBS	0.20	
FULAI	0.00044	60	8	JOLA	West Africa NC	0.00 0.10 1.02 1.15 4.59 0.40 0.79	0.00 0.10 1.02 1.15 4.59 0.40 0.79	0.00 0.10 1.02 1.15 4.59 0.40 0.79	0.00 0.10 1.02 1.15 4.59 0.40 0.79	0.00 0.00 0.63 1.11 3.42 0.35 0.17	0.00 0.00 0.63 1.11 3.42 0.35 0.17	0.10	CEU	0.10	CEU	0.10	CEU	0.10	CEU	0.10	CEU	0.10	
FULAI	0.00016	11	2	AKANS	Central West Africa NC	-0.19 0.00 0.87 1.17 4.72 0.26 0.72	-0.19 0.00 0.87 1.17 4.72 0.26 0.72	-0.19 0.00 0.87 1.17 4.72 0.26 0.72	-0.19 0.00 0.87 1.17 4.72 0.26 0.72	-0.19 0.00 0.87 1.17 4.72 0.26 0.72	-0.19 0.00 0.87 1.17 4.72 0.26 0.72	-0.19	IBS	-0.19	IBS	-0.19	IBS	-0.19	IBS	-0.19	IBS	-0.19	
MALINKE	8.9e-05	108	24	YORUBA	Central West Africa NC	0.03 0.00 0.27 0.90 0.55 0.24 0.11	0.03 0.00 0.27 0.90 0.55 0.24 0.11	0.03 0.00 0.27 0.90 0.55 0.24 0.11	0.03 0.00 0.27 0.90 0.55 0.24 0.11	0.03 0.00 0.27 0.90 0.55 0.24 0.11	0.03 0.00 0.27 0.90 0.55 0.24 0.11	0.03	PEL	0.03	PEL	0.03	PEL	0.03	PEL	0.03	PEL	0.03	
MALINKE	1.9e-05	8	2	AKANS	Central West Africa NC	-0.27 0.00 -0.18 0.30 0.10 -0.44	-0.27 0.00 -0.18 0.30 0.10 -0.44	-0.27 0.00 -0.18 0.30 0.10 -0.44	-0.27 0.00 -0.18 0.30 0.10 -0.44	-0.27 0.00 -0.18 0.30 0.10 -0.44	-0.27 0.00 -0.18 0.30 0.10 -0.44	0.10	TSI	0.10	TSI	0.10	TSI	0.10	TSI	0.10	TSI	0.10	
SERERE	2.5e-05	14	3	BANTU	Central West Africa NC	0.80 0.00 1.00 1.19 4.83 0.53 0.69	0.80 0.00 1.00 1.19 4.83 0.53 0.69	0.80 0.00 1.00 1.19 4.83 0.53 0.69	0.80 0.00 1.00 1.19 4.83 0.53 0.69	0.80 0.00 1.00 1.19 4.83 0.53 0.69	0.80 0.00 1.00 1.19 4.83 0.53 0.69	0.53	GBR	0.53	GBR	0.53	GBR	0.53	GBR	0.53	GBR	0.53	
MANJAGO	3.4e-05	3	1	AKANS	Central West Africa NC	0.24 0.00 0.92 1.73 6.49 0.25 0.86	0.24 0.00 0.92 1.73 6.49 0.25 0.86	0.24 0.00 0.92 1.73 6.49 0.25 0.86	0.24 0.00 0.92 1.73 6.49 0.25 0.86	0.24 0.00 0.92 1.73 6.49 0.25 0.86	0.24 0.00 0.92 1.73 6.49 0.25 0.86	0.24	GBR	0.24	GBR	0.24	GBR	0.24	GBR	0.24	GBR	0.24	
FULAI	2e-05	5	1	JU//HOANSI	Khoesan	0.09 0.20 1.16 2.17 5.69 0.36 0.00	0.09 0.20 1.16 2.17 5.69 0.36 0.00	0.09 0.20 1.16 2.17 5.69 0.36 0.00	0.09 0.20 1.16 2.17 5.69 0.36 0.00	0.09 0.20 1.16 2.17 5.69 0.36 0.00	0.09 0.20 1.16 2.17 5.69 0.36 0.00	0.09	GBR	0.09	GBR	0.09	GBR	0.09	GBR	0.09	GBR	0.09	
FULAI	7.6e-05	67	10	SEMI-BANTU	Central West Africa NC	0.18 0.00 0.78 1.74 5.30 -0.02 -0.01	0.18 0.00 0.78 1.74 5.30 -0.02 -0.01	0.18 0.00 0.78 1.74 5.30 -0.02 -0.01	0.18 0.00 0.78 1.74 5.30 -0.02 -0.01	0.18 0.00 0.78 1.74 5.30 -0.02 -0.01	0.18 0.00 0.78 1.74 5.30 -0.02 -0.01	0.18	TSI	0.18	TSI	0.18	TSI	0.18	TSI	0.18	TSI	0.18	
BAMBARA	2.2e-05	15	4	/GUI//GANA	Khoesan	0.62 0.31 0.69 0.67 2.94 0.25 0.00	0.62 0.31 0.69 0.67 2.94 0.25 0.00	0.62 0.31 0.69 0.67 2.94 0.25 0.00	0.62 0.31 0.69 0.67 2.94 0.25 0.00	0.62 0.31 0.69 0.67 2.94 0.25 0.00	0.62 0.31 0.69 0.67 2.94 0.25 0.00	0.25	IBS	0.25	IBS	0.25	IBS	0.25	IBS	0.25	IBS	0.25	
AKANS	0.00056	221	56	MALAWI	South Africa NC	0.05 0.15 0.06 0.72 0.88 0.00 0.65	0.05 0.15 0.06 0.72 0.88 0.00 0.65	0.05 0.15 0.06 0.72 0.88 0.00 0.65	0.05 0.15 0.06 0.72 0.88 0.00 0.65	0.05 0.15 0.06 0.72 0.88 0.00 0.65	0.05 0.15 0.06 0.72 0.88 0.00 0.65	0.05	GBR	0.05	GBR	0.05	GBR	0.05	GBR	0.05	GBR	0.05	
MOSSI	5.2e-06	8	2	/GUI//GANA	Khoesan	0.16 -0.62 -0.50 0.85 1.34 -0.43 0.00	0.16 -0.62 -0.50 0.85 1.34 -0.43 0.00	0.16 -0.62 -0.50 0.85 1.34 -0.43 0.00	0.16 -0.62 -0.50 0.85 1.34 -0.43 0.00	0.16 -0.62 -0.50 0.85 1.34 -0.43 0.00	0.16 -0.62 -0.50 0.85 1.34 -0.43 0.00	-0.43	CEU	-0.43	CEU	-0.43	CEU	-0.43	CEU	-0.43	CEU	-0.43	
YORUBA																							
KASEM																							
NAMKAM																							
SEMI-BANTU																							
BANTU	2.3e-05	56	12	MOSSI	Central West Africa NC	0.00 0.00 0.20 0.77 0.25 1.06	0.00 0.00 0.20 0.77 0.25 1.06	0.00 0.00 0.20 0.77 0.25 1.06	0.00 0.00 0.20 0.77 0.25 1.06	0.00 0.00 0.20 0.77 0.25 1.06	0.00 0.00 0.20 0.77 0.25 1.06	0.00	JU//HOANSI	0.00	JU//HOANSI	0.00	JU//HOANSI	0.00	JU//HOANSI	0.00	JU//HOANSI	0.00	
LUHYA	7e-05	25	2	AKANS	Central West Africa NC	0.97 0.00 0.00 7.52 11.82 1.74 2.63	0.97 0.00 0.00 7.52 11.82 1.74 2.63	0.97 0.00 0.00 7.52 11.82 1.74 2.63	0.97 0.00 0.00 7.52 11.82 1.74 2.63	0.97 0.00 0.00 7.52 11.82 1.74 2.63	0.97 0.00 0.00 7.52 11.82 1.74 2.63	0.97	TSI	0.97	TSI	0.97	TSI	0.97	TSI	0.97	TSI	0.97	
KAMBE	0.00014	15	2	YORUBA	Central West Africa NC	0.23 0.00 0.00 2.10 6.27 1.21 1.21	0.23 0.00 0.00 2.10 6.27 1.21 1.21	0.23 0.00 0.00 2.10 6.27 1.21 1.21	0.23 0.00 0.00 2.10 6.27 1.21 1.21	0.23 0.00 0.00 2.10 6.27 1.21 1.21	0.23 0.00 0.00 2.10 6.27 1.21 1.21	0.23	TSI	0.23	TSI	0.23	TSI	0.23	TSI	0.23	TSI	0.23	
CHONYI	0.00014	27	2	YORUBA	Central West Africa NC	0.17 0.00 0.00 3.28 9.02 1.89 1.65	0.17 0.00 0.00 3.28 9.02 1.89 1.65	0.17 0.00 0.00 3.28 9.02 1.89 1.65	0.17 0.00 0.00 3.28 9.02 1.89 1.65	0.17 0.00 0.00 3.28 9.02 1.89 1.65	0.17 0.00 0.00 3.28 9.02 1.89 1.65	0.17	TSI	0.17	TSI	0.17	TSI	0.17	TSI	0.17	TSI	0.17	
KAUMA	0.00015	26	2	SEMI-BANTU	Central West Africa NC	0.21 0.00 0.00 3.09 8.19 1.39	0.21 0.00 0.00 3.09 8.19 1.39	0.21 0.00 0.00 3.09 8.19 1.39	0.21 0.00 0.00 3.09 8.19 1.39	0.21 0.00 0.00 3.09 8.19 1.39	0.21 0.00 0.00 3.09 8.19 1.39	0.21	TSI	0.21	TSI	0.21	TSI	0.21	TSI	0.21	TSI	0.21	
WASAMBAA	0.00012	20	3	NAMKAM	Central West Africa NC	0.14 0.00 0.00 2.40 6.34 1.14 1.44	0.14 0.00 0.00 2.40 6.34 1.14 1.44	0.14 0.00 0.00 2.40 6.34 1.14 1.44	0.14 0.00 0.00 2.40 6.34 1.14 1.44	0.14 0.00 0.00 2.40 6.34 1.14 1.44	0.14 0.00 0.00 2.40 6.34 1.14 1.44	0.14	TSI	0.14	TSI	0.14	TSI	0.14	TSI	0.14	TSI	0.14	
GHIAAMA	0.00014	28	2	YORUBA	Central West Africa NC	0.63 0.00 0.00 4.08 9.30 2.09	0.63 0.00 0.00 4.08 9.30 2.09	0.63 0.00 0.00 4.08 9.30 2.09	0.63 0.00 0.00 4.08 9.30 2.09	0.63 0.00 0.00 4.08 9.30 2.09	0.63 0.00 0.00 4.08 9.30 2.09	0.63	TSI	0.63	TSI	0.63	TSI	0.63	TSI	0.63	TSI	0.63	
WABONDEI	1e-04	23	2	MOSSI	Central West Africa NC	0.16 0.00 0.00 2.94 8.04 1.41 1.97	0.16 0.00 0.00 2.94 8.04 1.41 1.97	0.16 0.00 0.00 2.94 8.04 1.41 1.97	0.16 0.00 0.00 2.94 8.04 1.41 1.97	0.16 0.00 0.00 2.94 8.04 1.41 1.97	0.16 0.00 0.00 2.94 8.04 1.41 1.97	0.16	TSI	0.16	TSI	0.16	TSI	0.16	TSI	0.16	TSI	0.16	
MZIGUA	0.00011	32	4	YORUBA	Central West Africa NC	0.18 0.00 0.00 1.98 5.11 1.12	0.18 0.00 0.00 1.98 5.11 1.12	0.18 0.00 0.00 1.98 5.11 1.12	0.18 0.00 0.00 1.98 5.11 1.12	0.18 0.00 0.00 1.98 5.11 1.12	0.18 0.00 0.00 1.98 5.11 1.12	0.18	TSI	0.18	TSI	0.18	TSI	0.18	TSI	0.18	TSI	0.18	
MAASAI	0.00037	74	15	SUDANESE	Nilo-Saharan	0.23 0.08 0.50 0.00 3.35 0.09 0.54	0.23 0.08 0.50 0.00 3.35 0.09 0.54	0.23 0.08 0.50 0.00 3.35 0.09 0.54	0.23 0.08 0.50 0.00 3.35 0.09 0.54	0.23 0.08 0.50 0.00 3.35 0.09 0.54	0.23 0.08 0.50 0.00 3.35 0.09 0.54	0.08	TSI	0.08	TSI	0.08	TSI	0.08	TSI	0.08	TSI	0.08	
MAASAI	8e-05	10	2	SEMI-BANTU	Central West Africa NC	0.44 0.00 0.47 -																	

Figure 3-Source Data 2. The evidence for multiple waves of admixture in African populations using MALDER and the African recombination map. Columns as in Figure 3-Source Data 1

Ethnic Group	amp	Date Date (CI)	Pop1	Pop1Anc	Pop1Anc-West Africa NC(Z)	Pop1Anc-Nilo-Saharan(Z)	Pop1Anc-East Africa NC(Z)	Pop1Anc-Afroasiatic(Z)	Pop1Anc-South Africa NC(Z)	Pop1Anc-Khoesan(Z)	Pop1Anc-Eurasia(Z)	Pop1Anc-(Z)	Pop2	Pop2Anc	Pop2Anc-West Africa NC(Z)	Pop2Anc-Central West Africa NC(Z)	Pop2Anc-East Africa NC(Z)	Pop2Anc-Nilo-Saharan(Z)	Pop2Anc-Afroasiatic(Z)	Pop2Anc-South Africa NC(Z)	Pop2Anc-Khoesan(Z)	Pop2Anc-Eurasia(Z)	Pop2Anc-(Z)	
JOLA	2.4e-05	121	36	/GUI//GANA	Khoean	0.39	0.32	0.44	0.34	0.32	0.00	0.32	CDX	Eurasia	2.53	2.82	2.81	2.36	1.54	2.61	0.00	1.54		
WOLOF	0.00017	153	38	MANDINKAI	West Africa NC	0.00	0.80	1.00	1.25	2.00	0.83	1.03	TSI	Eurasia		3.18	3.18	3.17	1.67	3.12	3.00	0.00	1.57	
WOLOF	2.4e-05	16	5	JOLA	West Africa NC	0.00	0.87	1.10	1.41	2.23	1.01	1.15	TSI	Eurasia								1.67		
MANDINKAI	3.7e-05	29	5	SEMI-BANTU	Central West Africa NC	0.53	0.00	0.93	1.12	4.07	0.39	0.12	IBS	Eurasia	6.33	6.33	7.70	7.02	4.51	7.79	7.22	0.00	4.51	
SEREHULE	5.7e-05	32	7	BANTU	Central West Africa NC	0.60	0.00	0.61	0.97	3.71	0.40	0.22	IBS	Eurasia	6.60	6.61	5.85	3.80	6.26	0.00	3.80	3.80		
FULAI	0.00041	86	11	SEMI-BANTU	Central West Africa NC	0.08	0.00	0.89	1.28	5.26	0.21	0.81	IBS	Eurasia	11.60	11.46	10.28	6.63	11.57	0.00	6.45	6.45		
FULAI	0.00015	16	2	AKANS	Central West Africa NC	-0.15	0.00	0.75	0.87	4.09	0.34	0.43	-0.15	Eurasia	8.62	8.53	7.68	5.13	8.59	0.00	5.13	5.13		
MALINKE	0.00021	200	45	BAMBARA	West Africa NC	0.00	0.58	0.63	0.80	0.99	0.39	0.57	JPT	Eurasia		1.21	1.15	0.78	1.15	0.00	0.78	0.78		
MALINKE	1.9e-05	13	2	SEMI-BANTU	Central West Africa NC	-0.24	0.00	0.03	0.56	1.26	-0.15	-0.22	-0.15	Eurasia	2.06	2.17	2.04	0.99	0.88	0.00	1.25	1.25		
SEREH	2.5e-05	22	4	BANTU	Central West Africa NC	0.60	0.00	1.08	1.38	4.94	0.83	0.92	0.60	Eurasia	6.65	6.82	7.37	4.63	8.29	7.90	0.00	4.63		
MANJAGO	3e-05	4	1	AKANS	Central West Africa NC	0.13	0.00	0.47	0.92	3.36	0.11	0.44	0.11	Eurasia	6.41	6.76	6.23	4.21	6.72	6.41	0.00	4.21		
FULAI	7.1e-05	91	22	JOLA	West Africa NC	0.00	0.11	0.52	0.68	2.07	0.27	0.53	0.11	Eurasia		3.16	3.08	2.88	1.97	3.12	2.89	0.00	1.94	
FULAI	1.8e-05	7	2	JU//HOANSI	Khoean	-0.93	-0.96	-0.41	-0.38	2.79	-0.84	0.00	-0.93	Eurasia	6.08	5.71	5.84	4.11	5.23	0.00	4.11	4.11		
BAMBARA	2.1e-05	22	6	/GUI//GANA	Khoean	0.45	0.15	0.56	0.60	2.78	0.01	0.00	0.01	Eurasia	5.36			3.17			0.00	3.17		
AKANS																								
MOSSI	7.3e-06	17	5	/GUI//GANA	Khoean	0.17	0.24	0.29	0.73	4.92	0.15	0.00	0.15	CEU	5.47			5.27	3.13		0.00	3.13		
YORUBA																								
KASEM	1.9e-06	25	7	GUMU'Z	Nilo-Saharan				0.00				4.91	JU//HOANSI	3.26			2.83		0.00		2.10		
NAMKAM																								
SEMI-BANTU																								
BANTU	2.3e-05	84	23	FULAI	West Africa NC	0.00	0.13	0.96	0.42			0.13	JU//HOANSI	Khoean	3.04	2.93	3.29	1.57	2.31	0.00	0.63	0.63		
LUHYA	8.7e-05	79	18	MALAWI	South Africa NC	0.89	0.18	0.61	1.41	4.02	0.00	0.65	0.18	TSI	4.30	3.87		3.29	1.57		3.09	0.00	1.57	
LUHYA	2.2e-05	18	6	YORUBA	Central West Africa NC	-0.18	0.00	-0.35	0.98	4.25	-1.00	0.85	-1.00	TSI	5.93	5.93		3.98	2.02	4.23	4.83	0.00	2.02	
KAMBE	0.00013	22	3	YORUBA	Central West Africa NC	0.08	0.00		1.77	5.34	0.24	0.69	0.08	TSI	8.56			7.99	4.93		8.52	0.00	4.93	
CHONYI	0.00013	40	3	YORUBA	Central West Africa NC	0.02	0.00		3.09	9.17	0.31	1.66	0.02	TSI	14.27			13.25	7.94		13.21	0.00	7.94	
KAUMA	0.00014	40	3	SEMI-BANTU	Central West Africa NC	0.33	0.00		3.02	7.44	0.21	0.75	0.21	TSI	11.04			10.01	5.92	11.85	10.72	0.00	5.92	
WASAMBAA	0.00013	58	14	BANTU	Central West Africa NC	0.12	0.00		1.26	3.93	0.03	0.46	0.03	TSI	5.25			4.39	2.47	4.90	4.54	0.00	2.29	
WASAMBAA	2.5e-05	11	3	YORUBA	Central West Africa NC	0.30	0.00		1.25	3.96	0.21	0.91	0.21	GBR	5.22			4.36	2.37	4.88	4.78	0.00	2.37	
GIRIAMA	0.00014	43	3	YORUBA	Central West Africa NC	0.49	0.00		3.88	9.92	0.35	1.66	0.35	TSI	12.40			10.92	6.16	12.29	11.83	0.00	6.16	
WABONDEI	0.00013	78	17	AMAXHOSA	South Africa NC	0.22	0.02		1.00	2.43	0.00	0.11	0.02	GBR	2.97	2.87		2.85	1.58	3.41	0.00	1.59		
WABONDEI	3e-05	14	4	BANTU	Central West Africa NC	0.25	0.00		1.21	2.50	-0.01	0.37	0.25	TSI	3.55			3.02	1.98	3.11	3.15	0.00	1.98	
MZIGUA	0.00011	26	4	YORUBA	Central West Africa NC	0.16	0.00	0.63	1.86	4.41	0.11	0.63	0.11	TSI	6.34			6.78	5.54	6.02	5.67	0.00	3.15	
MAASAI	0.00018	26	4	YORUBA	Central West Africa NC	0.05	0.00	1.37	1.24	5.63	0.76	1.37	0.05	TSI	8.66			9.16	8.42	4.89	8.92	8.73	0.00	4.89
SUDANESE	1.6e-05	10	3	/GUI//GANA	Khoean	0.42	0.10	0.69	0.53	3.42	0.35	0.00	0.10	GBR	5.33	5.53	5.66	4.46	2.56		0.00	2.56		
GUMU'Z	1.8e-05	7	1	SUDANESE	Nilo-Saharan	0.65	0.31	0.80	0.00	3.03	0.78	0.32	0.31	GBR	7.71	8.03		5.03	7.56	7.33	0.00	4.99		
GUMU'Z	7e-05	145	41	SUDANESE	Nilo-Saharan	-0.52	-0.91	0.21	0.00	1.99	1.08	-0.46	-0.91	PEL	9.39	9.92		5.41	9.06	9.04	0.00	8.86		
ANUAK	3.4e-05	74	20	SUDANESE	Nilo-Saharan	1.33	1.30	1.56	0.00	1.26	1.72	0.71	0.35	JU//HOANSI	1.06	1.05	1.12	0.78	1.37	0.00	1.27	0.35		
AFAR	0.00032	60	13	SUDANESE	Nilo-Saharan	0.43	0.25	0.64	0.00	2.16	0.35	0.39	0.25	FIN	5.50	5.46		4.22	5.50	5.55	0.00	4.20		
OROMO	0.00026	135	13	JU//HOANSI	Khoean	0.72	0.51	0.34	0.20	1.10	0.56	0.00	0.20	TSI	6.85	6.82	6.68	6.38	3.33	6.86	0.00	3.33		
OROMO	5.4e-05	11	2	GUMU'Z	Nilo-Saharan	-0.16	-0.29	0.05	0.00	0.55	-0.27	-0.85	-0.85	TSI	8.91	8.97	9.39	5.86	9.40	8.36	0.00	6.21		
SOMALI	0.00036	92	18	AKANS	Central West Africa NC	0.35	0.00	0.75	0.29	0.91	0.41	0.00	0.29	CEU	4.74			4.93	4.91	3.55	4.65	0.00	3.39	
WOLAYTA	0.00048	76	14	JU//HOANSI	Khoean	0.66	0.50	0.73	0.36	3.33	0.54	0.00	0.29	CEU	3.93			4.39	4.09	2.53	4.39	0.00	2.53	
ARI	0.00045	152	17	JU//HOANSI	Khoean	0.72	0.59	0.83	0.56	3.33	0.54	0.00	0.54	TSI				3.33	2.07		0.00	2.07		
AMHARA	0.00048	96	18	ANUAK	Nilo-Saharan	0.34	0.17	0.63	0.00	1.27	0.22	0.18	0.17	TSI	4.77	4.75	4.76	3.32	4.75		0.00	2.60		
TIGRAY	0.00064	115	12	SUDANESE	Nilo-Saharan	0.55	0.38	0.75	0.00	2.55	0.41	0.32	0.32	TSI	7.49			5.15	7.76		0.00	5.01		
MALAWI	3.7e-05	61	9	YORUBA	Central West Africa NC	0.36	0.00	0.70	1.36	1.22	1.10	0.50	0.36	JU//HOANSI	4.92			4.37	3.87		0.00	2.75	2.75	
SEBANTU																								
AMAXHOSA	0.00018	30	4	KASEM	Central West Africa NC	0.20	0.00	0.17	0.93	1.39	1.34	0.43	0.17	JU//HOANSI	7.76			7.90	7.56	6.89	0.00	5.54	5.43	
HERERO	0.00025	4	1	JU//HOANSI	Khoean	0.97	0.70	1.09	1.78	4.15	0.63	0.00	0.63	GBR	7.51	7.52	7.57	7.08	5.03	7.60	0.00	5.03		
KHWE	0.00026	154	34	SEMI-BANTU	Central West Africa NC	0.82	0.00	0.15	1.06	1.55	0.46	0.15	0.15	PEL	2.08	2.14	1.68	1.06	2.22		0.00	0.72		
KHWE	3.9e-05	25	7	YORUBA	Central West Africa NC	0.18	0.00	0.06	0.60	1.45	-0.30	-0.30	-0.30	TSI	2.66	2.71	2.26	1.81	2.92		0.00	1.36		
/GUI//GANA	3.7e-05	30	5	JOLA	West Africa NC	0.00	0.10	0.71	1.85	3.56	1.72	0.10	0.10	TSI				7.54	7.65	6.45	3.74	0.00	3.12	
KARRETHIE	0.00046	7	1	BANTU	Central West Africa NC	0.59	0.00	0.70	1.92	5.71		0.59	0.59	GBR	16.31			16.32	15.31	11.00		0.00	11.00	
NAMA	0.00029	8	1	MALAWI	South Africa NC	0.59	0.08	0.77	1.57	5.15	0.00	0.08	0.08	GBR	13.97			13.96	13.16	9.63		0.00	9.63	
NAMA	0.00045	139	35	MALAWI	South Africa NC	1.00	0.02	1.20	-1.95	5.99	0.00	1.20	1.20	PEL	13.80			13.79	12.62	9.88		0.00	9.88	
TXUN	1e-04	78	1	SEMI-BANTU	Central West Africa NC	0.13	0.00	0.32	1.33	2.59	0.02	0.02	0.02	TSI				5.37	5.90	4.81	2.25	0.00	2.17	
-KHOMANI	0.00048	8	1	BANTU	Central West Africa NC	0.55	0.00	0.93	1.56	5.30		0.55	0.55	GBR	15.25			15.26	15.03	10.47		0.00	10.47	
JU//HOANSI	1e-04	62	8	AMAXHOSA	South Africa NC	0.93	0.80	1.10	1.84	3.13	0.00	0.80	0.80	GBR				4.04	5.01	5.00	3.80	1.91	0.00	1.91

Figure 3-Source Data 3. The evidence for multiple waves of admixture in African populations using MALDER and the European recombination map. Columns as in Figure 3-Source Data 1

Ethnic Group	amp	Date Date (CI)	Popl	PoplAnc	PoplAnc-West Africa NC(Z)	PoplAnc-Central West Africa NC(Z)	PoplAnc-East Africa NC(Z)	PoplAnc-Nilo-Saharan(Z)	PoplAnc-Afroasiatic(Z)	PoplAnc-South Africa NC(Z)	PoplAnc-Khoesan(Z)	PoplAnc-Eurasia(Z)	PoplAnc-(Z)	Pop2	Pop2Anc	Pop2Anc-West Africa NC(Z)	Pop2Anc-Central West Africa NC(Z)	Pop2Anc-East Africa NC(Z)	Pop2Anc-Nilo-Saharan(Z)	Pop2Anc-Afroasiatic(Z)	Pop2Anc-South Africa NC(Z)	Pop2Anc-Khoesan(Z)	Pop2Anc-Eurasia(Z)	Pop2Anc-(Z)
JOLA	7.1e-05	152	27	JU//HOANSI	Central West Africa NC	0.43	0.39	0.92	1.09	2.97	0.21	0.00	0.21	GBR	Eurasia	3.10	3.10	3.10	3.10	3.10	3.10	3.10	3.10	3.10
WOLLOF	0.00011	96	24	SEMI-BANTU	Central West Africa NC	0.03	0.00	0.62	0.93	3.18	0.29	0.41	0.03	GBR	Eurasia	5.84	5.84	5.84	5.84	5.84	5.84	5.84	5.84	5.84
WOLLOF	2.4e-05	11	3	YORUBA	Central West Africa NC	0.45	0.00	0.67	1.01	3.10	0.18	0.31	0.45	TSI	Eurasia	5.73	5.73	5.73	5.73	5.73	5.73	5.73	5.73	5.73
MANDINKAI	3.9e-05	20	3	SEMI-BANTU	Central West Africa NC	0.51	0.00	1.04	0.99	4.40	0.23	0.19	0.19	IBS	Eurasia	8.38	8.38	8.38	8.38	8.38	8.38	8.38	8.38	8.38
SEREHULE	6e-05	23	4	BANTU	Central West Africa NC	0.43	0.00	0.60	1.05	3.94	0.07	0.07	0.07	IBS	Eurasia	7.41	7.41	7.41	7.41	7.41	7.41	7.41	7.41	7.41
FULAI	0.00043	61	8	SEMI-BANTU	Central West Africa NC	0.09	0.00	1.27	1.27	5.42	0.40	0.63	0.09	GBR	Eurasia	11.76	11.76	11.76	11.76	11.76	11.76	11.76	11.76	11.76
FULAI	0.00016	11	2	AKANS	Central West Africa NC	-0.20	0.00	0.80	0.94	4.22	0.03	0.50	-0.20	IBS	Eurasia	8.82	8.82	8.82	8.82	8.82	8.82	8.82	8.82	8.82
MALINKE	0.00012	123	31	SEREHULE	West Africa NC	0.00	0.30	0.23	0.64	0.87	0.22	0.27	0.22	PEL	Eurasia	1.38	1.38	1.38	1.38	1.38	1.38	1.38	1.38	1.38
MALINKE	1.9e-05	8	2	MALAWI	South Africa NC	-0.52	-0.13	-0.40	0.13	0.84	0.00	-0.67	-0.13	TSI	Eurasia	2.03	2.03	2.03	2.03	2.03	2.03	2.03	2.03	2.03
SERENJE	2.5e-05	15	3	BANTU	Central West Africa NC	0.68	0.00	0.75	1.74	6.45	0.41	0.25	0.25	GBR	Eurasia	7.98	7.98	7.98	7.98	7.98	7.98	7.98	7.98	7.98
MANJAGO	3.4e-05	3	1	YORUBA	Central West Africa NC	0.23	0.00	0.95	1.74	6.45	0.41	0.25	0.25	GBR	Eurasia	12.37	12.37	12.37	12.37	12.37	12.37	12.37	12.37	12.37
FULAI	2e-05	64	1	JU//HOANSI	Khoesan	0.11	0.21	1.13	2.11	5.37	0.38	0.00	0.11	GBR	Eurasia	9.78	9.78	9.78	9.78	9.78	9.78	9.78	9.78	9.78
FULAI	7e-05	5	10	SEMI-BANTU	Central West Africa NC	0.10	0.00	0.73	1.68	5.01	0.13	-0.08	0.13	TSI	Eurasia	9.51	9.51	9.51	9.51	9.51	9.51	9.51	9.51	9.51
BAMBARA	5.2e-05	86	29	YORUBA	Central West Africa NC	0.06	0.00	0.18	0.35	0.97	0.09	0.33	0.06	IBS	Eurasia	1.70	1.70	1.70	1.70	1.70	1.70	1.70	1.70	1.70
BAMBARA	1.3e-05	9	2	/GUI//GANA	Khoesan	-0.20	-0.56	-0.33	-0.01	0.57	-0.36	0.00	-0.56	FIN	Eurasia	1.48	1.48	1.48	1.48	1.48	1.48	1.48	1.48	1.48
AKANS	2.3e-05	138	35	MANDINKAI	West Africa NC	0.00							0.26	GBR	Eurasia	1.98	1.98	1.98	1.98	1.98	1.98	1.98	1.98	1.98
MOSSI																								
YORUBA																								
KASEM																								
NAMKAM																								
SEMI-BANTU	1e-05	46	12	AKANS	Central West Africa NC	0.23	0.00		1.10	1.00			0.23	JU//HOANSI	Khoesan	2.94	2.94	2.94	2.94	2.94	2.94	2.94	2.94	2.94
BANTU	2.1e-05	54	13	MOSSI	Central West Africa NC	0.01	0.00	0.41	0.54	1.79	2.07		0.01	JU//HOANSI	Khoesan	6.09	6.09	6.09	6.09	6.09	6.09	6.09	6.09	6.09
LUHYA	7.8e-05	59	18	NAMKAM	Central West Africa NC	0.12	0.00		0.61	3.90	0.09	0.51	0.09	TSI	Eurasia	4.88	4.88	4.88	4.88	4.88	4.88	4.88	4.88	4.88
LUHYA	2.6e-05	13	3	YORUBA	Central West Africa NC	0.06	0.00		0.56	3.58	-0.14	0.77	0.06	TSI	Eurasia	5.84	5.84	5.84	5.84	5.84	5.84	5.84	5.84	5.84
KAMBE	0.00014	15	2	BANTU	Central West Africa NC	0.20	0.00		1.95	5.82	1.10	0.94	0.20	TSI	Eurasia	9.29	9.29	9.29	9.29	9.29	9.29	9.29	9.29	9.29
CHONYI	0.00014	27	2	YORUBA	Central West Africa NC	0.05	0.00		3.21	8.89	1.91	1.50	0.05	TSI	Eurasia	13.70	13.70	13.70	13.70	13.70	13.70	13.70	13.70	13.70
KASUMBA	0.00015	27	2	SEMI-BANTU	Central West Africa NC	0.21	0.00		3.18	7.98	1.30		0.41	TSI	Eurasia	12.84	12.84	12.84	12.84	12.84	12.84	12.84	12.84	12.84
WASAMBA	0.00012	21	3	NAMKAM	Central West Africa NC	0.21	0.00		2.46	6.63	1.56		0.21	TSI	Eurasia	8.67	8.67	8.67	8.67	8.67	8.67	8.67	8.67	8.67
GRIAMA	0.00014	28	2	YORUBA	Central West Africa NC	0.58	0.00		3.89	8.81	2.05		0.58	TSI	Eurasia	12.37	12.37	12.37	12.37	12.37	12.37	12.37	12.37	12.37
WABONDEI	1e-04	23	2	NAMKAM	Central West Africa NC	0.36	0.00		3.78	9.03	1.59	2.31	0.36	TSI	Eurasia	11.60	11.60	11.60	11.60	11.60	11.60	11.60	11.60	11.60
MZIGUA	0.00011	32	4	YORUBA	Central West Africa NC	0.23	0.00		1.97	5.14	0.90	1.11	0.23	TSI	Eurasia	7.28	7.28	7.28	7.28	7.28	7.28	7.28	7.28	7.28
MAASAI	0.00039	81	17	SUDANESE	Nilo-Saharan	0.26	0.17		0.00	2.77	0.22	0.59	0.17	TSI	Eurasia	4.02	4.02	4.02	4.02	4.02	4.02	4.02	4.02	4.02
MAASAI	8.6e-05	11	2	SEMI-BANTU	Central West Africa NC	0.17	0.00		-0.23	2.85	0.12	0.43	0.12	TSI	Eurasia	4.24	4.24	4.24	4.24	4.24	4.24	4.24	4.24	4.24
SUDANESE	1.7e-05	6	2	/GUI//GANA	Khoesan	0.34	0.11		0.79	0.47	3.31	0.43	0.00	GBR	Eurasia	5.21	5.21	5.21	5.21	5.21	5.21	5.21	5.21	5.21
GUMUZ	1.8e-05	5	1	SUDANESE	Nilo-Saharan	0.50	0.25		0.64	0.00	2.23	0.59	0.24	GBR	Eurasia	5.58	5.58	5.58	5.58	5.58	5.58	5.58	5.58	5.58
GUMUZ	8.6e-05	104	26	SUDANESE	Nilo-Saharan	-0.55	-0.90		-0.07	0.00	1.22	-0.67	-0.99	PEL	Eurasia	6.14	6.14	6.14	6.14	6.14	6.14	6.14	6.14	6.14
ANUAK	3.3e-05	52	16	SUDANESE	Nilo-Saharan	1.56	1.53		1.82	0.00	1.45	2.49	0.53	JU//HOANSI	Khoesan	1.39	1.39	1.39	1.39	1.39	1.39	1.39	1.39	1.39
AFAR	0.00038	45	10	SUDANESE	Nilo-Saharan	0.23	0.09		0.54	0.00	2.01	0.17	0.37	TSI	Eurasia	4.84	4.84	4.84	4.84	4.84	4.84	4.84	4.84	4.84
OROMO	0.00033	86	9	JU//HOANSI	Khoesan	0.75	0.52		0.76	0.22	2.19	0.60	0.00	TSI	Eurasia	7.99	7.99	7.99	7.99	7.99	7.99	7.99	7.99	7.99
OROMO	5e-05	7	2	JU//HOANSI	Khoesan	0.75	0.52		0.76	0.22	2.19	0.60	0.00	TSI	Eurasia	8.02	8.02	8.02	8.02	8.02	8.02	8.02	8.02	8.02
SOMALI	0.00043	70	11	SUDANESE	Nilo-Saharan	0.30	0.04		0.76	0.86	2.19	0.65	0.00	TSI	Eurasia	5.26	5.26	5.26	5.26	5.26	5.26	5.26	5.26	5.26
WOLAYTA	0.00054	55	9	JU//HOANSI	Khoesan	0.88	0.69		1.00	1.46	3.32	0.59	0.00	CEU	Eurasia	4.79	4.79	4.79	4.79	4.79	4.79	4.79	4.79	4.79
ARI	0.00045	99	11	JU//HOANSI	Khoesan	0.83	0.67		1.03	1.46	3.32	0.59	0.00	TSI	Eurasia	5.17	5.17	5.17	5.17	5.17	5.17	5.17	5.17	5.17
AMHARA	0.00052	68	10	ANUAK	Nilo-Saharan	0.41	0.21		0.58	0.00	1.70	0.27	0.06	TSI	Eurasia	5.83	5.83	5.83	5.83	5.83	5.83	5.83	5.83	5.83
TIGRAY	0.00066	77	7	SUDANESE	Nilo-Saharan	0.65	0.44		0.86	0.00	2.83	0.47	0.44	TSI	Eurasia	9.02	9.02	9.02	9.02	9.02	9.02	9.02	9.02	9.02
MALAWI	4.1e-05	43	6	YORUBA	Central West Africa NC	0.42	0.00		1.53	1.34	1.12		0.51	JU//HOANSI	Khoesan									
SEBANTU	5.2e-05	18	4	BAMBARA	West Africa NC	0.00	0.25		1.74	5.17	1.83		0.25	GBR	Eurasia	6.78	6.78	6.78	6.78	6.78	6.78	6.78	6.78	6.78
HERERO	0.00024	2	0	/GUI//GANA	Khoesan	1.06	0.64		1.29	2.36	6.25	0.58	0.00	GBR	Eurasia	12.56	12.56	12.56	12.56	12.56	12.56	12.56	12.56	12.56
KHWE	0.00019	103	22	WASAMBA	East Africa NC	0.48	0.06		0.00	0.73	1.52		0.06	PEL	Eurasia	2.40	2.40	2.40	2.40	2.40	2.40	2.40	2.40	2.40
KHWE	4.4e-05	18	5	BANTU	Central West Africa NC	-0.04	0.00		0.09	0.48	1.04		-0.04	TSI	Eurasia	1.83	1.83	1.83	1.83	1.83	1.83	1.83	1.83	1.83
/GUI//GANA	3.8e-05	20	4	JOLA	West Africa NC	0.00	0.16		1.69	3.46	1.33		0.16	TSI	Eurasia	7.65	7.65	7.65	7.65	7.65	7.65	7.65	7.65	7.65
KARRETJIE	0.00047	5	0	BANTU	Central West Africa NC	0.76	0.00		2.33	6.97			0.76	GBR	Eurasia	19.78	19.78	19.78	19.78	19.78	19.78	19.78	19.78	19.78
NAMA	0.00027	5	0	MALAWI	South Africa NC	0.07	0.03		0.08	0.17	0.56	0.00	0.03	GBR	Eurasia	1.46	1.46	1.46	1.46	1.46	1.46	1.46	1.46	1.46
NAMA	0.00031	80	13	BANTU	Central West Africa NC	0.07	0.00		0.01	0.13	0.46	-0.01	0.07	FIN	Eurasia	1.24	1.24	1.24	1.24	1.24	1.24	1.24	1.24	1.24
XUN	0.00011	52	6	BANTU	Central West Africa NC	0.37	0.00		0.58	1.89	3.84	0.08	0.08	TSI	Eurasia	7.06	7.06	7.06	7.06	7.06	7.06	7.06	7.06	7.06
=KHOMANI	0.00049	6	0	BANTU	Central West Africa NC	0.73	0.00		2.10				0.73	GBR	Eurasia	20.41	20.41	20.41	20.41	20.41	20.41	20.41	20.41	20.41
JU//HOANSI	8.6e-05	41	5	MALAWI	South Africa NC	0.13	0.15	</																

Figure 3-Source Data 4. The evidence for multiple waves of admixture in African populations using MALDER and the HAPMAP recombination map and a mindis of 0.5cM. Columns are as in Figure 3-Source Data 1. Here we show the results for the MALDER analysis where we over-ride any short-range LD and define a minimum distance of 0.5cM from which to start computing admixture LD curves

Ethnic Group	amp	Date Date (CI)	Pop1	Pop1Anc	Pop1Anc-West Africa NC(Z)	Pop1Anc-Central West Africa NC(Z)	Pop1Anc-East Africa NC(Z)	Pop1Anc-Nilo-Saharan(Z)	Pop1Anc-Afroasiatic(Z)	Pop1Anc-South Africa NC(Z)	Pop1Anc-Khoesan(Z)	Pop1Anc-Eurasia(Z)	Pop1Anc-(Z)	Pop2	Pop2Anc	Pop2Anc-West Africa NC(Z)	Pop2Anc-Central West Africa NC(Z)	Pop2Anc-East Africa NC(Z)	Pop2Anc-Nilo-Saharan(Z)	Pop2Anc-Afroasiatic(Z)	Pop2Anc-South Africa NC(Z)	Pop2Anc-Khoesan(Z)	Pop2Anc-Eurasia(Z)	Pop2Anc-(Z)
JOLA	0.00016	225	24	MANDINKAI	West Africa NC	0.00	0.24	0.78	1.07	3.65	0.45	0.56	0.24	GBR	Eurasia	5.34	5.19	4.32	2.64	5.12	4.66	0.00	2.64	Pop2Anc-(Z)
WOLLOF	0.00019	141	16	JOLA	West Africa NC	0.00	0.36	1.83	2.89	8.97	0.90	0.78	0.36	TSI	Eurasia	13.75	14.09	12.53	7.70	13.63	12.55	0.00	7.70	Pop2Anc-Eurasia(Z)
WOLLOF	2.9e-05	12	3	JOLA	West Africa NC	0.00	0.55	1.83	2.90	8.99	0.90	1.29	0.55	TSI	Eurasia	13.89	14.13	12.53	7.70	13.94	13.50	0.00	7.70	Pop2Anc-Eurasia(Z)
MANDINKAI	7.5e-05	31	10	JOLA	West Africa NC	0.00	0.25	0.82	1.09	3.61	0.51	0.62	0.25	TSI	Eurasia	4.36	4.59	4.62	2.71	4.40	4.39	0.00	2.71	Pop2Anc-Khoesan(Z)
MANDINKAI	0.00016	170	37	JOLA	West Africa NC	0.00	0.15	0.59	1.03	3.61	0.25	0.37	0.15	IBS	Eurasia	5.25	5.42	4.62	2.65	4.76	4.58	0.00	2.65	Pop2Anc-Khoesan(Z)
MANDINKAI	2.6e-05	14	5	JOLA	West Africa NC	0.00	0.32	0.92	1.14	3.75	0.30	0.45	0.32	CEU	Eurasia	5.37	5.41	4.59	2.59	5.29	4.55	0.00	2.59	Pop2Anc-Khoesan(Z)
SERHULE	0.00018	116	19	JOLA	West Africa NC	0.00	0.28	1.03	1.96	6.12	0.50	0.83	0.28	TSI	Eurasia	9.70	9.95	8.98	5.73	9.64	9.42	0.00	5.73	Pop2Anc-Khoesan(Z)
SERHULE	2.6e-05	10	3	AKANS	Central West Africa NC	0.02	0.00	0.91	1.73	6.58	0.11	0.55	0.11	IBS	Eurasia	10.78	12.00	10.65	6.76	10.87	10.86	0.00	6.76	Pop2Anc-Khoesan(Z)
FULAI	6e-04	86	9	JOLA	West Africa NC	0.00	0.47	2.57	3.61	13.21	0.98	2.25	0.47	TSI	Eurasia	25.91	25.75	23.18	14.98	24.99	24.12	0.00	14.98	Pop2Anc-Khoesan(Z)
FULAI	0.00019	12	2	AKANS	Central West Africa NC	-0.50	0.00	2.08	3.08	12.59	0.51	1.76	-0.50	IBS	Eurasia	25.07	25.02	22.44	14.50	24.28	23.53	0.00	14.50	Pop2Anc-Khoesan(Z)
SERERE																								
MANJAGO	3.3e-05	3	0	AKANS	Central West Africa NC	0.08	0.00	0.97	1.80	6.48	0.23	0.92	0.08	GBR	Eurasia	12.26	12.97	12.00	8.15	12.87	12.31	0.00	8.15	Pop2Anc-Khoesan(Z)
MANJAGO	0.00024	279	40	JU//HOANSI	Khoesan	-1.13	-1.23	-0.25	0.73	5.76	-0.61	0.00	-1.13	TSI	Eurasia	12.86	13.37	13.41	12.78	9.42	12.92	0.00	9.42	Pop2Anc-Khoesan(Z)
FULAI	0.00013	106	11	AKANS	Central West Africa NC	0.04	0.00	1.03	1.60	6.88	0.40	0.92	0.04	TSI	Eurasia	10.92	12.61	10.98	6.87	11.37	11.18	0.00	6.87	Pop2Anc-Khoesan(Z)
FULAI	2.2e-05	6	1	JOLA	West Africa NC	0.00	-0.03	0.96	1.57	6.92	0.35	1.13	-0.03	GBR	Eurasia	10.70	11.25	11.07	6.48	10.75	10.04	0.00	6.48	Pop2Anc-Khoesan(Z)
BAMBARA	4.9e-05	37	11	AKANS	Central West Africa NC	0.33	0.00	0.50	0.81	2.70	0.13	0.21	0.13	TSI	Eurasia	3.76	4.81	4.12	2.52	4.45	4.41	0.00	2.52	Pop2Anc-Khoesan(Z)
AKANS	7.7e-05	174	38	JU//HOANSI	Khoesan	0.64	0.41	0.76	1.05	2.01	0.35	0.00	0.35	IBS	Eurasia	1.59	1.43	1.96	2.08	0.95	2.13	0.00	0.95	Pop2Anc-Khoesan(Z)
MOSSI	7.3e-05	118	26	BANTU	Central West Africa NC	0.09	0.00	0.40	0.80	2.53	0.12	0.13	0.09	GBR	Eurasia	3.50	3.94	3.22	2.01	3.45	3.23	0.00	1.78	Pop2Anc-Khoesan(Z)
YORUBA	0.00013	257	42	/GU//GANA	Khoesan	0.44	0.21	0.49	1.01	2.48	0.10	0.00	0.10	CEU	Eurasia	1.55	1.24	2.43	1.96	1.11	2.59	0.00	1.11	Pop2Anc-Khoesan(Z)
KASEM	0.00021	300	50	AMAXHOSA	South Africa NC	0.18	0.02	0.32	0.85	1.85	0.00	0.33	0.02	IBS	Eurasia	2.37	2.49	2.97	2.28	1.36	2.49	0.00	1.15	Pop2Anc-Khoesan(Z)
NAMKAM	0.00013	214	38	IXUN	Khoesan	0.22	0.11	0.48	0.62	2.14	0.00	0.00	0.00	TSI	Eurasia	3.11	3.05	3.22	2.77	1.47	3.33	0.00	1.47	Pop2Anc-Khoesan(Z)
SEMI-BANTU	1e-04	197	36	JU//HOANSI	Khoesan	1.60	1.30	1.55	1.50		0.89	0.00	0.89	IBS	Eurasia	1.34	1.05	1.38	1.74	0.87	2.71	0.00	0.87	Pop2Anc-Khoesan(Z)
BANTU	0.00017	311	49	JU//HOANSI	Khoesan	0.36	0.30	0.43		0.07	0.00		0.36	GBR	Eurasia	1.85	1.45	1.92	1.16	1.06	2.66	0.00	1.06	Pop2Anc-Khoesan(Z)
BANTU	2.1e-05	51	14	YORUBA	Central West Africa NC	0.57	0.00	0.27	0.15	0.14	1.39	-1.81	0.57	JU//HOANSI	Khoesan	3.01	2.62	2.05	1.25	1.02	0.00	-1.42	1.88	Pop2Anc-Khoesan(Z)
LUHYA	0.00013	103	24	JU//HOANSI	Khoesan	0.52	0.49	0.78	0.46	4.66	0.37	0.00	0.37	TSI	Eurasia	5.35	5.32	5.58	5.45	3.13	7.19	0.00	3.13	Pop2Anc-Khoesan(Z)
LUHYA	4.1e-05	16	3	MALAWI	South Africa NC	0.39	0.21	0.65	0.16	5.61	0.00	-0.06	0.21	TSI	Eurasia	10.24	9.26	11.20	9.63	5.17	9.85	0.00	5.33	Pop2Anc-Khoesan(Z)
KAMBE	0.00018	20	2	MALAWI	South Africa NC	0.98	0.32	0.88	3.69	9.41	0.00	1.23	0.32	TSI	Eurasia	14.36	14.30	15.34	13.51	8.35	14.33	0.00	8.00	Pop2Anc-Khoesan(Z)
CHONYI	0.00018	123	28	JU//HOANSI	Khoesan	0.94	0.75	1.06	1.90	4.24	0.21	0.00	0.21	IBS	Eurasia	4.59	4.60	4.48	4.96	3.02	5.73	0.00	3.02	Pop2Anc-Khoesan(Z)
CHONYI	0.00012	31	2	MALAWI	South Africa NC	0.46	0.31	0.56	1.23	3.88	0.00	-0.38	0.31	TSI	Eurasia	5.64	5.49	6.72	5.92	3.30	5.63	0.00	3.51	Pop2Anc-Khoesan(Z)
KUMA	0.00019	23	2	MALAWI	South Africa NC	1.52	0.67	1.28	5.87	14.76	0.00	1.75	0.67	TSI	Eurasia	21.44	21.18	23.06	19.97	12.10	21.42	0.00	11.46	Pop2Anc-Khoesan(Z)
WASAMBAA	0.00018	28	4	MALAWI	South Africa NC	1.08	0.50	0.95	3.70	8.45	0.00	1.54	0.50	TSI	Eurasia	11.02	10.81	12.27	10.10	5.83	10.99	0.00	5.83	Pop2Anc-Khoesan(Z)
GIRIAMA	0.00017	32	1	MALAWI	South Africa NC	2.33	0.78	2.11	9.54	22.80	0.00	3.41	0.78	TSI	Eurasia	31.59	31.09	34.24	29.08	17.33	31.40	0.00	16.28	Pop2Anc-Khoesan(Z)
WABONDEI	0.00016	33	3	MALAWI	South Africa NC	1.35	0.70	1.45	5.90	13.06	0.00	2.46	0.70	TSI	Eurasia	16.92	16.61	19.24	16.10	9.40	16.92	0.00	8.67	Pop2Anc-Khoesan(Z)
MZIGUA	0.00015	130	33	AMAXHOSA	South Africa NC	1.26	0.75	0.98	2.75	6.66	0.00	0.20	0.20	GBR	Eurasia	8.19	7.88	9.60	8.48	4.81	8.67	0.00	4.11	Pop2Anc-Khoesan(Z)
MZIGUA	8.8e-05	26	5	MALAWI	South Africa NC	0.50	0.24	0.71	1.54	4.32	0.00	-0.25	0.24	TSI	Eurasia	5.66	5.67	6.72	5.37	3.18	5.95	0.00	3.23	Pop2Anc-Khoesan(Z)
MAASAI	0.00053	104	9	SUDANESE	Nilo-Saharan	1.17	0.65	1.85	0.00	10.69	0.60	1.05	0.60	TSI	Eurasia	19.62	21.14	21.11	15.32	21.17	19.21	0.00	11.84	Pop2Anc-Khoesan(Z)
MAASAI	9.9e-05	11	2	MALAWI	South Africa NC	0.82	0.10	1.26	-0.60	9.60	0.00	0.42	0.10	TSI	Eurasia	18.61	18.96	20.17	19.72	11.30	18.97	0.00	10.88	Pop2Anc-Khoesan(Z)
SUDANESE																								
GUMUZ	1.9e-05	5	1	SUDANESE	Nilo-Saharan	0.61	0.33	0.68	0.00	2.65	0.68	0.30	0.30	GBR	Eurasia	6.71	6.96	6.95	5.24	7.21	6.25	0.00	4.16	Pop2Anc-Khoesan(Z)
GUMUZ	8.2e-05	110	25	SUDANESE	Nilo-Saharan	0.81	0.43	1.25	0.00	2.58	0.65	0.29	0.29	TSI	Eurasia	6.52	6.76	6.90	5.10	7.00	6.27	0.00	4.41	Pop2Anc-Khoesan(Z)
ANUAK																								
AFAR	0.001	166	31	ANUAK	Nilo-Saharan	0.59	0.37	0.64	0.00	2.55	0.44	0.56	0.37	TSI	Eurasia	6.03	6.28	6.40	5.12	6.33	6.20	0.00	4.15	Pop2Anc-Khoesan(Z)
AFAR	0.00013	21	6	SUDANESE	Nilo-Saharan	0.72	0.60	0.73	0.00	2.69	0.57	0.74	0.60	TSI	Eurasia	6.15	6.48	6.53	5.20	6.45	6.33	0.00	4.38	Pop2Anc-Khoesan(Z)
OROMO	0.00073	93	7	SUDANESE	Nilo-Saharan	1.58	1.05	2.14	0.00	7.99	1.10	0.42	0.42	TSI	Eurasia	22.22	22.84	23.25	18.51	22.96	22.13	0.00	15.00	Pop2Anc-Khoesan(Z)
OROMO	5.4e-05	7	2	JU//HOANSI	Khoesan	1.44	0.96	1.77	1.51	7.87	1.05	0.00	1.51	TSI	Eurasia	22.98	24.38	24.89	23.06	15.63	24.84	0.00	15.63	Pop2Anc-Khoesan(Z)
SOMALI	0.00075	100	5	SUDANESE	Nilo-Saharan	1.15	0.87</																	

Figure 3-Source Data 4 Continued from previous page

Ethnic Group	amp	Date	Date (CI)	Pop1	Pop1Anc	Pop1Anc-West Africa NC(Z)	Pop1Anc-Central West Africa NC(Z)	Pop1Anc-East Africa NC(Z)	Pop1Anc-Nilo-Saharan(Z)	Pop1Anc-Afroasiatic(Z)	Pop1Anc-South Africa NC(Z)	Pop1Anc-Khoesan(Z)	Pop1Anc-Eurasia(Z)	Pop2Anc(Z)	Pop2Anc-West Africa NC(Z)	Pop2Anc-Central West Africa NC(Z)	Pop2Anc-East Africa NC(Z)	Pop2Anc-Nilo-Saharan(Z)	Pop2Anc-Afroasiatic(Z)	Pop2Anc-South Africa NC(Z)	Pop2Anc-Khoesan(Z)	Pop2Anc-Eurasia(Z)	Pop2Anc(Z)
!XUN =KHOMANI =KHOMANI !XUN	0.00013	10	1	MOSSI	Central West Africa NC	0.23	0.00	0.82	-0.12	-1.61	8.78	-2.87	0.23	JU//HOANSI	10.14	10.37	10.14	9.62	4.91	0.00	7.78	4.91	
	0.00094	5	1	JU//HOANSI	Khoesan	4.80	4.52	4.84	5.27	6.81	1.90	0.00	1.90	GBR	6.96	8.06	8.03	7.10	4.74	9.79	0.00	4.74	
	0.00051	49	13	JU//HOANSI	Khoesan	6.12	5.75	6.14	6.45	7.77	1.93	0.00	1.93	TSI	6.81	7.94	7.96	6.94	4.53	9.75	0.00	4.53	
	0.00029	39	3	!XUN	Khoesan	11.56	11.12	11.64	12.76	14.32	2.67	0.00	2.67	TSI	7.98	9.92	10.1	7.52	3.67	15.73	0.00	3.67	

Figure 4-Source Data 1. Results of the main GLOBETROTTER analysis. *Analysis* refers to whether the main or masked analysis was used to produce the final result. Admixture P -values are based on 100 bootstrap replicates of the NULL procedure. Our resulting inferences, *res* can be: 1D (two admixing sources at a single date); 1MW (multiple admixing sources at a single date); 2D (admixture at multiple dates); NA (no-admixture); U (uncertain). $\max(R_1)$ refers to the R^2 goodness-of-fit for a single date of admixture, taking the maximum value across all inferred coancestry curves. FQ_1 is the fit of a single admixture event (i.e. the first principal component, reflecting admixture involving two sources) and FQ_2 is the fit of the first two principal components capturing the admixture event(s) (the second component might be thought of as capturing a second, less strongly-signalled event. M is the additional R^2 explained by adding a second date versus assuming only a single date of admixture; we use values above 0.35 to infer multiple dates (although see Supplementary Text for details). As well as the final result, for each event we show the inferred dates, α s and best matching sources for 1D, 1MW, and 2D inferences. Inferred dates are in years(+ 95% CI; B=BCE, otherwise CE); the proportion of admixture from the minority source (source 1) is represented by α . Date confidence intervals are based on 100 bootstrap replicates of the date inference

Cluster	Analysis	P	res	max(R ₁)	FQ ₁	FQ ₂	M	1D	1D α	1D src1	1D src2	1MW α	1MW src3	1MW src4	2D1	2D1 α	2D1 src1	2D1 src2	2D2	2D2 α	2D2 src1	2D2 src2
SEREHULE	main	<0.01	1D	0.943	1	1	0.36413 (514B-921)	0.11	GBR	JOLA												
MANDINKAI	main	<0.01	1D	0.925	1	1	0.3951B (1244B-1066)	0.1	GBR	JOLA												
WOLLOF	main	<0.01	1D	0.938	1	1	0.63254B (779B-355)	0.09	GBR	JOLA												
MANDINKAI	main	<0.01	1D	0.806	1	1	0.47312B (718B-416)	0.15	GBR	JOLA												
SERERE	main	<0.01	1D	0.885	1	1	0.41776B (1740B-254)	0.08	GBR	JOLA												
FULAI	main	<0.01	1D	0.969	1	1	0.82239 (199B-486)	0.19	IBS	WOLLOF												
BAMBARA	main	<0.01	1D	0.901	1	1	0.36152 (675B-906)	0.06	CEU	MALINKE												
FULAI	main	<0.01	1D	0.939	1	1	0.57765 (253B-876)	0.1	GBR	MALINKE												
MANJAGO	main	<0.01	1D	0.779	0.985	1	0.50413 (1590B-1718)	0.21	FULAI	JOLA												
JOLA	main	<0.01	1D	0.642	1	1	0.08834B (2287B-169)	0.17	FULAI	SERERE												
MALINKE	main	<0.01	1D	0.91	1	1	0.50326 (544B-865)	0.11	GBR	BAMBARA												
YORUBA	main	<0.01	1D	0.465	0.999	1	0.03848 (338-1184)	0.48	SEMI-BANTU	AKANS												
KASEM	main	<0.01	1D	0.305	0.997	1	0.03819 (456-1329)	0.1	SEMI-BANTU	MOSSI												
NAMKAM	main	<0.01	1D	0.197	0.996	0.990	0.2616 (184B-1197)	0.11	SEMI-BANTU	MOSSI												
SEMI-BANTU	main	<0.01	1D	0.43	0.999	1	0.03674 (192-1126)	0.2	MZIGUA	YORUBA												
BANTU	main	<0.01	1D	0.447	0.995	0.990	0.057 (617B-602)	0.34	MALAWI	YORUBA												
MOSSI	main	<0.01	1D	0.41	0.997	1	0.05355 (97B-952)	0.21	YORUBA	KASEM												
AKANS	main	<0.01	1MW	0.623	0.791	0.990	0.141399 (717-1675)	0.03	MALAWI	KASEM	0.29	KASEM	NAMKAM									
MAASAI	main	<0.01	2D	0.918	0.992	1	0.61								1660 (1573-1747)	0.06	TIGRAY	LUHYA	254B (764B-239)	0.35	AFAR	LUHYA
CHONYI	main	<0.01	1D	0.989	0.981	1	0.181138 (1080-1182)	0.08	KHV	WASAMBAA												
MZIGUA	main	<0.01	1D	0.988	0.999	1	0.191080 (1007-1138)	0.11	AFAR	WABONDEI												
KAMBE	main	<0.01	2D	0.989	0.987	1	0.46								1544 (1370-1776)	0.14	KAUMA	MZIGUA	761 (461B-1053)	0.06	GBR	MZIGUA
WABONDEI	main	<0.01	2D	0.985	0.998	0.990	0.48								1573 (1326-1834)	0.28	WASAMBAA	MZIGUA	703 (19-936)	0.1	TIGRAY	MZIGUA
LUHYA	main	<0.01	2D	0.993	0.998	0.990	0.59								1486 (1428-1573)	0.25	SUDANESE	MZIGUA	85 (400B-616)	0.29	WASAMBAA	MZIGUA
GRIAMA	main	<0.01	1D	0.989	0.999	1	0.211196 (1138-1254)	0.1	OROMO	MZIGUA												
WASAMBAA	main	<0.01	1D	0.988	1	1	0.341312 (1254-1341)	0.14	TIGRAY	MZIGUA												
KAUMA	main	<0.01	1D	0.992	0.982	1	0.291225 (1167-1254)	0.06	GIH	MZIGUA												
MALAWI	main, null	<0.01	1MW	0.894	0.967	0.970	0.12471 (340-631)	0.17	SEMI-BANTU	MZIGUA	0.16	AMAXHOSA	MZIGUA									
ANUAK	main	<0.01	1MW	0.575	0.961	0.970	0.04703 (427-1037)	0.17	YORUBA	SUDANESE	0.33	SUDANESE	SUDANESE									
ARI	main, null	<0.01	1D	0.866	0.995	0.990	0.0689B (965B-297B)	0.14	TSI	GUMUZ												
SUDANESE	main	<0.01	1MW	0.789	0.782	0.990	0.211341 (1225-1660)	0.27	GUMUZ	ANUAK	0.25	ANUAK	ANUAK									
GUMUZ	main	<0.01	1D	0.785	0.987	0.990	0.261544 (1384-1718)	0.24	ARI	ANUAK												
AFAR	main, null	<0.01	1D	0.895	1	1	0.14326 (7-587)	0.23	TSI	SOMALI												
OROMO	main, null	<0.01	2D	0.922	1	1	0.50								1834 (1674-1892)	0.17	IBS	WOLAYTA	283B (617B-51B)	0.27	TSI	ARI
SOMALI	main, null	<0.01	2D	0.939	1	1	0.40								1573 (870-1805)	0.04	TSI	WOLAYTA	863B (1887B-399B)	0.47	TIGRAY	GUMUZ
WOLAYTA	main, null	<0.01	1D	0.825	1	1	0.18268 (8B-602)	0.22	TSI	ARI												
TIGRAY	main, null	<0.01	1D	0.947	1	1	0.2036 (196B-240)	0.35	TSI	ARI												
Continued on next page																						

Continued on next page

[illegible]

Figure 4-Source Data 2. Results of the main GLOBETROTTER analysis. *Analysis* refers to whether the main or masked analysis was used to produce the final result. Admixture P -values are based on 100 bootstrap replicates of the NULL procedure. Our resulting inferences, *res* can be: 1D (two admixing sources at a single date); 1MW (multiple admixing sources at a single date); 2D (admixture at multiple dates); NA (no-admixture); U (uncertain). $\max(R_1)$ refers to the R^2 goodness-of-fit for a single date of admixture, taking the maximum value across all inferred coancestry curves. FQ_1 is the fit of a single admixture event (i.e. the first principal component, reflecting admixture involving two sources) and FQ_2 is the fit of the first two principal components capturing the admixture event(s) (the second component might be thought of as capturing a second, less strongly-signalised event. M is the additional R^2 explained by adding a second date versus assuming only a single date of admixture; we use values above 0.35 to infer multiple dates (although see Supplementary Text for details). As well as the final result, for each event we show the inferred dates, α s and best matching sources for 1D, 1MW, and 2D inferences. Inferred dates are in years(+ 95% CI; B=BCE, otherwise CE); the proportion of admixture from the minority source (source 1) is represented by α . Date confidence intervals are based on 100 bootstrap replicates of the date inference

Cluster	Analysis	P	$\max(R_1)$	FQ_1	FQ_2	M	1D	1D arc1	1D arc2	1MW	1MW arc3	1MW arc4	2D1	2D1 arc1	2D2	2D2 arc1	2D2 arc2
AFAR	main	<0.01	0.935	1	1	0.23558	0.22	TSI	SOMALI	0.43	AMHARA	OROMO	1196	TSI	WOLAYTA	1443B	ARI
AFAR	main.null	<0.01	0.895	1	1	0.14926	0.23	TSI	SOMALI	0.37	OROMO	AMHARA	1167	TSI	WOLAYTA	1530B	ARI
AKANS	main	<0.01	0.623	0.7910	0.9990	141399	0.03	MALAWI	KASEM	0.29	KASEM	NAMKAM	1805	MOSSI	NAMKAM	84	KASEM
AKANS	main.null	<0.01	0.253	0.9960	0.9980	06935	0.04	MALAWI	KASEM	0.1	MOSSI	NAMKAM	1805	MOSSI	NAMKAM	283B	MOSSI
AMAXHOSA	main	<0.01	0.986	1	1	0.431225	0.31	KARRETJIE	MALAWI	0.34	SEBANTU	SEBANTU	1312	KARRETJIE	MALAWI	1112B	MALAWI
AMAXHOSA	main.null	<0.01	0.982	1	1	0.49116	0.32	KARRETJIE	MALAWI	0.32	SEBANTU	SEBANTU	1312	KARRETJIE	MALAWI	2458B	MALAWI
AMHARA	main	<0.01	0.954	1	1	0.397	0.35	TSI	ARI	0.25	TIGRAY	AFAR	1573	TSI	OROMO	631B	ARI
AMHARA	main.null	<0.01	0.947	1	1	0.26138B	0.35	TSI	ARI	0.3	TIGRAY	TIGRAY	1631	TSI	OROMO	370B	ARI
ANUAK	main	<0.01	0.575	0.9610	0.9970	04703	0.17	YORUBA	SUDANESE	0.33	SUDANESE	SUDANESE	1892	GUMUZ	SUDANESE	471	SUDANESE
ANUAK	main.null	<0.01	0.449	0.9620	0.9960	05387	0.14	YORUBA	SUDANESE	0.25	GUMUZ	SUDANESE	1225	SUDANESE	SUDANESE	297	SUDANESE
ARI	main	<0.01	0.885	0.9940	0.9990	08602B	0.15	TSI	GUMUZ	0.24	WOLAYTA	SOMALI	1399	WOLAYTA	MAASAI	1123B	GUMUZ
ARI	main.null	<0.01	0.866	0.9950	0.9990	0689B	0.14	TSI	GUMUZ	0.21	WOLAYTA	SOMALI	1312	SOMALI	GUMUZ	1327B	GUMUZ
BAMBARA	main	<0.01	0.901	1	1	0.361370	0.11	GBR	MALINKE	0.3	SERERE	MALINKE	1747	FULAI	MALINKE	152	MALINKE
BAMBARA	main.null	<0.01	0.897	1	1	0.311341	0.1	GBR	MALINKE	0.22	SERERE	MALINKE	1718	FULAI	MALINKE	51B	MALINKE
BANTU	main	<0.01	0.447	0.9950	0.9990	05057	0.34	MALAWI	YORUBA	0.31	MALAWI	MALAWI	1776	MZIGUA	MALAWI	109B	YORUBA
BANTU	main.null	<0.01	0.588	0.998	1	0.02747B	0.5	MALAWI	YORUBA	0.29	MALAWI	MALAWI	1283	MZIGUA	MALAWI	312B	YORUBA
CHONYI	main	<0.01	0.989	0.981	1	0.181138	0.08	KHV	WASAMBAA	0.24	LUHYA	MALAWI	1254	MALAWI	WASAMBAA	1356B	MZIGUA
CHONYI	main.null	<0.01	0.988	0.975	1	0.241109	0.07	CDX	GIRIAMA	0.3	LUHYA	MALAWI	1283	CDX	WASAMBAA	254B	MZIGUA
FULAI	main	<0.01	0.969	1	1	0.82239	0.19	IBS	WOLLOF	0.44	WOLLOF	WOLLOF	1660	BAMBARA	IBS	1225newline(10026	SEREHUJI
FULAI	main.null	<0.01	0.966	1	1	0.82297	0.2	IBS	WOLLOF	0.4	WOLLOF	WOLLOF	1689	BAMBARA	WOLLOF	1254	SEREHUJI
FULAI	main	<0.01	0.939	1	1	0.371399	0.13	GBR	MALINKE	0.39	SERERE	BAMBARA	1660	FULAI	MALINKE	65	MALINKE
FULAI	main.null	<0.01	0.939	1	1	0.561399	0.13	GBR	MALINKE	0.14	SERERE	MALINKE	1660	FULAI	MALINKE	51B	MALINKE
GIRIAMA	main	<0.01	0.989	0.999	1	0.211196	0.1	OROMO	MZIGUA	0.18	SEMI-BANTU	MALAWI	1370	WASAMBAA	MZIGUA	210	MZIGUA
GIRIAMA	main.null	<0.01	0.989	1	1	0.281196	0.11	OROMO	MZIGUA	0.2	SEMI-BANTU	MALAWI	1399	WASAMBAA	MZIGUA	22B	MZIGUA
/GUT//GANA	main	<0.01	0.904	0.995	1	0.401544	0.25	MALAWI	KARRETJIE	0.15	KHWE	AMAXHOSA	1747	MALAWI	JU//HOANSI	877	KARRETJIE
/GUT//GANA	main.null	<0.01	0.865	0.995	1	0.331544	0.24	MALAWI	JU//HOANSI	0.12	KHWE	AMAXHOSA	1834	MALAWI	JU//HOANSI	935	KARRETJIE
GUMUZ	main	<0.01	0.785	0.9870	0.9990	261544	0.24	ARI	ANUAK	0.42	ANUAK	ANUAK	1747	ARI	ANUAK	1385B	ANUAK
GUMUZ	main.null	<0.01	0.748	0.9740	0.9960	141573	0.19	ARI	ANUAK	0.46	ANUAK	ANUAK	1718	ARI	ANUAK	1124B	ANUAK
JOLA	main	<0.01	0.642	1	1	0.08225B	0.18	FULAI	SERERE	0.49	MANDINKAI	MALINKE	1892	MANJAGO	WOLLOF	834B	SERERE
JOLA	main.null	<0.01	0.652	1	1	0.041907B	0.11	GBR	SERERE	0.47	MANDINKAI	FULAI	1834	SERERE	WOLLOF	2487B	SERERE
JU//HOANSI	main	<0.01	0.866	0.9740	0.9980	10732	0.15	SOMALI	KARRETJIE	0.33	NAMA	KARRETJIE	1892	NAMA	KARRETJIE	887	KARRETJIE
JU//HOANSI	main.null	<0.01	0.832	0.986	1	0.08558	0.11	SOMALI	KARRETJIE	0.48	IXUN	KARRETJIE	1805	KARRETJIE	IXUN	413	KARRETJIE
KAMBE	main	<0.01	0.989	0.987	1	0.461283	0.07	GBR	MZIGUA	0.34	LUHYA	MALAWI	1544	KAUMA	MZIGUA	761	MZIGUA
KAMBE	main.null	<0.01	0.988	0.986	1	0.491225	0.07	GBR	MZIGUA	0.32	LUHYA	MALAWI	1602	KAUMA	MZIGUA	790	MZIGUA

Continued on next page

Figure 4-Source Data 2 Continued from previous page

	Cluster	Analysis	P	res	max(R ₁)	FQ ₁	FQ ₂	M	ID	1D α_1	1D α_2	1D α_3	1D α_4	2D1 α	2D1 α_1	2D1 α_2	2D2 α	2D2 α_1	2D2 α_2				
	KARRETIJE	main	<0.01	1D(2D)	0.998	0.999	1	0.381776	0.1	GBR	/GUI//GANA	0.19	CDX	=KHOMANI	1805	(1747-1878)	0.1	GBR	/GUI//GANA	1602	(1123-1776)	0.15	AMAXHOSA =KHOMANI
	KARRETIJE	main.null	<0.01	1D	0.997	0.999	1	0.331776	0.1	GBR	/GUI//GANA	0.21	CDX	=KHOMANI	1805	(1747-1892)	0.09	GBR	/GUI//GANA	1660	(1092-1776)	0.37	AMAXHOSA =KHOMANI
	KASEM	main	<0.01	1D	0.305	0.997	1	0.03819	0.1	SEMI-BANTU	MOSSI	0.37	MOSSI	NAMKAM	1892	(1282-1892)	0.45	MOSSI	MOSSI	790	(10518B-1414)	0.07	SEMI-BANTU MOSSI
	KASEM	main.null	<0.01	1D	0.278	0.997	1	0.02645	0.16	YORUBA	MOSSI	0.37	MOSSI	NAMKAM	1399	(1138-1892)	0.45	AKANS	MOSSI	616	(6103B-1312)	0.06	MALAWI MOSSI
	KAUMA	main	<0.01	1D	0.992	0.982	1	0.291225	0.06	GIH	MZIGUA	0.42	LUHYA	MALAWI	1515	(1341-1878)	0.2	KAMBE	MZIGUA	874	(138B-1080)	0.07	GIH MZIGUA
	KAUMA	main.null	<0.01	1D	0.991	0.986	1	0.311196	0.08	GIH	MZIGUA	0.44	WASAMBAA	MALAWI	1573	(1370-1834)	0.17	KAMBE	MZIGUA	761	(166-1022)	0.08	GIH MZIGUA
	=KHOMANI	main	<0.01	1D	0.998	0.999	1	0.181776	0.13	CEU	KARRETIJE	0.27	HERERO	KARRETIJE	1776	(1747-1820)	0.13	CEU	KARRETIJE	312B	(356B-1762)	0.19	WOLAYTA KARRETIJE
	=KHOMANI	main.null	<0.01	1D	0.998	0.999	1	0.091776	0.13	CEU	KARRETIJE	0.27	HERERO	KARRETIJE	1776	(1747-1892)	0.13	CEU	KARRETIJE	1472B	(1142B-1776)	0.21	IBS KARRETIJE
	KHWE	main	<0.01	1D	0.928	0.9810.9980.21	1341	(1225-1428)	0.41	/GUI//GANASEMI-BANTU	SEBANTU	0.28	MZIGUA	SEBANTU	1486	(1370-1675)	0.43	/GUI//GANASEMI-BANTU	SEBANTU	457B	(2490B-703)	0.45	/GUI//GANASEMI-BANTU
	KHWE	main.null	<0.01	1D	0.911	0.9850.9990.16	1312	(1152-1399)	0.4	JU//HOANSI SEMI-BANTU	SEBANTU	0.27	MZIGUA	SEBANTU	1573	(1370-1820)	0.43	IXUN	MALAWI	268	(2009B-979)	0.42	/GUI//GANASEMI-BANTU
	LUHYA	main	<0.01	2D	0.993	0.9980.9990.59	1370	(1341-1428)	0.28	SUDANESE	MZIGUA	0.46	WASAMBAA	MZIGUA	1486	(1428-1573)	0.25	SUDANESE	MZIGUA	85	(400B-616)	0.29	WASAMBAA MZIGUA
	LUHYA	main.null	<0.01	2D	0.991	0.9980.9990.59	1341	(1341-1399)	0.29	SUDANESE	MZIGUA	0.44	WASAMBAA	MZIGUA	1486	(1442-1544)	0.25	SUDANESE	MZIGUA	123	(313B-486)	0.26	ANUAK MZIGUA
	MALAWI	main	<0.01	1MW	0.895	0.9290.9920.08	587	(413-761)	0.21	SEMI-BANTU	MZIGUA	0.16	SEBANTU	MZIGUA	1225	(703-1892)	0.41	MZIGUA	MZIGUA	167B	(2506B-500)	0.14	YORUBA MZIGUA
	MALAWI	main.null	<0.01	1MW	0.894	0.9670.9970.12	471	(340-631)	0.17	SEMI-BANTU	MZIGUA	0.16	AMAXHOSA	MZIGUA	1863	(1297-1892)	0.38	MZIGUA	MZIGUA	355	(21-457)	0.18	SEMI-BANTU MZIGUA
	MALINKE	main	<0.01	1D(2D)	0.91	1	1	0.501457	0.14	GBR	BAMBARA	0.31	SERERE	BAMBARA	1718	(1674-1834)	0.24	FULAI	FULAI	326	(544B-865)	0.11	GBR BAMBARA
	MALINKE	main.null	<0.01	1D(2D)	0.916	1	1	0.431370	0.12	GBR	BAMBARA	0.39	JOLA	BAMBARA	1718	(1645-1834)	0.23	FULAI	BAMBARA	355	(591B-747)	0.08	GBR BAMBARA
	MANDINKAI	main	<0.01	1D(2D)	0.806	1	1	0.471573	0.16	FULAI	JOLA	0.44	JOLA	SEREHULE	1805	(1747-1878)	0.19	FULAI	JOLA	312B	(718B-416)	0.15	GBR JOLA
	MANDINKAI	main.null	<0.01	1D(2D)	0.796	1	1	0.421544	0.16	FULAI	JOLA	0.49	SEREHULE	JOLA	1805	(1747-1863)	0.19	FULAI	JOLA	428B	(690B-256)	0.15	GBR JOLA
	MANDINKAI	main	<0.01	1D(2D)	0.925	1	1	0.391225	0.14	GBR	JOLA	0.3	JOLA	MALINKE	1631	(1486-1892)	0.2	FULAI	JOLA	51B	(1244B-1066)	0.1	GBR JOLA
	MANDINKAI	main.null	<0.01	1D(2D)	0.923	1	1	0.381196	0.19	FULAI	JOLA	0.43	JOLA	FULAI	1573	(1500-1892)	0.24	FULAI	JOLA	109B	(843B-1138)	0.17	GBR JOLA
	MANJAGO	main	<0.01	1D(2D)	0.779	0.985	1	0.501747	0.18	FULAI	JOLA	0.39	SERERE	SEREHULE	1892	(1790-1892)	0.23	FULAI	JOLA	413	(1590B-1718)	0.21	FULAI JOLA
	MANJAGO	main.null	<0.01	1D(2D)	0.766	0.989	1	0.481776	0.19	FULAI	JOLA	0.38	SERERE	WOLLOF	1892	(1805-1892)	0.17	FULAI	JOLA	326	(880B-1660)	0.2	FULAI SERERE
	MAASAI	main	<0.01	2D	0.918	0.992	1	0.611312	0.49	LUHYA	SOMALI	0.27	WABONDEI	LUHYA	1660	(1573-1747)	0.06	TIGRAY	LUHYA	254B	(764B-239)	0.35	AFAR LUHYA
	MAASAI	main.null	<0.01	2D	0.911	0.992	1	0.561254	0.49	SOMALI	LUHYA	0.27	WABONDEI	LUHYA	1631	(1515-1747)	0.06	TIGRAY	LUHYA	109B	(603B-297)	0.37	SOMALI LUHYA
	MOSSI	main	<0.01	1D	0.41	0.997	1	0.05855	0.21	YORUBA	KASEM	0.2	AKANS	KASEM	1892	(1485-1892)	0.27	KASEM	KASEM	42	(840B-979)	0.15	SEMI-BANTU KASEM
	MOSSI	main.null	<0.01	1D	0.417	0.998	1	0.0465	0.28	YORUBA	BAMBARA	0.18	NAMKAM	KASEM	1892	(1399-1892)	0.24	KASEM	KASEM	181	(501B-776)	0.2	SEMI-BANTU KASEM
	MZIGUA	main	<0.01	1D	0.988	0.999	1	0.191080	0.11	AFAR	WABONDEI	0.32	LUHYA	MALAWI	1196	(1109-1733)	0.13	WASAMBAA	WABONDEI	370B	(1156B-993)	0.12	WASAMBAA WABONDEI
	MZIGUA	main.null	<0.01	1D	0.986	0.999	1	0.181022	0.09	AFAR	WABONDEI	0.23	SEMI-BANTU	MALAWI	1341	(1108-1791)	0.17	WASAMBAA	WABONDEI	268	(385B-863)	0.1	AFAR WABONDEI
	NAMA	main	<0.01	1D(2D)	0.992	0.982	1	0.701805	0.3	HERERO	=KHOMANI	0.21	CEU	JU//HOANSI	1834	(1834-1863)	0.31	HERERO	=KHOMANI	210	(152-935)	0.15	HERERO =KHOMANI
	NAMA	main.null	<0.01	1MW(2D)	0.991	0.93	0.9990.52	1805	(1805-1834)	0.12	CEU	JU//HOANSI	HERERO	=KHOMANI	1834	(1805-1892)	0.11	CEU	JU//HOANSI	373B	(425B-1153)	0.29	GBR KARRETIJE
	NAMKAM	main	<0.01	1D	0.197	0.9900.9990.02	616	(184B-1197)	0.11	SEMI-BANTU	MOSSI	0.15	MOSSI	KASEM	1892	(1458B-1196)	0.4	AKANS	MOSSI	268	(1458B-1196)	0.12	SEMI-BANTU MOSSI
	NAMKAM	main.null	<0.01	1D	0.18	0.999	1	0.03123	0.17	SEMI-BANTU	MOSSI	0.2	MOSSI	KASEM	1892	(920-1892)	0.17	YORUBA	MOSSI	515B	(2417B-820)	0.28	YORUBA MOSSI
	OROMO	main	<0.01	2D	0.931	1	1	0.61355	0.28	TSI	ARI	0.2	AMHARA	AFAR	1805	(1660-1892)	0.15	IBS	WOLAYTA	312B	(690B-36B)	0.28	TSI ARI
	OROMO	main.null	<0.01	2D	0.922	1	1	0.50239	0.27	TSI	ARI	0.21	AMHARA	AFAR	1834	(1674-1892)	0.17	IBS	WOLAYTA	283B	(617B-51B)	0.27	TSI ARI
	SEBANTU	main	<0.01	1D(2D)	0.986	1	1	0.551167	0.29	KARRETIJE	MALAWI	0.28	AMAXHOSA	AMAXHOSA	1254	(1196-1399)	0.27	KARRETIJE	MALAWI	2951B	(3604B-848B)	0.31	KARRETIJE MALAWI
	SEBANTU	main.null	<0.01	1D(2D)	0.982	1	1	0.571109	0.3	KARRETIJE	MALAWI	0.34	AMAXHOSA	AMAXHOSA	1254	(1196-1641)	0.28	KARRETIJE	MALAWI	2777B	(4025B-390)	0.31	KARRETIJE MALAWI
	SEMI-BANTU	main	<0.01	1D	0.43	0.999	1	0.03874	0.2	MZIGUA	YORUBA	0.36	YORUBA	YORUBA	1167	(192-1126)	0.19	MZIGUA	YORUBA	7	(1153B	0.42	MALAWI YORUBA
	SEMI-BANTU	main.null	<0.01	1D	0.535	1	1	0.02167B	0.27	MZIGUA	YORUBA	0.48	YORUBA	YORUBA	NA	(965B-474)	0.11	MZIGUA	YORUBA	NA	1153B	0.27	MZIGUA YORUBA
	SEREHULE	main	<0.01	1D(2D)	0.943	1	1	0.361109	0.12	GBR	JOLA	0.46	MANJAGO	BAMBARA	1689	(514B-921)	0.25	FULAI	SERERE	413	(514B-921)	0.11	GBR JOLA
	SEREHULE	main.null	<0.01	1D	0.94	1	1	0.351080	0.09	GBR	JOLA	0.41	JOLA	MALINKE	1544	(1204B-822)	0.22	FULAI	JOLA	22B	(1740B-254)	0.08	GBR JOLA
	SERERE	main	<0.01	1D(2D)	0.885	1	1	0.411080	0.14	GBR	JOLA	0.47	MANJAGO	FULAI	1602	(761-1356)	0.24	FULAI	JOLA	776B	(1740B-254)	0.08	GBR JOLA
	SERERE	main.null	<0.01	1D(2D)	0.878	1	1	0.381051	0.14	GBR	JOLA	0.34	FULAI	MANJAGO	1602	(1956B-211)	0.24	FULAI	JOLA	950B	(1956B-211)	0.08	GBR JOLA
																							Continued on next page

Continued on next page

Figure 4-Source Data 2 Continued from previous page

Figure 4-Source Data 2 Continued from previous page																				
Cluster	Analysis	P	res	max(R ₁) FQ ₁ FQ ₂	M	ID	1D α 1	1D α 1	1D α 2	1MW α 1	1MW α 2	1MW α 3	2D1	2D1 α	2D1 α 1	2D1 α 2	2D2	2D2 α	2D2 α 1	2D2 α 2
SOMALI	main	<0.01	2D	0.933	1	0.60268 (36-529)	0.38	ANUAK	TIGRAY	0.12	WOLAYTA	WOLAYTA	1573 (1370-1791)	0.18	MAASAI	WOLAYTA	921B (1458B-413B)	0.46	TIGRAY	GUMUZ
SOMALI	main.null	<0.01	2D	0.939	1	0.40836 (196B-268)	0.39	ANUAK	TIGRAY	0.07	WASAMBAA	WOLAYTA	1573 (876-1805)	0.04	TSI	WOLAYTA	863B (1887B-399B)	0.47	TIGRAY	GUMUZ
SUDANESE	main	<0.01	1MW	0.789	0.7820	0.9990.211341 (1225-1660)	0.27	GUMUZ	ANUAK	0.25	ANUAK	ANUAK	1660 (1469-1892)	0.36	ANUAK	ANUAK	254B (979B-1284)	0.28	GUMUZ	ANUAK
SUDANESE	main.null	<0.01	1MW	0.658	0.8560	0.9990.131138 (1196-1689)	0.31	GUMUZ	ANUAK	0.19	ANUAK	ANUAK	1892 (1674-1892)	0.15	ANUAK	ANUAK	790 (8B-1515)	0.23	GUMUZ	ANUAK
HERERO	main	<0.01	1D(2D)	0.851	0.998	1 0.461631 (1747-1863)	0.41	SEMI-BANTU	NAMA	0.19	KAMBE	AMAXHOSA	1834 (1834-1892)	0.26	NAMA	AMAXHOSA	558 (298B-935)	0.43	NAMA	MALAWI
HERERO	main.null	<0.01	1D	0.816	0.998	1 0.341631 (1718-1863)	0.41	SEMI-BANTU	NAMA	0.19	KAMBE	AMAXHOSA	1834 (1805-1892)	0.24	NAMA	AMAXHOSA	874 (124B-979)	0.44	NAMA	SEMI-BANTU
TIGRAY	main	<0.01	1D	0.956	1	0.26152 (51B-370)	0.32	TSI	ARI	0.31	AMHARA	AMHARA	1341 (819-1776)	0.21	IBS	OROMO	802B (1968B-138B)	0.32	TSI	ARI
TIGRAY	main.null	<0.01	1D	0.947	1	0.20836 (196B-240)	0.35	TSI	ARI	0.42	AMHARA	AMHARA	1573 (862-1892)	0.16	IBS	AFAR	399B (1562B-94B)	0.35	TSI	ARI
WABONDEI	main	<0.01	2D	0.985	0.9980	0.9990.481138 (1080-1240)	0.11	TIGRAY	MZIGUA	0.5	MALAWI	WASAMBAA	1573 (1326-1834)	0.28	WASAMBAA	MZIGUA	703 (19-936)	0.1	TIGRAY	MZIGUA
WABONDEI	main.null	<0.01	2D	0.983	0.9960	0.9990.531109 (1051-1196)	0.1	AFAR	MZIGUA	0.49	WASAMBAA	MALAWI	1573 (1266-1820)	0.28	WASAMBAA	MZIGUA	587 (52B-863)	0.1	AFAR	MZIGUA
WASAMBAA	main	<0.01	1D	0.988	1	0.341312 (1254-1341)	0.14	TIGRAY	MZIGUA	0.3	LUHYA	MALAWI	1370 (1341-1834)	0.12	TIGRAY	MZIGUA	631B (1226B-1138)	0.16	WOLAYTA	MZIGUA
WASAMBAA	main.null	<0.01	1D	0.988	1	0.341254 (1225-1341)	0.15	AMHARA	MZIGUA	0.29	LUHYA	MALAWI	1486 (1384-1863)	0.15	AMHARA	MZIGUA	210 (283B-1066)	0.13	OROMO	MZIGUA
WOLAYTA	main	<0.01	1D	0.891	1	0.29268 (138B-602)	0.22	TSI	ARI	0.26	OROMO	SOMALI	1312 (1138-1834)	0.14	TSI	SOMALI	1182B (1893B-144)	0.27	TSI	ARI
WOLAYTA	main.null	<0.01	1D	0.825	1	0.18268 (8B-602)	0.22	TSI	ARI	0.26	OROMO	SOMALI	1457 (1137-1892)	0.15	TIGRAY	SOMALI	573B (1330B-388)	0.24	TSI	ARI
WOLLOF	main	<0.01	1D(2D)	0.938	1	0.631225 (1138-1341)	0.13	GBR	JOLA	0.4	MALINKE	JOLA	1631 (1544-1733)	0.19	FULAI	JOLA	254B (779B-355)	0.09	GBR	JOLA
WOLLOF	main.null	<0.01	1D(2D)	0.942	1	0.581167 (1080-1283)	0.12	GBR	JOLA	0.36	FULAI	MANDINKAI	1573 (1515-1762)	0.19	FULAI	JOLA	399B (951B-370)	0.08	GBR	JOLA
XUN	main	<0.01	1D(2D)	0.957	0.999	1 0.391341 (1254-1399)	0.27	SEMI-BANTU	JU//HOANSI	0.08	SOMALI	/GUI//GANA	1602 (1312-1892)	0.21	SEMI-BANTU	JU//HOANSI	819 (1008B-1196)	0.17	SEMI-BANTU	JU//HOANSI
XUN	main.null	<0.01	1D	0.951	0.999	1 0.211312 (1254-1385)	0.28	SEMI-BANTU	JU//HOANSI	0.07	SOMALI	/GUI//GANA	1805 (1341-1892)	0.15	SEMI-BANTU	JU//HOANSI	1080 (329B-1167)	0.27	SEMI-BANTU	JU//HOANSI
YORUBA	main	<0.01	1D	0.465	0.999	1 0.03848 (338-1184)	0.48	SEMI-BANTU	AKANS	0.29	AKANS	AKANS	1167 (1036-1892)	0.49	AKANS	SEMI-BANTU	1501B (2896B-1182)	0.25	MOSSI	SEMI-BANTU
YORUBA	main.null	<0.01	1D	0.498	1	0.04855 (429B-761)	0.37	SEMI-BANTU	AKANS	0.26	AKANS	AKANS	1109 (905-1892)	0.5	SEMI-BANTU	AKANS	2342B (4107B-778)	0.16	BANTU	AKANS

Supplementary Tables

Supplementary Table 1. Evidence for admixture across the ethnic groups included in the study using f_3 tests and ALDER. For each group, we report the f_3 test with the most negative value. Source.1. f_3 and Source.2. f_3 refer to the two groups between which gene flow must have occurred in the ancestors of the ethnic group under consideration. Groups without results are those where no negative f_3 statistic was found are shown with no results. For the ALDER analysis, we show the two sources of admixture, Source.2.LD and Source.1.LD, with the lowest reported admixture P -values. Dates for the ALDER events involving these groups are shown in generations, Date.gens, and in years, Date.

Ethnic.Group	Source.1. f_3	Source.2. f_3	f_3 (Z-score)	Source.1.LD	Source.2.LD	P.value.LD	Date.gens	Date
JOLA								
MANJAGO	JOLA	GBR	-11.091	MAASAI	GIH	5.3e-08	3.27 +/- 0.52	1826CE (1811-1841CE)
SEREHULE	JOLA	TSI	-17.778	WOLAYTA	IBS	1.6e-05	30.98 +/- 5.81	1023CE (854-1191CE)
SERERE	JOLA	TSI	-12.360	SOMALI	TSI	0.013	21.61 +/- 5.46	1294CE (1136-1453CE)
WOLLOF	JOLA	TSI	-23.384	ARI	IBS	6.3e-15	23.00 +/- 2.73	1254CE (1175-1333CE)
MANDINKAI	JOLA	FULAI	-20.677					
MANDINKAII	JOLA	FULAI	-18.910	MZIGUA	AFAR	1.1e-09	24.15 +/- 3.51	1221CE (1119-1322CE)
FULAI	JOLA	TSI	-55.429	ARI	AMAXHOSA	2e-21	16.06 +/- 1.21	1455CE (1420-1490CE)
FULAI	FULAI	YORUBA	-25.483	MAASAI	AFAR	1.8e-10	8.69 +/- 1.22	1669CE (1634-1704CE)
BAMBARA	FULAI	AKANS	-10.291	MANDINKAII	IBS	0.049	15.27 +/- 4.21	1478CE (1356-1600CE)
MALINKE	FULAI	AKANS	-8.553					
MOSSI	AKANS	TSI	-3.905					
AKANS	NAMIKAM	JU/'HOANSI	-9.038					
KASEM	KASEM	CHONYI	-1.890					
NAMKAM	NAMKAM	JU/'HOANSI	-2.862					
YORUBA	MOSSI	JU/'HOANSI	-13.894					
BANTU	MOSSI	JU/'HOANSI	-14.844					
SEMI-BANTU	NAMKAM	JU/'HOANSI	-7.927					
MALAWI	MALAWI	TIGRAY	-32.492	YORUBA	JU/'HOANSI	5.9e-07	38.47 +/- 5.71	805CE (640-971CE)
LUHYA	MALAWI	TSI	-14.987	MANJAGO	JU/'HOANSI	1.2e-12	42.03 +/- 5.26	702CE (550-855CE)
CHONYI	MALAWI	TSI	-34.562	OROMO	TSI	4.6e-25	30.42 +/- 2.81	1039CE (957-1120CE)
KAUMA	MALAWI	TSI	-44.547	GUMUZ	CEU	2.4e-27	29.55 +/- 2.43	1064CE (994-1135CE)
KAMBE	MALAWI	TSI	-21.929	WOLAYTA	IBS	6.7e-17	15.07 +/- 1.69	1484CE (1435-1533CE)
GIRIAMA	MALAWI	TSI	-34.835	ARI	AMHARA	1.4e-08	23.28 +/- 3.16	1246CE (1154-1338CE)
MZIGUA	MALAWI	TSI	-42.835	MANDINKAI	IBS	6.9e-21	32.96 +/- 3.33	965CE (869-1062CE)
WABONDEI	MALAWI	TSI	-50.513	ANUAK	CEU	9.8e-20	30.71 +/- 3.19	1030CE (938-1123CE)
WASAMBAA	MALAWI	TIGRAY	-66.343	KASEM	ARI	1.9e-17	15.16 +/- 1.67	1481CE (1433-1530CE)
MAASAI	SUDANESE	TSI	-67.795	GUMUZ	GBR	9.8e-21	31.33 +/- 3.17	1012CE (920-1104CE)
SUDANESE	SUDANESE	TSI	-23.214	MANDINKAI	TIGRAY	0.036	7.76 +/- 2.10	1696CE (1635-1757CE)
SOMALI	JU/'HOANSI	TSI	-5.653	GUMUZ	FIN	3.7e-11	68.34 +/- 9.14	61BCE (326B-204CE)
ARI	SUDANESE	ARI	-64.721	SOMALI	TSI	0.00019	5.78 +/- 1.19	1753CE (1719-1788CE)
ANUAK	SUDANESE	TSI	-86.578	SEMI-BANTU	GBR	7e-31	79.21 +/- 6.59	376BCE (567B-185BCE)
GUMUZ	SUDANESE	TSI	-87.571	BAMBARA	GBR	5.6e-18	76.91 +/- 7.81	309BCE (536B-83BCE)
AFAR	SUDANESE	TSI	-66.705	/GUI //GANA	CEU	6e-08	50.92 +/- 8.12	444CE (209-680CE)
TIGRAY	ANUAK	TSI	-90.232	BANTU	GBR	1.5e-45	59.19 +/- 4.08	204CE (86-323CE)
AMHARA	SUDANESE	FIN	-5.953	AMHARA	GBR	9.1e-08	2.29 +/- 0.37	1855CE (1844-1865CE)
WOLAYTA	MALAWI	CEU	-81.913	SUDANESE	JPT	8e-32	5.43 +/- 0.45	1764CE (1750-1777CE)
OROMO	JU/'HOANSI	KARRETJIE	-46.315					
HERERO	MOSSI	JU/'HOANSI	-39.259	ARI	WOLAYTA	0.019	27.01 +/- 6.35	1138CE (954-1322CE)
NAMA	MOSSI	GBR	-46.547	SUDANESE	CDX	2.3e-28	4.88 +/- 0.42	1779CE (1767-1792CE)
SEBANTU	JU/'HOANSI	FIN	-86.201	NAMKAM	CHS	6.8e-47	5.54 +/- 0.38	1760CE (1749-1771CE)
AMAXHOSA	MOSSI	FIN	-55.836	SEREHULE	SUDANESE	0.002	21.99 +/- 5.02	1283CE (1138-1429CE)
KARRETJIE	JU/'HOANSI	JU/'HOANSI	-55.053	SUDANESE	ARI	0.021	7.92 +/- 2.06	1691CE (1632-1751CE)
#KHOMANI	YORUBA	JU/'HOANSI	-36.212	ARI	GIH	2.5e-07	34.11 +/- 5.63	932CE (769-1095CE)
KHWE	MALAWI	JU/'HOANSI						
'XUN								
/GUI //GANA								
JU/'HOANSI								