

1     **Gene-dense autosomal chromosomes show evidence for increased selection**

2

3     M. Reza Jabalameli, Clare Horscroft, Alejandra Vergara Lope, \*Reuben J. Pengelly, \*Andrew Collins

4     \*Contributed equally

5

6     Genetic Epidemiology and Bioinformatics,

7     Faculty of Medicine,

8     University of Southampton,

9     Duthie Building (808),

10    Tremona Road, Southampton, SO16 6YD, UK.

11    Email: [arc@soton.ac.uk](mailto:arc@soton.ac.uk)

12    Running title: Increased selection on gene dense autosomes

13    Word count: 3789

## 14    **Abstract**

15    Purifying selection tends to reduce nucleotide and haplotype diversity leading to increased linkage  
16    disequilibrium. However, detection of evidence for selection is difficult as the signature is confounded  
17    by wide variation in the recombination rate which has a complex relationship with selection. The  
18    effective bottleneck time (the ratio of the linkage disequilibrium map to the genetic map in Morgans)  
19    controls for variability in recombination rate. Reduced effective bottleneck times indicate stronger  
20    residual linkage disequilibrium, consistent with increased selection. Using whole genome sequence  
21    data from one European and three Sub-Saharan African human populations we find, in the African  
22    samples, strong correlations between high gene densities and reduced effective bottleneck time for  
23    autosomal chromosomes. This suggests that gene-dense autosomes have been subject to increased  
24    purifying selection reducing effective bottleneck times compared to gene-poor autosomes. Although  
25    previous studies have shown unusually strong linkage disequilibrium for the sex chromosomes  
26    variation within the autosomes has not been recognised. The strongest relationship is between  
27    effective bottleneck time and the density of essential genes, which are likely targets of greater  
28    selective pressure ( $p = 0.006$ , for the 22 autosomes). The magnitude of the reduction in chromosome-  
29    specific effective bottleneck times from the least to the most gene-dense autosomes is ~17-21% for  
30    Sub-Saharan African populations. The effect size is greater in Sub-Saharan African populations,  
31    compared to a European sample, consistent with increased efficiency of selection in populations with  
32    larger effective population sizes which have not been subject to intense population bottlenecks as  
33    experienced by populations of European ancestry. The findings highlight the value of deeper analyses  
34    of selection within Sub-Saharan African populations.

35

## 36 **Introduction**

37 An understanding of the impact of selection on genomes is important for the interpretation of  
38 population history, genome function and disease genomics. Patterns of linkage disequilibrium (LD)  
39 have been used extensively to find signatures of positive and purifying selection by recognition of  
40 reduced diversity (Voight *et al.*, 2006; Huff *et al.*, 2010). Tests include the identification of long-  
41 range haplotypes through, for example, extended haplotype homozygosity (EHH) which relies on the  
42 relationship between the frequency of an allele and the extent of linkage disequilibrium (LD)  
43 surrounding it. In such regions a pattern of unusually long-range LD, given the allele's population  
44 frequency, is found after correcting for the recombination rate (Sabeti *et al.*, 2002).

45 Signatures of increased selection have also been demonstrated for whole chromosomes, specifically  
46 the sex chromosomes. The unique mode of inheritance of the X chromosome is likely to expose it to  
47 more intense selective pressure compared to the autosomal chromosomes. There is evidence that the  
48 X chromosome has been subject to strong selective sweeps (Dutheil *et al.*, 2015) reducing nucleotide  
49 diversity around genes compared to autosomal chromosomes. Megabase-size regions spanning about  
50 one third of the X chromosome show a reduction in nucleotide diversity to less than 20% of the  
51 chromosome average (Nam *et al.*, 2015). Several very strong selective sweeps are considered to have  
52 operated independently in these regions and testis-expressed ampliconic genes may have been the  
53 target: a mechanism to correct for distortions in the sex ratio is proposed (Dutheil *et al.*, 2015).  
54 Gottipati *et al.* (2011) concluded that diversity on chromosome X has been impacted to a greater  
55 extent by selection at linked sites than the autosomes, probably through the increased exposure in  
56 males when damaging X-linked recessive variants are present in single copy. Veeramah *et al.* (2014)  
57 present evidence indicating that the X chromosome has been subject to elevated levels of both  
58 purifying and positive selection.

59 The Y chromosome also has unusually low levels of genetic diversity. Because most of the Y  
60 chromosome does not recombine all sites are effectively linked and selection operating on any  
61 individual site will affect all other sites. The Y chromosome is subject to purifying selection removing

62 new deleterious mutations and reducing diversity at linked neutral sites (Wilson Sayres *et al.*, 2014).  
63 Similarly, positive selection, operating on any beneficial mutations may decrease diversity at linked  
64 neutral sites. The Y chromosome is gene-poor and the number of sites which are subject to purifying  
65 selection exceeds the number of Y-linked coding sites (there are ~100,000 single copy coding sites).  
66 The highly repetitive ampliconic sequences span a much larger region (5.7 Mb) and genes in these  
67 regions are expressed exclusively in the testis and possibly subject to selection in relation to male  
68 fertility. The ampliconic region plays an important role in the restoration of mutation-free gene copies  
69 through intrachromosomal gene conversion (Bachtrog, 2013).

70 The complex relationships between recombination and natural selection have received considerable  
71 attention over many years. It is known that selection impacts genetic diversity over greater physical  
72 distances in chromosome regions which have lower rates of recombination, through ‘hitch-hiking’  
73 (Smith and Haigh, 1974; Hudson and Kaplan, 1995). However, recombination also improves the  
74 efficiency of selection whenever multiple linked loci are simultaneously under selective pressure (Hill  
75 and Robertson, 1966; Castellano *et al.*, 2016; Pengelly *et al.*, 2019). If selection is more intense for  
76 genic sequences then genomic regions with a high density of coding sequences will contain more  
77 potential targets of selection compared to regions with low gene density. After controlling for the  
78 recombination rate, which varies widely along the length and between chromosomes, reductions in  
79 nucleotide (and hence haplotype) diversity in these regions suggest positive or purifying selection  
80 (Payseur and Nachman, 2002; Cutter and Payseur, 2013). Therefore it is expected that longer blocks  
81 of strong LD will surround substitutions at sites subject to selection. Autosomal chromosomes show  
82 very marked variations in gene density (for example, chromosome 19 has ~24 genes/Mb, compared to  
83 chromosome 13 which has only ~3 genes/Mb (Spataro *et al.*, 2017)). Payseur and Nachman (2002)  
84 show that the number of genes in a genomic region is a reasonable predictor of selection intensity.

85 Covariance between gene density and recombination rates can obscure the signatures of linked  
86 selection (Cutter and Payseur, 2013). However, a useful measure is the effective bottleneck time (EBT)  
87 defined as the ratio of the LD unit (LDU) map length to the linkage map length in Morgans (Zhang *et al.*  
88 *et al.*, 2004). The EBT corresponds to the number of generations since an “effective” bottleneck (a

single bottleneck which reflects the combined impact of multiple bottlenecks on haplotype diversity over time). For example, Tapper et al (2005) determined that the EBT for the X chromosome in the HapMap CEU population (representing Utah Residents with Northern and Western European ancestry) was substantially reduced compared to the autosomes consistent with increased purifying selection reducing diversity. Analysis of whole genome sequence (WGS) data from European and Nigerian population samples demonstrated a substantially higher resolution of LD structure compared to array-based panels (Pengelly *et al.*, 2015). Using WGS data from individuals in four samples (representing European, Ethiopian, Zulu and Baganda populations) we describe patterns of variation in EBTs for autosomal chromosomes and relationship with gene density.

## **Methods**

### ***Processing of samples with European origin***

SNP genotypes were obtained from WGS data from the Scripps Welllderly Genome Resource. Samples comprise 454 unrelated individuals with self-declared ethnically European origin (Erikson *et al.*, 2016). Following Pengelly et al (2015) SNPs with > 5% missing genotypes and a significant Hardy-Weinberg deviation p-value (<0.001) were excluded (Wigginton *et al.*, 2005). Rare SNPs are uninformative for LD and all SNPs with alternative minor allele frequencies (MAF) of <0.01 were excluded from the maps. Completed LDU maps of chromosomes 1-22 contain 7,162,973 SNPs (Supplementary Table 1).

### ***Processing of samples with African origin***

WGS data from 320 healthy individuals from seven Sub-Saharan African (SSA) ethnolinguistic groups recruited to the African Genome Variation Project (AGVP) were considered (Gurdasani *et al.*, 2015). To develop LD maps representing major ethnic populations, principal component analysis (PCA) was carried out using SNPRelate (Zheng *et al.*, 2012). Individuals with >10% missing genotypes, SNPs with >10% missing genotypes and SNPs with significant Hardy-Weinberg deviation were removed. Only SNPs with MAF >5% were retained for PCA analysis. Identity by descent (IBD) was measured within each population and related individuals (IBD > 0.05) were removed. Squared

correlation coefficients between allele counts were used to prune SNPs in strong LD ( $r^2 > 0.5$ ) with other SNPs. PCA analysis confirmed three major clusters mirroring distinct linguistic groupings across these populations (Supplementary Figure 1). One Somali individual was excluded from further consideration due to apparent ethnic divergence from all the major groups. The Gumuz population formed a discrete cluster but comprises only 24 individuals so was not considered further. Samples from Amhara (n=24), Oromo (n= 24), Somali (n= 23) and Wolayta (n=24) sub-populations were considered together as a broader Ethiopian sub-population (n= 95). The two other populations with large sample size (n=100 each), namely Zulu and Baganda, were considered independently.

Population specific multivariant VCF files were generated for 295 individuals according to the PCA clustering pattern as described. Data from each population were pre-processed to exclude SNPs with >5% missing genotypes, SNPs with MAF <0.01 in the respective sub-population and SNPs deviating from Hardy-Weinberg equilibrium at a significance threshold of  $p < 0.001$ .

#### ***Linkage disequilibrium map construction***

For each of the four populations we constructed LD maps in linkage disequilibrium units for chromosomes 1-22 using the LDMAP program (Lau *et al.*, 2007; Pengelly *et al.*, 2015). The program constructs LD maps according to the Malécot-Morton model:

$$\hat{\rho} = (1 - L)Me^{-\epsilon d} + L$$

where  $\hat{\rho}$  is the association between a pair of SNPs, the asymptote L is ‘background’ association which is not due to linkage, M reflects association at zero distance with values  $\sim 1$  consistent with monophyletic haplotypes and  $< 1$  with polyphyletic inheritance,  $\epsilon$  is the rate of LD decline, and d is the physical distance in kilobases between SNPs. Parameters  $\epsilon$ , L and M are estimated iteratively for each adjacent SNP-SNP interval. Map distances in LDUs describe the rate of decline of LD across each SNP interval computed as the product  $\epsilon d$  with cumulative map distances similar to the centimorgan scale. One LDU corresponds to the (highly variable) physical distance along the DNA sequence over which LD declines to background levels. LDUs plotted against chromosome location reveal ‘steps’ which reflect recombination hotspots and plateaus aligning with blocks of low haplotype diversity.

The LDU map is constructed iteratively with the fit of the pairwise SNP association data to the kilobase map established through composite likelihood. Previous studies have shown LDU maps are largely insensitive to variation in the SNP marker density used for their construction, even when using much lower density SNP array data (Ke *et al.*, 2004). Similarly, LDU map lengths have been shown to be largely stable and representative of population LD structure even for small population samples (Pengelly *et al.*, 2015).

#### ***Computation of effective bottleneck time***

The effective bottleneck time was defined for LDU maps by Zhang et al (2004). Considering the  $i^{\text{th}}$  interval between SNPs map distances in the corresponding linkage map can be expressed as  $w_i$  Morgans (typically given as  $100w_i$  cM). Map distances along an LDU map correspond to  $w_i t$ , where  $t$  is the number of generations since LD began to decline from an effective hypothetical bottleneck, reflecting the cumulative impact of multiple bottlenecks. For populations these multiple bottlenecks drive reduced effective population size through mortality, migration, selective sweeps or other factors. Differences between chromosomes suggest variation in intensity of purifying selection impacting the strength of LD after correcting for variable recombination rates. The LDU/ $w$  ratio describes the effective bottleneck time in generations: evidence that some chromosomes within a population show relatively reduced EBTs suggests they have been subject to elevated selective pressure.

#### ***Alternative genetic recombination maps***

The Kong et al (2002) linkage map was constructed using 5,136 microsatellite markers genotyped in 146 Icelandic families (869 individuals) and includes 1,257 meiotic events spanning an autosomal chromosome length of 3,436 cM (Supplementary Tables 2 and 3). Bh  rer et al (2017) used recombination data from six different sources and four data sets representing European populations to construct a ‘refined’ European linkage map spanning 3,351 cM and representing 97,723 meioses. The majority of the samples used were from Icelandic families (Kong *et al.*, 2014). Hinch et al (2011) constructed a recombination map from African Americans (individuals with ~80% West African and 20% European ancestry) from 29,589 apparently unrelated African Americans genotyped on SNP

arrays for genome-wide association studies. However, the map is constructed from population data and not families and has therefore been normalised by total map length. However, any systematic differences in relative chromosome lengths between a (largely) African population compared to European map should be evident.

### ***Gene density***

We considered gene densities for human chromosomes presented by Mayer et al (2005) in their figure 1. We also computed gene densities using the total gene counts in Spataro et al (2017) and the chromosome physical lengths in megabases from the human genome hg19 assembly (Supplementary Tables 2 and 4). The Spataro (2017) study provides gene groups categorised according to essentiality and disease relationships. We computed gene densities within these gene groups (Supplementary Table 4).

### ***Chromosome regional analysis***

To establish relationships between EBTs and gene densities in chromosome sub-regions the genome was firstly segmented into contiguous 10 Mb regions. Regions containing < 30,000 SNP markers in the LD maps were excluded from further analysis, as these represent small regions truncated by the end of chromosomes or are regions containing extensive centromeric/heterochromatic sequences or other unsequenceable regions. Regional EBTs were calculated from the Bhérer (2017) linkage map and the Zulu LDU map (Gurdasani *et al.*, 2015). The regional densities of genes were computed counting a particular gene as within a region if the region contained at least half of the gene span.

### ***Chromosome simulations***

To independently validate real-data relationships between EBTs and gene densities we simulated chromosome populations with known parameters using the SLiM V3.3 software (Haller and Messer, 2019). SLiM V3.3 is a highly flexible simulation software for modelling the chromosomes of individuals in a population. The length of the chromosomes simulated was 10 Mb. The simulations began from a population of 10,000 and ran for 5,000 generations. The overall mutation rate was given



as  $10^{-8}$  per site per generation and the recombination rate was set at  $10^{-8}$  per adjacent bases per generation. The model simulated two types of mutation: neutral mutations which did not affect fitness, and deleterious mutations, with a dominance coefficient of 0.5 and a fixed fitness effect of -0.03. Chromosomes contained genic and non-genic regions. Genes were fixed at 10,000 bp in length, and arranged evenly across the region, with rates of between five and 20 gene per Mb. Mutations arising in the non-genic regions were always neutral. Mutations arising within genes were designated neutral or deleterious at a ratio of 2:8. The final output in VCF format was sub-sampled with 100 chromosomes taken randomly from each simulation. VCF files were converted into .tped format using PLINK v1.90 (Purcell *et al.*, 2007). SNPs with a minor allele frequency of  $< 0.01$  and markers failing the Hardy-Weinberg Equilibrium at  $P < 0.001$  significance were excluded. The simulated datasets were then passed to the LDMAP program. Given the fixed recombination rate of  $10^{-8}$  per site (i.e. 0.01 recombination events per MB per generation) and 10 Mb chromosomes, a linkage length on 0.1w was used for EBT calculation.

## Results

### *LD maps*

LD maps were constructed from the Wellderly sample (Erikson *et al.*, 2016) and samples from three SSA populations: Ethiopia, Zulu and Baganda (Gurdasani *et al.*, 2015). The latter were combined from smaller population groups defined by principal component analysis. Supplementary Figure 1 shows the PCA plot made using single nucleotide polymorphism (SNP) data from seven SSA sub-populations. Although PCA decomposition of genotypes demonstrates very close alignment between Zulu and Baganda populations we constructed independent LD maps for both populations. LD maps in linkage disequilibrium units constructed for the 22 autosomal chromosomes include ~7.2 million SNPs from the Wellderly sample and ~13.6-14.4 million SNPs for each of the SSA populations (Supplementary Table 1). Figure 1 shows the LDU contour for a representative chromosome (chromosome 22) for Wellderly, Ethiopia and Baganda populations demonstrating close alignment across populations between regions of LD breakdown (“steps”) and flatter plateau regions of lower

haplotype diversity (strong LD). The overall increased LDU map lengths for SSA populations compared to Wellderly, reflect, as expected for these populations, increased times since an effective population bottleneck. The Ethiopian map is intermediate consistent with increased admixture with populations of European and Middle Eastern origin perhaps through multiple events within the last 3,000 years (Hodgson *et al.*, 2014; Busby *et al.*, 2016).

Map lengths for the autosomes total ~63K LDUs for Wellderly and 107-130K LDUs for the three SSA populations (Supplementary Table 2). The Wellderly map length exceeds a previous estimate by Tapper *et al* (2005), using HapMap array genotype data from the CEU population, of 57,819 LDUs. However, although the Tapper *et al* (2005) maps also included the X chromosome which was excluded here and SNPs were selected using a minor allele frequency (MAF) cut-off of >0.05, unlike the MAF of >0.01 used here. Inclusion of a higher SNP density in the current analysis increases LDU map length through greater resolution of chromosome regions poorly represented in lower density maps (Pengelly *et al.*, 2015).

### ***Effective bottleneck times***

Figure 2 and Supplementary Table 3 give effective bottleneck times by population and chromosome. EBTs were computed using three alternative linkage maps: the Kong *et al* (2002) map, made from Icelandic pedigrees (“Kong map”); the Bhérer *et al* (2017) map which includes a larger number of meioses and is derived from a combined European data set which includes many Icelandic families (“Bhérer map”), and the Hinch *et al* (2011) map made using data from African American samples (“Hinch map”). The Hinch map is slightly longer than the Bhérer and Kong maps (at ~3,523 cM, Supplementary Table 2) and therefore contributes relatively reduced EBTs for all four populations although the impact is across all autosomes and not biased towards specific chromosomes. There is relatively limited difference in the pattern of EBTs defined using alternate linkage maps within a population (Figure 2). Using the Bhérer map to define EBTs indicates 1,876, 3,073, 3,801 and 3,793 generations to an effective bottleneck for Wellderly, Ethiopia, Zulu and Baganda populations respectively. The observed pattern is consistent with the expectation from numerous other studies of

reduced effective population sizes experienced by European populations. Figure 2 demonstrates the much increased variance in EBTs for different chromosomes for SSA populations relative to the Welldeby group. The Ethiopian population is intermediate with relatively reduced variance in EBTs compared to the more southerly located SSA populations but with increased variance compared to the Welldeby population. Populations in the Horn of Africa have experienced substantial gene-flow from Middle Eastern and European populations during the past 3,000 years (Hodgson *et al.*, 2014) and their genomes show evidence consistent with back-to-Africa migratory events (Pickrell *et al.*, 2014; Busby *et al.*, 2016). As a result, Eurasian admixture over a variety of timescales has generated relatively increased LD demonstrated by reduced LDU map lengths compared to the other SSA populations considered here.

### ***Effective bottleneck time and gene density***

The calculation of gene density for each chromosome depends in part on how the coding regions are defined. We considered alternative gene densities given by Mayer *et al* (2005) and Spataro *et al* (2017) (Supplementary Table 4). The latter provides gene counts in different gene groups: non disease non-essential (NDNE) which are genes not known to be involved in disease phenotypes or classed as ‘essential’; complex non-Mendelian (CNM) genes known to contain common variation involved in complex but not currently implicated in Mendelian phenotypes; complex-Mendelian (CM) genes containing variants known to be involved in both complex and Mendelian phenotypes; Mendelian non-complex (MNC) genes known to contain Mendelian, but not complex trait variation and essential non-disease (END) essential genes not known to contain disease variation. END genes were defined as having a mouse ortholog showing pre-natal, peri-natal or post-natal lethality in mouse knockouts. Gene densities from both Mayer *et al* and Spataro *et al* indicate that chromosomes 13, 4 and 18 have the lowest density of genes (Mayer *et al.*, 2005; Spataro *et al.*, 2017). Mayer *et al* (2005) has chromosomes 19, 17 and 22 as the most gene dense and Spataro *et al* (2017) has 19, 17 and 11 as most gene dense. Table 1, supplementary Table 5 and Figure 3 show relationships between gene densities and effective bottleneck times. Negative correlations between gene densities and EBTs are found for all four populations although only as a trend for the Welldeby population (Table 1).

Correlations are stronger for the Mayer et al gene densities than the Spataro et al ‘all’ gene densities although the density of END genes shows the strongest correlation with EBTs. From the linear fit the reduction in EBT from chromosomes with the lowest density of END genes to the highest density for the SSA populations (Figure 3) is ~17-21% corresponding to ~514-748 generations. The density of genes associated with only complex phenotypes (CNM, Table 1) shows non-significant correlations.

Figure 4 shows the relationship between EBTs in sub-chromosomal regions across the genome (computed using the Bh  rer linkage map and the Zulu LDU map, Supplementary Table 2) and the density of END genes in each region. The pattern is consistent with the relationship observed for whole chromosomes of reduced EBTs associated with higher gene densities (for the 10 Mb regions, correlation  $r = -0.328$ ,  $P < 0.00001$ ).

The data are supportive of the hypothesis that chromosomes and their sub-regions which have a high density of essential genes are subject to stronger purifying selection reflected in consistently reduced EBTs. The trend is weak for the Wellderly sample and strongest for both the Zulu and Baganda populations (Figure 3). The impact of defining EBT using the alternative Kong and Hinch linkage maps has a variable impact on the significance of correlations but trends remain consistent (Supplementary Table 5).

Chromosome simulations using known input parameters were performed to investigate the relationship between gene density and EBT. Figure 5 shows LDU maps from simulated chromosomes where only gene densities are varied. All other parameters were fixed between populations, including recombination and mutation rates, however mutations falling in genic regions were assumed to have a higher chance of being deleterious with non-genic regions subject only to neutral mutation. The results demonstrate a relationship between increasing gene density and decreasing LDU map length. Given the fixed recombination rate per chromosome the LDU map lengths are readily converted into EBTs (Figure 6) demonstrating a strong negative correlation with gene density ( $p < 0.001$ ). The simulations provide support to the relationship shown here in real data between increased gene density and reduced EBTs reflecting purifying selection in gene-dense chromosomes.

## 296 Discussion

297 Selection which acts to increase the frequency of beneficial mutations and eliminate detrimental  
298 mutations is more efficient in large populations (Cutter and Payseur, 2013). Therefore the intensity of  
299 hitch-hiking and the impact of background selection is increased within populations with greater  
300 effective population size ( $N_e$ ) (Charlesworth, 2009). Given increased  $N_e$  in SSA populations,  
301 compared to populations with European ancestry, this expectation is consistent with differences in  
302 EBTs which indicate chromosome-specific patterns of selection. It is known that intense population  
303 bottlenecks, such as those experienced by European populations on leaving Africa, may erase  
304 historically-present signatures of selection at linked sites (Cutter and Payseur, 2013). This might  
305 account for the much reduced evidence (trend only) for a relationship between chromosome EBTs and  
306 gene density in the Welllderly sample which has European ancestry.

307 Accurate computation of EBTs depends on availability of population-specific genetic linkage maps,  
308 unless the pattern of recombination is broadly conserved across populations. The Kong linkage map  
309 from Icelandic families and the Bh  rer map, from a combined data set representing European  
310 populations, yield consistently similar EBTs. The lack of a recombination map from a pure SSA  
311 population presents difficulties. The Hinch map represents an African-American population, with a  
312 high proportion of SSA ancestry, although it is constructed from population rather than family data  
313 and has therefore been normalised for total map length. There is little indication that recombination  
314 rates differ by chromosome sufficiently to impact the pattern of chromosome-specific EBTs and there  
315 is good evidence that broad patterns of recombination are conserved across human populations  
316 (Jorgenson *et al.*, 2005; Serre *et al.*, 2005). The ratio of recombination rates for two populations has  
317 been shown to be constant along the chromosomes such that genetic linkage maps made from  
318 European populations are considered valid to make inferences in other populations (Serre *et al.*, 2005).  
319 However, on finer scales there is evidence that more of the genome is recombinationally active in  
320 West Africans (Hinch *et al.*, 2011) which have a larger number of recombination hotspots making  
321 crossovers more evenly distributed across chromosomes compared to Europeans. Only limited  
322 variation in EBTs is evident when defined using the three alternative genetic linkage maps, consistent

with the suggestion that European-based linkage maps are adequate for making wider inferences in other populations.

Increased haplotype diversity, and therefore overall reduced LD, amongst SSA populations reflects their extended population history. Detection of selection events is enhanced by both the extended history of SSA populations, providing sufficient time for selective forces to create a detectable signal, and the reduced impact of population bottlenecks which may erase signatures of selection. Although some signatures of recent adaptation since the divergence of African and non-African populations have been identified (Hamblin and Di Rienzo, 2000; Sabeti *et al.*, 2002; Bersaglieri *et al.*, 2004; Lamason *et al.*, 2005; Pickrell *et al.*, 2009) there is limited evidence for frequent hard selective sweeps since migration out of Africa (Teshima *et al.*, 2006) and identification of such events is confounded by more prominent evolutionary forces such as population bottlenecks, drift and stratification (Coop *et al.*, 2009; Berg *et al.*, 2018).

Evidence established over a number of years demonstrates that selection has impacted whole chromosomes to different degrees. The sex chromosomes show strong signatures of increased selection relative to the autosomes as a whole. X chromosome hemizyosity in males results in accelerated adaptive evolution across X-linked genes (Vicoso and Charlesworth, 2006). Recessive X-linked mutations in males are subject to selection and become readily fixed in the population even when their deleterious burden in females surmounts their advantageous utility in males (Nam *et al.*, 2015). As a result, the X-chromosome is enriched for sexually antagonistic alleles that are pertinent to male and female reproductive fitness (Sangrithi and Turner, 2018). Although gene density on the X is generally low compared to autosomes there is evidence from diversity that the targets of selection include its coding regions. Hammer *et al.* (2010) considered the relationship between normalised nucleotide diversity and genetic distance from genes on both the X chromosome and autosomes. Considering nongenic regions in 0.1 cM bins, diversity was substantially lower within bins located close to genes compared to bins further away. Although the magnitude of the effect was more marked

on the X chromosome compared to the autosomes a broadly similar relationship was observed for both.

The evidence here indicates a strong relationship between chromosome-specific EBTs and the density of essential genes. Classical examples of positive selection in the human genome have been demonstrated for individual genes and include non-essential genes such as *SLC24A5* on chromosome 15 (Sabeti *et al.*, 2002) and *LCT* (Bersaglieri *et al.*, 2004) and *EDAR* (Sabeti *et al.*, 2007) which are both on chromosome 2. However, evidence for a relationship between variation in chromosome-specific EBTs and the density of all genes (not just essential genes) is consistent with purifying selection impacting numerous targets across the genome. It is known that chromosomes and regions with high gene densities are subject to weaker positive selection due to Hill-Robertson interference (Hill and Robertson, 1966) but are under increased purifying selection acting on functionally deleterious mutations (Castellano *et al.*, 2016). Our findings show that the impact of these processes can be observed at the level of whole chromosomes, at least in SSA populations.

Only small differences in EBTs are evident between different autosomes in the European population studied here. Historical emphasis on studies in European populations may have made chromosome-specific signatures of selection difficult to detect. The greatly increased availability of samples derived from European compared to SSA populations has reduced the possibility of detecting chromosome-specific differences in selection intensity for different autosomes. The evidence presented here supports deeper analysis of signatures of selection within SSA populations and the expectation that studies in these populations offer greater power to detect chromosome (and gene) specific signatures of selection.

## Conclusions

The increased lengths of LDU maps constructed from SSA human populations, compared to a map from a population of European ancestry are consistent with the expectation of increased EBTs in these populations. Previous focus on European populations has suggested that EBTs are relatively constant for all autosomal chromosomes, although reduced for the X chromosome in line with evidence that it

376 has undergone unusually intense selective sweeps. However, analysis of SSA populations shows  
377 variability in EBTs across the autosomes with reduced EBTs for some chromosomes, such as  
378 chromosomes 17 and 19, which have high a density of essential genes. This is consistent with a  
379 pattern of strong purifying selection (and weaker positive selection due to Hill-Robertson interference)  
380 reducing the diversity of gene-dense chromosomes. The findings strongly support efforts to analyse  
381 patterns of selection in SSA populations since the power to detect signals is enhanced through their  
382 extended population history and reduced impact of intense population bottlenecks.

383

384



## References

- Bachtrog D (2013). Y-chromosome evolution: emerging insights into processes of Y-chromosome degeneration. *Nat Rev Genet* **14**: 113–124.
- Berg JJ, Harpak A, Sinnott-Armstrong N, Jørgensen AM, Mostafavi H, Field Y, *et al.* (2018). Reduced signal for polygenic adaptation of height in UK Biobank. *bioRxiv*: 354951.
- Bersaglieri T, Sabeti PC, Patterson N, Vanderploeg T, Schaffner SF, Drake JA, *et al.* (2004). Genetic Signatures of Strong Recent Positive Selection at the Lactase Gene. *Am J Hum Genet* **74**: 1111–1120.
- Bhérier C, Campbell CL, Auton A (2017). Refined genetic maps reveal sexual dimorphism in human meiotic recombination at multiple scales. *Nat Commun* **8**: 14994.
- Busby GB, Band G, Si Le Q, Jallow M, Bougama E, Mangano VD, *et al.* (2016). Admixture into and within sub-Saharan Africa. *Elife* **5**.
- Castellano D, Coronado-Zamora M, Campos JL, Barbadilla A, Eyre-Walker A (2016). Adaptive Evolution Is Substantially Impeded by Hill-Robertson Interference in *Drosophila*. *Mol Biol Evol* **33**: 442–55.
- Charlesworth B (2009). Effective population size and patterns of molecular evolution and variation. *Nat Rev Genet* **10**: 195–205.
- Coop G, Pickrell JK, Novembre J, Kudaravalli S, Li J, Absher D, *et al.* (2009). The Role of Geography in Human Adaptation (MH Schierup, Ed.). *PLoS Genet* **5**: e1000500.
- Cutter AD, Payseur BA (2013). Genomic signatures of selection at linked sites: unifying the disparity among species. *Nat Rev Genet* **14**: 262–274.
- Dutheil JY, Munch K, Nam K, Mailund T, Schierup MH (2015). Strong Selective Sweeps on the X Chromosome in the Human-Chimpanzee Ancestor Explain Its Low Divergence (NH Barton, Ed.). *PLOS Genet* **11**: e1005451.

409 Erikson GA, Bodian DL, Rueda M, Molparia B, Scott ER, Scott-Van Zeeland AA, *et al.* (2016).  
410 Whole-Genome Sequencing of a Healthy Aging Cohort. *Cell* **165**: 1002–1011.

411 Gottipati S, Arbiza L, Siepel A, Clark AG, Keinan A (2011). Analyses of X-linked and autosomal  
412 genetic variation in population-scale whole genome sequencing. *Nat Genet* **43**: 741–743.

413 Gurdasani D, Carstensen T, Tekola-Ayele F, Pagani L, Tachmazidou I, Hatzikotoulas K, *et al.* (2015).  
414 The African Genome Variation Project shapes medical genetics in Africa. *Nature* **517**: 327–332.

415 Haller BC, Messer PW (2019). SLiM 3: Forward Genetic Simulations Beyond the Wright–Fisher  
416 Model (R Hernandez, Ed.). *Mol Biol Evol* **36**: 632–637.

417 Hamblin MT, Di Rienzo A (2000). Detection of the Signature of Natural Selection in Humans:  
418 Evidence from the Duffy Blood Group Locus. *Am J Hum Genet* **66**: 1669–1679.

419 Hammer MF, Woerner AE, Mendez FL, Watkins JC, Cox MP, Wall JD (2010). The ratio of human X  
420 chromosome to autosome diversity is positively correlated with genetic distance from genes. *Nat*  
421 *Genet* **42**: 830–831.

422 Hill WG, Robertson A (1966). The effect of linkage on limits to artificial selection. *Genet Res* **8**: 269–  
423 94.

424 Hinch AG, Tandon A, Patterson N, Song Y, Rohland N, Palmer CD, *et al.* (2011). The landscape of  
425 recombination in African Americans. *Nature* **476**: 170–175.

426 Hodgson JA, Mulligan CJ, Al-Meerri A, Raaum RL (2014). Early Back-to-Africa Migration into the  
427 Horn of Africa (SM Williams, Ed.). *PLoS Genet* **10**: e1004393.

428 Hudson RR, Kaplan NL (1995). Deleterious background selection with recombination. *Genetics* **141**:  
429 1605–17.

430 Huff CD, Harpending HC, Rogers AR (2010). Detecting positive selection from genome scans of  
431 linkage disequilibrium. *BMC Genomics* **11**: 8.

432 Jorgenson E, Tang H, Gadde M, Province M, Leppert M, Kardia S, *et al.* (2005). Ethnicity and  
433 Human Genetic Linkage Maps. *Am J Hum Genet* **76**: 276–290.

434 Ke X, Hunt S, Tapper W, Lawrence R, Stavrides G, Ghorri J, *et al.* (2004). The impact of SNP density  
435 on fine-scale patterns of linkage disequilibrium. *Hum Mol Genet.*

436 Kong A, Gudbjartsson DF, Sainz J, Jonsdottir GM, Gudjonsson SA, Richardsson B, *et al.* (2002). A  
437 high-resolution recombination map of the human genome. *Nat Genet* **31**: 241–247.

438 Kong A, Thorleifsson G, Frigge ML, Masson G, Gudbjartsson DF, Villemoes R, *et al.* (2014).  
439 Common and low-frequency variants associated with genome-wide recombination rate. *Nat*  
440 *Genet* **46**: 11–16.

441 Lamason RL, Mohideen M-APK, Mest JR, Wong AC, Norton HL, Aros MC, *et al.* (2005). SLC24A5,  
442 a Putative Cation Exchanger, Affects Pigmentation in Zebrafish and Humans. *Science (80- )* **310**:  
443 1782–1786.

444 Lau W, Kuo T-YTY, Tapper W, Cox S, Collins A (2007). Exploiting large scale computing to  
445 construct high resolution linkage disequilibrium maps of the human genome. *Bioinformatics* **23**:  
446 517–519.

447 Mayer R, Brero A, von Hase J, Schroeder T, Cremer T, Dietzel S (2005). Common themes and cell  
448 type specific variations of higher order chromatin arrangements in the mouse. *BMC Cell Biol* **6**:  
449 44.

450 Nam K, Munch K, Hobolth A, Dutheil JY, Veeramah KR, Woerner AE, *et al.* (2015). Extreme  
451 selective sweeps independently targeted the X chromosomes of the great apes. *Proc Natl Acad*  
452 *Sci U S A* **112**: 6413–8.

453 Payseur BA, Nachman MW (2002). Gene Density and Human Nucleotide Polymorphism. *Mol Biol*  
454 *Evol* **19**: 336–340.

455 Pengelly RJ, Tapper W, Gibson J, Knut M, Tearle R, Collins A, *et al.* (2015). Whole genome

456 sequences are required to fully resolve the linkage disequilibrium structure of human  
 457 populations. *BMC Genomics* **16**: 666.

458 Pengelly RJ, Vergara-Lope A, Alyousfi D, Jabalameli MR, Collins A (2019). Understanding the  
 459 disease genome: gene essentiality and the interplay of selection, recombination and mutation.  
 460 *Brief Bioinform* **20**: 267–273.

461 Pickrell JK, Coop G, Novembre J, Kudaravalli S, Li JZ, Absher D, *et al.* (2009). Signals of recent  
 462 positive selection in a worldwide sample of human populations. *Genome Res* **19**: 826–837.

463 Pickrell JK, Patterson N, Loh P-R, Lipson M, Berger B, Stoneking M, *et al.* (2014). Ancient west  
 464 Eurasian ancestry in southern and eastern Africa. *Proc Natl Acad Sci* **111**: 2632–2637.

465 Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MARR, Bender D, *et al.* (2007). PLINK: a  
 466 tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet*  
 467 **81**: 559–575.

468 Sabeti PC, Reich DE, Higgins JM, Levine HZP, Richter DJ, Schaffner SF, *et al.* (2002). Detecting  
 469 recent positive selection in the human genome from haplotype structure. *Nature* **419**: 832–837.

470 Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, Cotsapas C, *et al.* (2007). Genome-wide  
 471 detection and characterization of positive selection in human populations. *Nature* **449**: 913–918.

472 Sangrithi MN, Turner JMA (2018). Mammalian X Chromosome Dosage Compensation: Perspectives  
 473 From the Germ Line. *BioEssays* **40**: 1800024.

474 Serre D, Nadon R, Hudson TJ (2005). Large-scale recombination rate patterns are conserved among  
 475 human populations. *Genome Res* **15**: 1547–52.

476 Smith JM, Haigh J (1974). The hitch-hiking effect of a favourable gene. *Genet Res* **23**: 23–35.

477 Spataro N, Rodríguez JA, Navarro A, Bosch E (2017). Properties of human disease genes and the role  
 478 of genes linked to Mendelian disorders in complex disease aetiology. *Hum Mol Genet* **26**: 489–  
 479 500.

480 Tapper W, Collins A, Gibson J, Maniatis N, Ennis S, Morton NE (2005). A map of the human  
481 genome in linkage disequilibrium units. *Proc Natl Acad Sci U S A* **102**: 11835–11839.

482 Teshima KM, Coop G, Przeworski M (2006). How reliable are empirical genomic scans for selective  
483 sweeps? *Genome Res* **16**: 702–712.

484 Veeramah KR, Gutenkunst RN, Woerner AE, Watkins JC, Hammer MF (2014). Evidence for  
485 Increased Levels of Positive and Negative Selection on the X Chromosome versus Autosomes in  
486 Humans. *Mol Biol Evol* **31**: 2267–2282.

487 Vicoso B, Charlesworth B (2006). Evolution on the X chromosome: unusual patterns and processes.  
488 *Nat Rev Genet* **7**: 645–653.

489 Voight BF, Kudaravalli S, Wen X, Pritchard JK (2006). A Map of Recent Positive Selection in the  
490 Human Genome (L Hurst, Ed.). *PLoS Biol* **4**: e72.

491 Wigginton JE, Cutler DJ, Abecasis GRGR (2005). A note on exact tests of Hardy-Weinberg  
492 equilibrium. *Am J Hum Genet* **76**: 887–893.

493 Wilson Sayres MA, Lohmueller KE, Nielsen R (2014). Natural Selection Reduced Diversity on  
494 Human Y Chromosomes (BA Payseur, Ed.). *PLoS Genet* **10**: e1004064.

495 Zhang W, Collins A, Gibson J, Tapper WJ, Hunt S, Deloukas P, *et al.* (2004). Impact of population  
496 structure, effective bottleneck time, and allele frequency on linkage disequilibrium maps. *Proc*  
497 *Natl Acad Sci* **101**: 18075–18080.

498 Zheng X, Levine D, Shen J, Gogarten SM, Laurie C, Weir BS (2012). A high-performance computing  
499 toolset for relatedness and principal component analysis of SNP data. *Bioinformatics* **28**: 3326–  
500 3328.

501

502

## Figure legends

**Figure 1. LDU maps of chromosome 22 for three populations.** The cumulative LDU scale is plotted against the physical location in kilobases for three populations. Map contours show close alignment across populations between steeper areas of rapid LD breakdown and flatter regions of strong LD. Map lengths reflect differences in effective bottleneck times between populations with a longer time to an effective bottleneck for the African Baganda and Ethiopian populations compared to the Wellderly sample of European origin.

**Figure 2. Box plot of effective bottleneck times of autosomes for four populations and alternative linkage maps.** Effective bottleneck times computed using the Kong, Bhérer and Hinch linkage maps show small differences with relatively reduced effective bottleneck times when using the African-American Hinch map compared to the Bhérer linkage map from a European sample. The magnitude of variation in chromosome-specific effective bottleneck times follows the trend European < Ethiopian < Zulu/Baganda and is consistent whichever linkage map is used.

**Figure 3. Effective bottleneck times (using Bhérer map) against the density of Essential Non-Disease genes (END).** Chromosome numbers corresponding to the points plotted are given in circles above the graph. A strong trend following reduced effective bottleneck times for chromosomes with increasing density of essential genes is shown for the SSA populations (Ethiopia, Zulu, Baganda) with a weaker and non-significant trend for the European (Wellderly) population. Linear trends indicate ~17, 20 and 21% reduction in effective bottleneck times for Ethiopia, Baganda and Zulu populations respectively for chromosomes with lowest to highest density of essential genes. The difference is only ~5% for the European (Wellderly) population.

**Figure 4. Effective bottleneck times in 10 Mb regions computed from the Bhérer map and the Zulu LDU map versus the regional density of Essential Non-Disease genes (END).** The trend towards reduced EBTs with an increasing regional density of essential genes is evident. The linear trend is consistent with a ~57% reduction in EBTs for regions with the lowest to highest density of essential genes.

**Figure 5. LDU map lengths of gene-dense and gene-sparse chromosomes simulated using SLiM software.** LDU maps made from a sample of 100 individuals from a simulated population using SLiM V3.3. The same parameters were applied to all populations, including recombination rates, with the only variation being in gene density, ranging between five and twenty genes per Mb. Reduced LDU lengths (and correspondingly reduced EBTs in the gene-dense chromosome populations) is consistent with real-data findings in SSA populations.

**Figure 6. Relationship between EBT and gene density for chromosomes simulated using SLiM software.** LDU maps made from a sample of 100 individuals from a simulated population using SLiM V3.3. Increased gene density is negatively correlated with EBT as predicted and demonstrated in the SSA data. This relationship is statistically significant, calculated using the default F-test in R.

544

545 **Table 1 Pearson correlations (P values) between chromosome-specific effective bottleneck times (Bh  rer map) and gene density.**

Population	Mayer et al (2005) genes/Mb	Spataro et al (2017) all genes/Mb	NDNE genes/Mb	CNM genes/Mb	CM genes/Mb	MNC genes/Mb	END genes/Mb
Wellderly	-0.251 (0.2604)	-0.217 (0.3311)	-0.201 (0.3696)	-0.207 (0.3549)	<b>-0.481 (0.0234*)</b>	-0.116 (0.6069)	-0.298 (0.1779)
Ethiopia	<b>-0.533 (0.0107 *)</b>	<b>-0.471 (0.0271 *)</b>	<b>-0.448 (0.0367 *)</b>	-0.314 (0.1547)	<b>-0.444 (0.0384*)</b>	<b>-0.507 (0.0161*)</b>	<b>-0.570 (0.0057 **)</b>
Zulu	<b>-0.502 (0.0173 *)</b>	<b>-0.444 (0.0383 *)</b>	-0.422 (0.0507)	-0.283 (0.2024)	<b>-0.474 (0.0258*)</b>	<b>-0.453 (0.0344*)</b>	<b>-0.540 (0.0095 **)</b>
Baganda	<b>-0.540 (0.0094 **)</b>	<b>-0.480 (0.0237 *)</b>	<b>-0.461 (0.0310 *)</b>	-0.314 (0.1550)	<b>-0.496 (0.0190*)</b>	<b>-0.485 (0.0221*)</b>	<b>-0.566 (0.0061**)</b>

546 **Significance level: P= \* < 0.05, \*\* <0.01.**

547

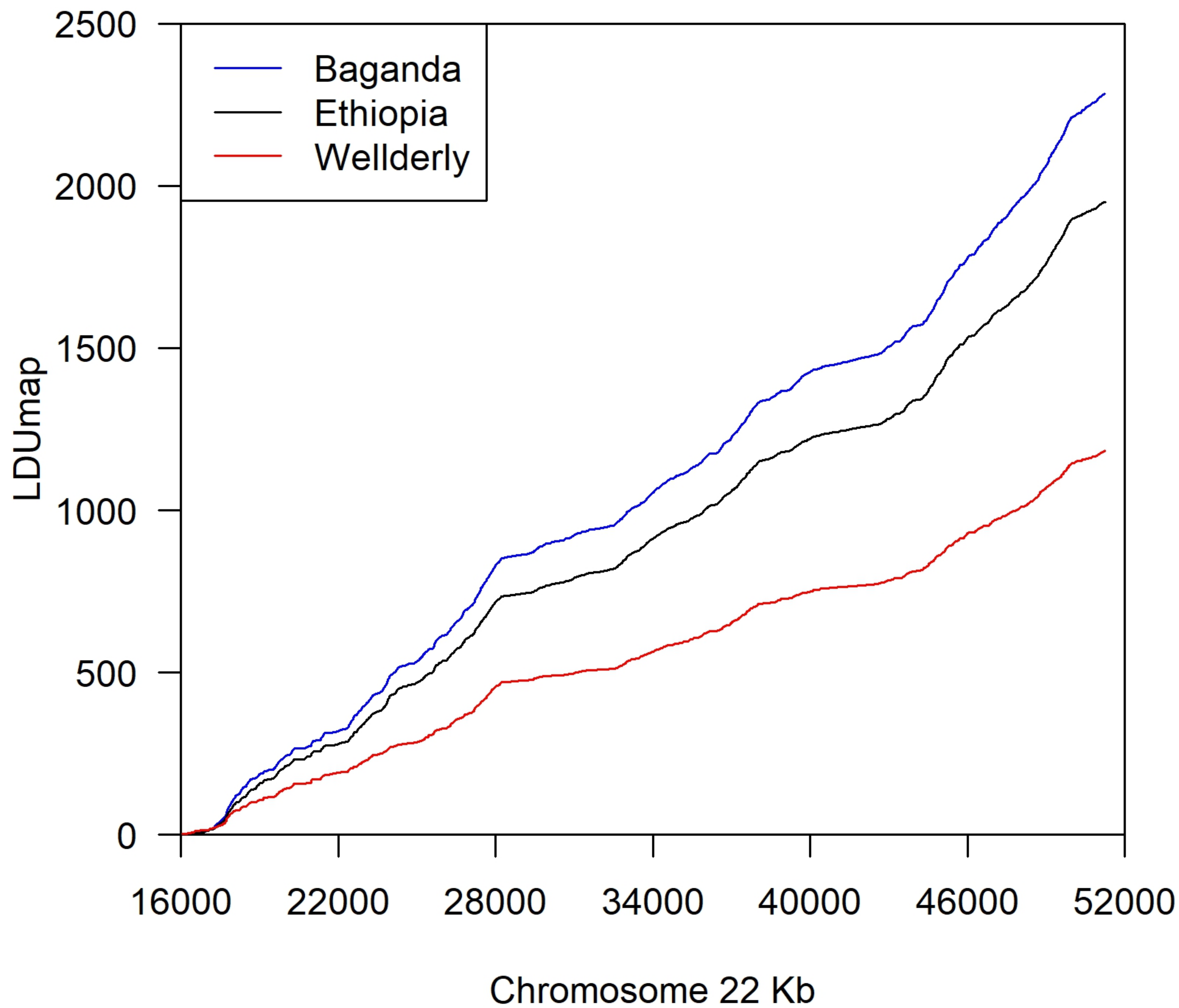
548

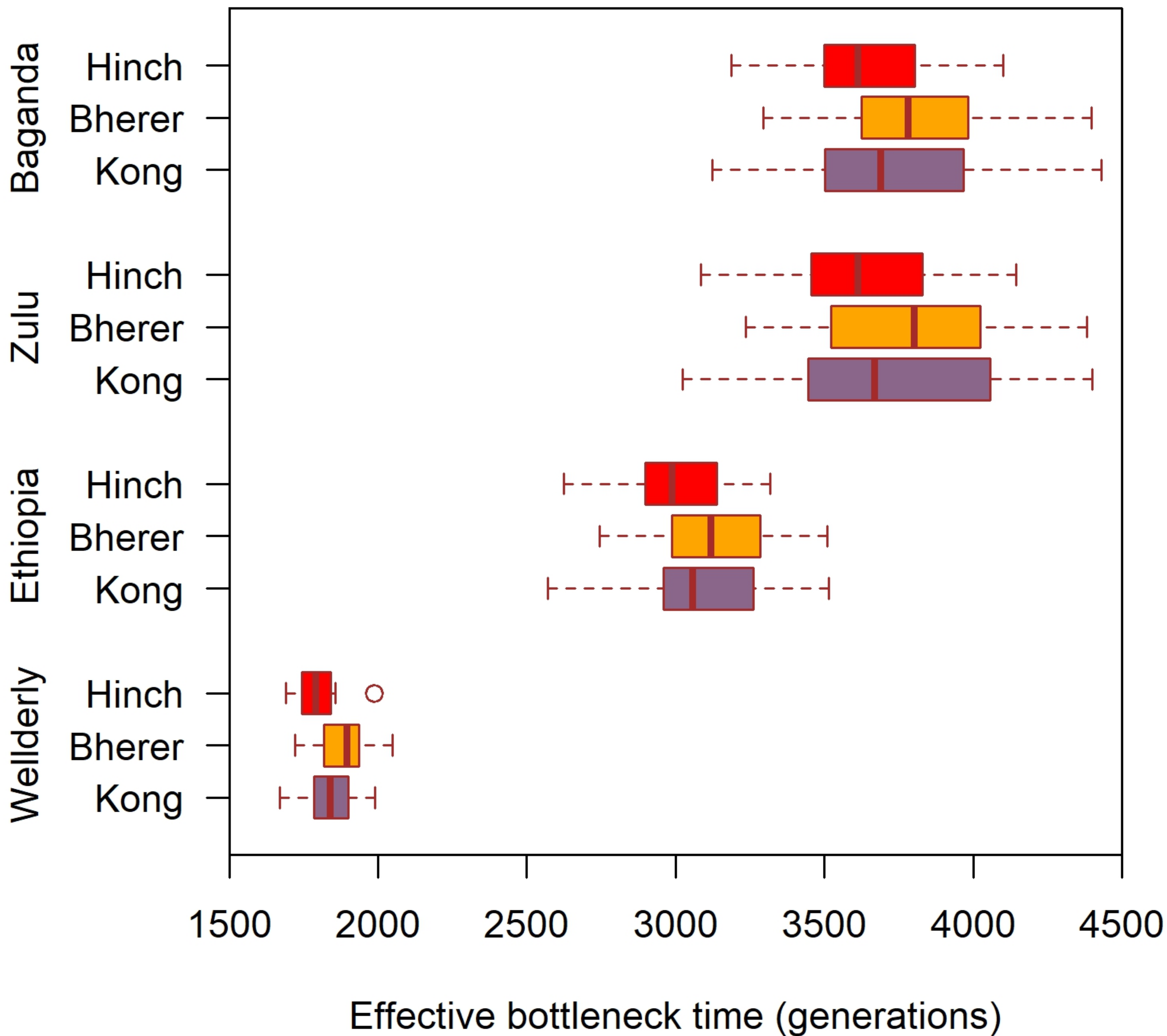
549

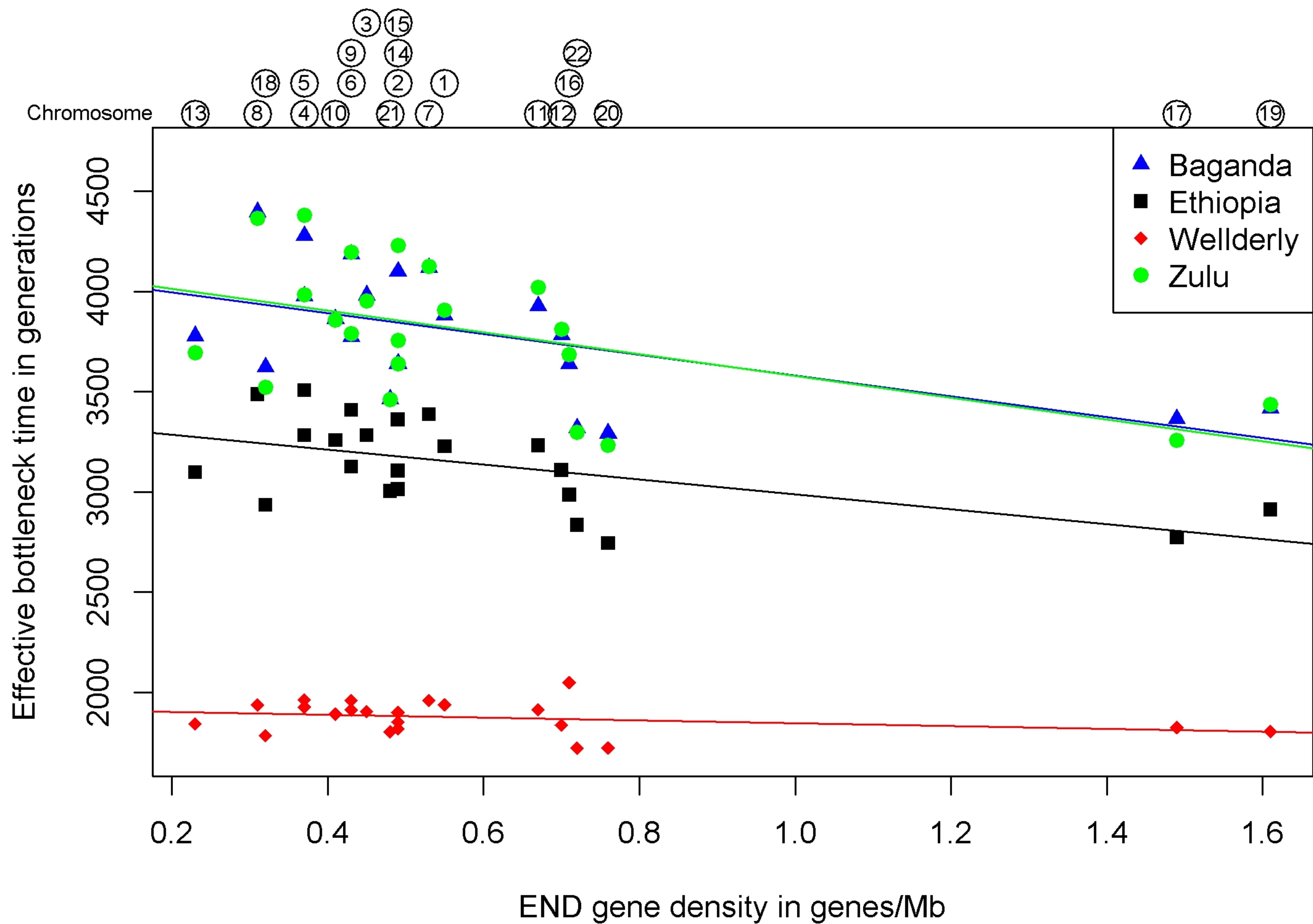
550

551











Effective bottleneck time in generations

