



Genome-wide detection of natural selection in African Americans pre-and post-admixture

Wenfei Jin, Shuhua Xu, Haifeng Wang, et al.

Genome Res. published online November 29, 2011

Access the most recent version at doi:[10.1101/gr.124784.111](https://doi.org/10.1101/gr.124784.111)

P<P Published online November 29, 2011 in advance of the print journal.

Accepted Manuscript Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.

License

Email Alerting Service Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Copyright © 2011, Cold Spring Harbor Laboratory Press

1 **Abstract**

2 It is particularly meaningful to investigate natural selection in African Americans (AfA)
3 due to the high mortality their African ancestry has experienced in history. In this study, we
4 examined 491,526 autosomal SNPs genotyped in 5,210 individuals and conducted a
5 genome-wide search for selection signals in 1,890 AfA. Several genomic regions showing
6 excess of African or European ancestry, which were thought as the footprints of selection
7 since population admixture, were detected based on a commonly used approach. However,
8 we also developed a new strategy to detect natural selection both pre-and post-admixture
9 by reconstructing an ancestral African population (AAF) from inferred African
10 components of ancestry in AfA and comparing it with indigenous African populations
11 (IAF). Interestingly, many selection-candidate genes identified by the new approach were
12 associated with AfA specific high-risk diseases such as prostate cancer and hypertension,
13 suggesting an important role these disease-related genes might have played in adapting to
14 new environment. *CD36* and *HBB*, whose mutations confer a degree of protection against
15 malaria, were also located in the highly differentiated regions between AAF and IAF.
16 Further analysis showed that the frequencies of alleles protecting against malaria in AAF
17 were lower than that in IAF, which consists with the relaxed selection pressure of malaria
18 in the New World. There is no overlap between the top candidate genes detected by the two
19 approaches, indicating the different environmental pressures AfA experienced pre-and
20 post-population-admixture. We suggest that the new approach is reasonably powerful and
21 can also be applied to other admixed populations such as Latinos and Uyghurs.

22 Supplementary material is available online at <http://www.genome.org>

1 **Introduction**

2 Although the vast majority of human genetic variations evolve neutrally, some
3 parts, however, have been shaped by natural selection ([Kimura 2003](#); [Balaesque et al.](#)
4 [2007](#)). Studies on recent natural selection have led to the discovery of genes showing
5 population differences in adapting to pathogens, diet, climate and other environmental
6 challenges. These discoveries have greatly enriched our understanding about the origins
7 and also the evolutionary history of human species, identified many genes with important
8 biological functions, and will lead to further elucidation of the genetic basis of some human
9 diseases ([Balaesque et al. 2007](#); [Nielsen et al. 2007](#); [Sabeti et al. 2007](#); [Hancock et al.](#)
10 [2008](#); [Akey 2009](#)). The recent availability of high-density SNPs has provided essential
11 resources for genome-wide detection of natural selection, especially in ethnically
12 well-defined populations with little admixture ([Akey 2009](#); [Barreiro et al. 2008](#); [Pickrell et](#)
13 [al. 2009](#); [Sabeti et al. 2007](#)). Although there have been several studies on recently admixed
14 populations ([Tang et al. 2007](#); [Basu et al. 2008](#); [Bryc et al. 2010](#)), no study so far has
15 particularly investigated the locus-specific population differentiation between the ancestral
16 components of the admixed population and its ancestral parental population, which might
17 reflected the natural selection since the two split. In this case, we used African Americans,
18 a well-studied admixed population, as the object of our study.

19 African Americans (AfA) are residents of the United States with partial recent
20 Sub-Saharan African ancestry. The majority of them are descendents of the Africans,
21 probably 500,000 to 650,000 in number, who were forcibly brought to North America
22 during the Middle Passage ([Thomas 1999](#); [Zakharia et al. 2009](#)). However, many of those
23 captive Africans died either during the Atlantic shipment to America due to the severe

1 conditions, or soon upon their arrival in the New World as a result of exposure to foreign
2 pathogens and/or the poor living conditions. Although the exact amount of life lost in this
3 process remains a mystery, it may equal or exceed the actual amount enslaved ([Stannard](#)
4 [1993](#)). Therefore, the high mortality of the AfA in the whole slavery era could be attributed
5 to the overwhelming environmental challenges. These persistent selection pressures might
6 make the frequencies of ‘beneficial’ alleles increase continually, which should lead to
7 higher population differentiation between African components of ancestry in AfA and
8 indigenous Africans at these loci, in contrast to those evolved neutrally.

9 Meanwhile, before Africans and Europeans migrated to the New World, their
10 ancestral parental populations had evolved independently in distinct environments for tens
11 of thousands of years ([Basu et al. 2008](#)). The completely new environment in the New
12 World constituted a challenge for both populations and the populations of their admixture.
13 Therefore, it is very likely that some genomic regions in AfA show excess of a particular
14 ancestry as a result of selection pressures after the population admixture ([Tang et al. 2007](#);
15 [Basu et al. 2008](#)). For example, there were various studies in detecting signatures of
16 selection in AfA by examining admixture proportion using a small number of available loci
17 ([WORKMAN et al. 1963](#); [Reed 1969](#); [Blumberg and Hesser 1971](#); [Adams and Ward 1973](#);
18 [Long 1991](#)). Recently, [Bryc et al. \(2010\)](#) identified three autosomal regions showing
19 excessive or reduced African ancestry (by at least three standard deviation of the mean) as
20 natural selection candidates based on approximately 500K SNPs genotyped in 365 AfA.

21 Here, we analyzed the genotype data from approximately 500,000 autosomal SNPs
22 shared by 1,890 AfA and 3,320 non-AfA for detecting the signatures of selection in AfA
23 that were classified into pre-admixture and post-admixture according to the different

1 environments they experienced (**Figure S1**). Firstly, we detected the natural selections that
2 were more likely to occur after admixture by examining the genome-wide distribution of
3 ancestry in AfA, Then we developed a new strategy by reconstructing an ancestral African
4 population (AAF) from inferred African components of ancestry in AfA and comparing it
5 with indigenous African populations (IAF) (**Figure S2**), which reflect natural selection
6 since the African ancestors of AfA left Africa (including both pre-and post-admixture).
7 Many candidate genes identified by our new approach could explain the challenges that
8 AfA and their ancestry had experienced. Thus we suggest that our new approach is
9 reasonably powerful and can also be applied to the studies of other admixed populations.

10 **Results**

11 *African and European ancestries in AfA*

12 The populations from West Africa and Europe, who had undoubtedly contributed
13 to the current gene pool of AfA, were considered as ancestral parental populations of AfA
14 in this study. However, it was practically difficult to select proper populations as genetic
15 donor to the gene pool of AfA since the extant populations in West Africa and Europe
16 might not necessarily represent those 300 years ago, given the possible influence of genetic
17 drift, selection, and demographic history. In this study, we chose the samples of YRI and
18 CEU to respectively represent those who contributed to the formation of AfA, largely due
19 to the availability of phased genetic data for these two populations. Besides, an exploratory
20 analysis including Human Genome Diversity Project (HGDP) populations ([Li et al. 2008](#))
21 also suggested that YRI and CEU are better choice for ancestral parental populations of

1 AfA than the other populations with available genome-wide data (see [Supplemental Text](#)
2 [and Figure S3](#)).

3 FRAPPE ([Tang et al. 2005](#)) was applied to the genome-wide high-density SNPs
4 data (by taking $K = 2$), and the estimated European contribution to AfA was 21.65% at
5 population level ([Figure S4, S5](#)). STRUCTURE ([Pritchard et al. 2000](#); [Falush et al. 2003](#))
6 analyses yielded virtually identical results. Based on the thinned data with 341,672 SNPs,
7 the European contribution to the 1,890 AfA was estimated to be 21.61% ([Figure S3D](#)).
8 These estimations were essentially consistent with those in previous studies ([Bryc et al.](#)
9 [2010](#); [Smith et al. 2004](#); [Xu et al. 2007](#)), although datasets analyzed were different.

10 *Identification of genomic regions with biased ancestry in AfA*

11 To estimate the distribution of genetic contribution of European and African ancestry
12 to AfA across the genome, we used haplotypes of 88 CEU and 88 YRI (the phased parents
13 of trios, presumably unrelated, from HapMap3 data) to represent their ancestral parental
14 populations. The identical monomorphic SNPs in CEU and YRI samples were removed for
15 they could not provide valuable information in the local ancestry inference. HAPMIX
16 ([Price et al. 2009](#)) was employed to estimate the locus-specific genetic contributions of the
17 ancestral parental populations to AfA (see [Methods](#)) by taking 21.65% as the European
18 contribution, as estimated by FRAPPE. Based on the likelihood given by HAPMIX,
19 generation since admixture (λ) was estimated to be $\lambda = 7$ (essentially hybrid isolation
20 model), which was similar to those based on other AfA datasets ([Smith et al. 2004](#); [Price et](#)
21 [al. 2009](#)). Detailed analysis showed that λ values for most individuals ranged from 1 to 12,
22 and could be explained by continuous-gene-flow (CGF) model (see [Supplemental Text and](#)
23 [Figure S6](#)). The locus-specific European ancestry proportion across the genome of AfA

1 was estimated to be $21.68\% \pm 0.75\%$ (mean \pm SD). The standard deviation (SD) of
2 locus-specific ancestral genetic contributions in this study is lower than those in any
3 previous studies ([WORKMAN et al. 1963](#); [Reed 1969](#); [Blumberg and Hesser 1971](#); [Tang](#)
4 [et al. 2007](#); [Basu et al. 2008](#); [Bryc et al. 2010](#)), which was not beyond our expectation since
5 we used a much larger sample size.

6 The genomic regions showing excessive or reduced ancestry in admixed population
7 are likely to be signatures of natural selection ([Tang et al. 2007](#); [Basu et al. 2008](#); [Bryc et](#)
8 [al. 2010](#); [Oleksyk et al. 2010](#)). The loci showing strong deviation of European ancestry (3
9 SDs above or below the genome-wide average) were therefore identified as candidates of
10 natural selection in this study. Four regions (2p22, 3q13, 6q26, 16q21) with excessive
11 European influence and two regions (1p36, 2q37) with excessive African influence were
12 observed in AfA genomes ([Figure 1](#)). Each of the six regions was significantly different
13 from the genome-wide average of ancestral contributions ($P < 2.2 \times 10^{-16}$, t test). The
14 detailed annotations of the six candidate regions are presented in [Table 1](#). Most genomic
15 regions showing ancestry deviation can be replicated by the analysis with LAMP-ANC
16 (see [Supplemental Text and Figure S7](#)), although it has lower accuracy than HAPMIX
17 based on our simulated data.

18 A close examination of SNPs in the six regions showing biased ancestry revealed
19 neither significant deviation from Hardy-Weinberg expectation nor unusual fraction of
20 missing data, suggesting that the genotyping quality is unlikely the cause of the bias. In
21 addition, the six regions did not overlap with previously reported long-ranged LD blocks
22 including inversions which may confound genome-wide scans for selection signals in
23 admixed populations ([Price et al. 2008](#)), although three short inversions have been found in

1 one European individual in 1p36. Although the estimated locus-specific ancestral
2 contributions across the genome are generally consistent with the study by Bryc *et al.*
3 (2010), the three regions showing biased ancestry identified by them were only in
4 moderate excess of African or European ancestry in this study, possibly due to the different
5 samples and/or sample sizes between the two studies. Our estimations were also generally
6 consistent with the genome-wide distribution of ancestry calculated in previous admixture
7 mappings (Reich *et al.* 2005; Kao *et al.* 2008).

8 The region showing the strongest bias is located on 1p36, where African ancestry is
9 over-represented. In this region, *IGSF21* and *AKR7A2* are located next to the SNPs that
10 show strongest signals, respectively. *IGSF21* belongs to immunoglobulin superfamily, while
11 *AKR7A2* is involved in the detoxification of aldehydes and ketones, and is implicated in
12 various cancers such as pancreatic cancer (Prabl *et al.* 2008; Cui *et al.* 2009). Another
13 region with over-represented African ancestry is 2q37, in which *PDCDI* is involved in
14 signaling of immune system and various diseases. The region showing highest excessive
15 European ancestry is located on 3q13 which only harbors *LSAMP*, a candidate of tumor
16 suppressor gene in human osteosarcomas, and is associated with coronary artery disease
17 (Kresse *et al.* 2009; Yen *et al.* 2009). Interestingly, *EIF2AK2*, involved in influenza
18 infection pathways (McAllister *et al.* 2010; Pereira *et al.* 2010), is located ~200 Kb away
19 from the peak (2p22) showing the second highest excessive European ancestry, which
20 suggested the possible difference between African and European in resisting to influenza.
21 *PACRG*, located in 6q26, is associated with Parkinson's disease. However, the region in
22 16q21 showing excessive European ancestry did not contain any genes or known function
23 elements.

1 *African/European components of ancestry in AfA*

2 The segments of African ancestry and those of the European ancestry in AfA,
3 inferred by HAPMIX, were collectively referred to as African components of ancestry and
4 European components of ancestry, respectively, each of which could be considered as an
5 ancestral African population (AAF) or an ancestral European population (AEU) residing in
6 America before admixture. The inference of AAF and AEU are credible given the accuracy
7 of HAPMIX that is over 98% based on simulated data ([see Supplemental Text](#)).

8 Then population differentiation between AEU and each putative European parental
9 population was calculated, with CEU showing the lowest F_{ST} with AEU ($F_{ST [AEU-CEU]} =$
10 0.0005) among all putative parental populations. When 2,648 Caucasian from GWAS data
11 (referred to as CAU-GWAS) were considered, $F_{ST [AEU-CAU-GWAS]}$ was 0.0006, which was
12 the second lowest among all values. Among all indigenous African populations (IAF), YRI
13 showed the lowest F_{ST} with AAF ($F_{ST [AAF-YRI]} = 0.0007$). When the observed F_{ST} 's were
14 compared with those simulated under neutrality ([see Methods](#)), different pattern emerged
15 between European ancestry and African ancestry. For European ancestry, the observed F_{ST}
16 between AEU and CEU (0.0005) was lower than that simulated (simulated $F_{ST [AEU-CEU]} =$
17 0.0006), and the genome-wide distribution of observed locus-specific F_{ST} did not deviate
18 much from those simulated ($P = 0.042$, Kolmogorov–Smirnov test).

19 However, for African ancestry, the observed F_{ST} between AAF and YRI (0.0007)
20 was higher than that simulated (simulated $F_{ST [AAF-YRI]} = 0.0006$), and genome-wide
21 distribution of locus-specific $F_{ST [AAF-YRI]}$ was significantly different from those simulated
22 ($P < 2.2 \times 10^{-16}$, Kolmogorov–Smirnov test). In particular, Q-Q plot between the observed
23 and simulated locus-specific $F_{ST [AAF-YRI]}$ showed an enrichment of SNPs with high F_{ST} at

1 the tail of the observed locus-specific F_{ST} (**Figure S8**), suggesting a possible role of natural
2 selection. The results remain essentially unchanged when bottlenecks for African ancestry
3 was assumed in our simulations.

4 *Identification of regions highly differentiated between AAF and African*

5 Since the African immigrants left for America, they might have experienced a
6 completely different population history compared with indigenous Africans. The high
7 mortality of AfA during the slavery era ([Meltzer 1993](#); [Stannard 1993](#); [Thomas 1999](#))
8 suggested a possible presence of strong selective pressures. According to the theory of
9 neutrality, F_{ST} is largely influenced by demographic history which affect all loci similarly
10 ([Weir and Cockerham 1984](#); [Kimura 2003](#)), if not equally. By contrast, the force of
11 positive selection acts in a locus-specific manner and tends to increase F_{ST} ([Nielsen 2005](#)),
12 which has been widely used to detect the positive selection in various studies ([Akey et al.](#)
13 [2002](#); [Oleksyk et al. 2010](#)). In this study, we calculated the locus-specific F_{ST} between
14 AAF and YRI (**Figure S9**). However, loci with a very low minor allele frequency (MAF)
15 may be subjected to sampling error and statistical error, therefore, the SNPs with $MAF <$
16 0.05 in AAF or YRI were removed from further analyses.

17 The genome-wide distribution of F_{ST} between AAF and YRI for 401,599
18 autosomal SNPs were presented in **Figure 2A**. In spite of the low differentiation between
19 AAF and YRI, a substantial proportion of SNPs were located at the right tail of the
20 distribution. The functional annotations of the most differentiated SNP clusters (99.99th
21 percentile; $F_{ST} > 0.0452$) were listed in **Table 2**. Although it is reported that F_{ST} of
22 individual-marker was too variable ([Weir et al. 2005](#)), the four most significant regions
23 (7q21, 6p21-22, 1q22 and 11p15) carrying multiple highly differentiated SNPs should be

1 indicative. The 7q21 region harbors two genes: *CD36* and *SEMA3C*. *CD36* directly
2 mediates cytoadherence of *Plasmodium falciparum* parasitized erythrocytes and it binds
3 long chain fatty acids and may function in the transport and/or as a regulator of fatty acid
4 transport ([Oquendo et al. 1989](#); [Baruch et al. 1996](#); [Erdman et al. 2009](#)). It was reported
5 that *CD36* has been subjected to positive selection for malaria or some unknown selection
6 pressures ([Aitman et al. 2000](#); [Omi et al. 2002](#); [Omi et al. 2003](#); [Fry et al. 2009](#)). *SEMA3C*,
7 induced by *ADAMTS1*, promotes the migration of cancer cells ([Esselens et al. 2010](#)). The
8 6p21-22 region, harboring human major histocompatibility complex (MHC), also showed
9 an over-representation of African ancestry in this study (although less than 3SD) and has
10 been reported under selection in various studies ([Garrigan and Hedrick 2003](#); [Tang et al.](#)
11 [2007](#)). *MUC1*, located in 1q22, is involved in signaling by PDGF and serves a protective
12 function by binding to pathogens ([Davila et al. 2010](#); [Li et al. 2010](#)). *HBB* and *HBD*,
13 located in 11p15, have been subjected to balancing selections because their mutations
14 protected against malaria according to numerous studies ([Ashley-Koch et al. 2000](#); [Wood](#)
15 [et al. 2005](#)).

16 Ingenuity pathway analysis (IPA) was particularly helpful in exploring the function
17 and pathways of the selection-candidate genes in the context of higher-order cellular and
18 molecular mechanisms. The 402 SNPs with the highest F_{ST} (99.90th percentile; F_{ST}
19 >0.0287) were subjected to IPA analysis, whose results showed that genes involved in
20 metabolic diseases were the most significantly enriched ($P = 1.51 \times 10^{-16}$) among all
21 function classes, followed by the genes involved in endocrine system disorder ($P = 2.23$
22 $\times 10^{-16}$), immunological diseases ($P = 9.30 \times 10^{-12}$) and genetic disorder ($P = 5.67 \times 10^{-11}$).
23 Antigen presentation pathway was most significantly enriched ($P = 1.95 \times 10^{-4}$), followed

1 by allograft rejection signaling ($P = 4.69 \times 10^{-3}$), Graft-versus-host disease signaling ($P =$
2 4.69×10^{-3}) and autoimmune thyroid diseases signaling ($P = 5.35 \times 10^{-3}$). All of the four
3 aforementioned pathways are related to immune system, which might reflect a great
4 environmental differentiation between Sub-Saharan Africa and North America. We also
5 conducted IPA on 4,011 SNPs showing the highest F_{ST} (99.00th percentile; $F_{ST} > 0.0162$),
6 which yielded results similar to that with $F_{ST} > 0.0287$ (99.90th percentile). And two
7 additional pathways emerged: IL-9 signaling pathway ($P = 8.01 \times 10^{-3}$) and EGF signaling
8 pathway ($P = 6.38 \times 10^{-3}$).

9 *Reconstituted African Americans and its difference with AfA*

10 We then compared the AfA genome with that of reconstituted African American
11 (rAfA) using genotypes of YRI and CEU. The rationale behind is that the former have been
12 subjected to possible natural selection, while the latter have not. In particular, the allele
13 frequencies of rAfA were estimated at each locus using YRI and CEU for the given level of
14 admixture (21.68%) assuming no natural selection after admixture. The candidates
15 identified under selection in this way could avoid potential errors introduced in the
16 inference of ancestry. The genome-wide distribution of F_{ST} between AfA and rAfA was
17 shown in **Figure 2B**. Overall 81% of the SNPs showing highest difference between AAF
18 and IAF could be validated in the comparison between AfA and rAfA. In particular, these
19 main F_{ST} peaks between AAF and YRI are essentially the same as those between AfA and
20 rAfA, which indicated that almost all the selection signals identified by comparing AfA
21 and rAfA were originated from AAF or IAF. We also used CAU-GWAS instead of CEU to
22 construct rAfA, with the genomic distribution of F_{ST} similar to that using CEU (**Figure**
23 **S10**).

1 Next, we reconstructed a rAfA population using a set of putative parental
2 populations. First, we constructed an African parental population of AfA (APP) using 64%
3 Yoruba, 19% Mandenka and 14% Bantu according to previous reported ancestral
4 contribution to AAF ([Zakharia et al. 2009](#)). Then we used genotypes of the APP and
5 CAU-GWAS to reconstruct rAfA. We obtained the genome-wide distribution of F_{ST}
6 between AfA and rAfA ([Figure S11](#)). Although the population differentiation between this
7 rAfA and AfA is higher than that based on two aforementioned rAfAs, the genome-wide
8 distribution of F_{ST} is similar to that using only two pure parental populations.

9 *Further evidences for positive selection in AAF and African*

10 Since the force of positive selection acts in a locus-specific manner and tends to
11 increase F_{ST} ([Nielsen 2005](#)), we hypothesized that positive selection preferentially acted
12 upon functionally important loci over the others in the genome, thus leading to an
13 enrichment of functional SNPs in the high F_{ST} bin ([Barreiro et al. 2008](#)). Here, we
14 investigated the enrichment of different SNP classes among the high F_{ST} bin (top 1st
15 percentile among all SNPs, $F_{ST} > 0.0164$) between AAF and YRI.

16 The SNPs, based on their location and function relative to the genes, were classified
17 into nongenic, genic, intronic, 3'UTR, 5'UTR, synonymous, non-synonymous, coding,
18 transcriptonic, near-gene-3 and near-gene-5 according to UCSC annotation. We found that
19 F_{ST} distributions of each SNP category were not significantly different from those of
20 non-genic SNPs. However, the proportion of genic SNPs among high F_{ST} bin (top 1st
21 percentile among all SNPs, $F_{ST} > 0.0164$) is significantly higher than that of nongenic
22 SNPs (χ^2 test, $P=0.046$; [Figure 3](#)). Notably, this excess is particularly marked for
23 transcriptonic SNPs (χ^2 test, $P=0.004$; [Figure 3](#)). The proportion of synonymous SNPs in

1 high F_{ST} bin was 1.22-fold higher than the expectation under neutrality, which could
2 attribute to the linkage disequilibrium of those SNPs with loci under selection (a
3 phenomenon known as hitchhiking). Similar observations could be made when the
4 thresholds for high F_{ST} bin were set at top 5% ($F_{ST} > 0.0083$) or top 0.1% ($F_{ST} > 0.0304$)
5 (Figure S12, S13), and the conclusions still held when the SNPs with MAF > 0.05 in both
6 YRI and AAF were examined (Figure S14). The significant enrichment of high F_{ST} loci in
7 the SNP categories with genetic functions supports the presence of positive selection either
8 in AAF or in YRI, or both.

9 **Discussion**

10 Admixed populations such as African Americans provide a unique opportunity to
11 study very recent natural selection, as their genomes are donated by long-diverged
12 continental ancestries and may have been subjected to novel environmental challenges.
13 The first strategy, detecting excessive or decreased ancestry contribution from its ancestral
14 parental populations, has been used in several recent studies (Tang et al. 2007; Basu et al.
15 2008; Bryc et al. 2010). However, natural selections before admixture cannot be detected
16 by this approach because the distribution of ancestry across the admixed genomes would
17 not be affected by such selections. Therefore, we developed a new strategy, which
18 examined selection since the African ancestry left for America (including selection both
19 pre-and post-admixture). The candidate genes under selection identified by our new
20 strategy were also confirmed in an analysis of allele frequency differentiation between
21 rAfA and AfA (although the two analyses are not completely independent).

22 Above all, we have identified six regions showing excess of African or European
23 ancestry using the first strategy. The highest ancestry deviation among all regions showing

1 excess of African and European ancestry is <0.026 , which is lower than the values of any
2 previous studies on the admixed populations in the New World ([WORKMAN et al. 1963](#);
3 [Reed 1969](#); [Blumberg and Hesser 1971](#); [Tang et al. 2007](#); [Basu et al. 2008](#); [Bryc et al.](#)
4 [2010](#)). For example, the African ancestry deviation even exceeds 0.14 in Puerto Rican
5 according to a study by Tang *et al.* ([2007](#)), which is more than five-fold higher than that in
6 this study. We attribute the low ancestry deviation to the much larger sample size, denser
7 markers, more accurate phased data using trios, and more powerful statistical methods
8 used in this study. Based on the maximum deviation regions, we estimated that the highest
9 selection coefficient is approximate to be 0.002 (assuming 12 generations since admixture
10 under CGF model). In fact, the real selection coefficient could be much lower than 0.002
11 considering the statistics and evolutionary noises. These results are reasonable considering
12 the fact that there is much lower mortality rate of AfA in recent 200 years compared with
13 that during Atlantic slave trade.

14 Secondly, we identified a large number of genes with highly differentiated allele
15 frequencies between AAF and YRI using our new approach. These genes do not overlap
16 with those identified by the former approach, suggesting the different environmental
17 pressures AfA experienced before and after population admixture. IPA Analysis of SNPs
18 with high population differentiation between AAF and YRI showed that genes involved in
19 metabolic diseases ($P = 1.51 \times 10^{-16}$), endocrine system disorder ($P = 2.23 \times 10^{-16}$),
20 immunological diseases ($P = 9.30 \times 10^{-12}$) and genetic disorder ($P = 5.67 \times 10^{-11}$) were
21 significantly enriched. Especially, we found that genes such as *PSCA*, *ZP4*, *AKAP12* were
22 associated with AfA specific high-risk diseases such as hypertension and prostate cancer
23 ([Smith and O'Brien 2005](#); [Goran 2008](#)). Five genes (*CD36*, *HBB*, *HBD*, *HLA-B*, *HLA-DR*),

1 whose mutations protect against malaria, are also located in the highly differentiated
2 regions between AAF and IAF.

3 Compared with Caucasians, African Americans have higher mortality rate for all
4 cancers combined and for most major cancers ([Jemal et al. 2006](#)), as well as higher risk of
5 obesity-related disease such as diabetes, hypertension and prostate cancer ([Goran 2008](#)).
6 Interestingly, many genes located in selection candidate regions identified by the novel
7 approach are associated with AfA ethnic high-risk diseases such as hypertension, prostate
8 cancer and systemic sclerosis. Especially, One of the most significantly differentiated
9 SNPs (rs2294008; $F_{ST} = 0.04561$) between AAF and YRI, located in 8q24, is a missense
10 mutation c.57T>C (p.Met1Thr) in *PSCA* and was reported to be associated with gastric and
11 bladder cancer ([Sakamoto et al. 2008](#); [Matsuo et al. 2009](#); [Wu et al. 2009](#)). Many studies
12 also reported that multiple loci in 8q24 were associated with prostate cancer in AfA
13 ([Freedman et al. 2006](#); [Al Olama et al. 2009](#); [Yeager et al. 2009](#)). We proposed a
14 hypothesis that most of the genes associated with AfA ethnic diseases may have played an
15 important role in AfA's adaptation to local environment and thus show higher population
16 differentiation between AAF and IAF. Further analysis of the 8q24 region would provide
17 new insights into the etiology and evolutionary history of these cancers.

18 Among the selection candidate genes detected by genome-wide locus specific F_{ST}
19 between AAF and YRI, five genes (*CD36*, *HBB*, *HBD*, *HLA-B*, *HLA-DR*) have been
20 reported subjected to natural selection probably due to malaria ([Kwiatkowski 1999](#);
21 [Kwiatkowski 2005](#)). Because of the strong selection pressure of malaria, loss-of-function
22 or abnormality of these genes was supposed to increase the survival rate of individual
23 living in Africa. Some mutations in these genes have reached much higher frequencies in

1 Africans compared with that in areas of low incidence of malaria. However, these
2 mutations that defend malaria could become disadvantage in AfA because the malaria was
3 no longer a strong selection pressure in North America and these mutations could even lead
4 to morbidity or mortality (Platt et al. 1994). We hypothesize that frequencies of these
5 mutations would have decreased in AAF compared with those in indigenous African due to
6 their disadvantage in AfA.

7 Next, we examined this hypothesis in the empirical data. Since the functional
8 mutations in these genes were not genotyped in this study, we examined the SNPs strongly
9 linked (linkage disequilibrium) with these mutations instead of these mutations
10 themselves. It is well known that rs3211938 is a nonsense mutation c.1389T>G
11 (p.Tyr325X) in *CD36* and has been subjected to natural selection because of malaria or
12 some other environmental factors in African (Aitman et al. 2000; Omi et al. 2003; Erdman
13 et al. 2009; Fry et al. 2009). We found three SNPs that are highly differentiated between
14 AAF and YRI (Table S1) in 7q21 are strongly linked with rs3211938 (each with $r^2 > 0.4$ in
15 YRI). Interestingly, we did observe that the frequencies of the alleles linked with
16 rs3211938 (G), which is the derived allele, were much lower in AAF compared with that in
17 YRI (Table S1). Another example, rs334 is a missense mutation c.70A>T (p.Glu7Val) in
18 *HBB*, which leads to sickle cell anemia [MIM 603903] (Ashley-Koch et al. 2000;
19 Winichagoon et al. 2000; Wood et al. 2005), one of the most well-studied genetic disorder.
20 rs7952293, one of the SNPs showing high F_{ST} between AAF and YRI, was strongly linked
21 with rs334 ($r^2 = 0.237$ in YRI). In particular, the haplotype constructed by rs7952293(A)
22 and rs334(T) accounted for 86.67% of the haplotypes containing rs334(T) in YRI. We
23 observed that the frequency of rs7952293(A) in AAF (0.2261) was lower compared with

1 that in YRI (0.3172), which also supports the hypothesis that frequencies of alleles
2 protecting against malaria in AAF are lower than those in indigenous African. The other
3 three genes were not examined using the same procedure because the frequencies of
4 mutations on these genes are too low to find strong-linked representative SNPs.

5 Our study takes advantage of both large sample size and high-density genome-wide
6 data. However, our analysis demonstrated that the maximum deviation showing excess of
7 African or European ancestry was small (<2.6%). Detecting such weak selection signals in
8 admixed population such as AfA is a big challenge, which needs large sample size as in
9 this study, or even larger sample size, to distinguish the real signals from the ancestry
10 deviation caused by genetic drift and sampling error. Therefore, we propose that any study
11 in the future trying to detect such weak selection signals in AfA or other recently admixed
12 populations should collect at least thousands of samples. With the new strategy, we
13 detected a lot of genes associated with AfA specific high-risk diseases such as
14 hypertension and prostate cancer, and we also detected five genes whose mutations are
15 against malaria. This new approach is powerful in detecting natural selection both before
16 and after the establishment of AfA and can be applied to other admixed populations such as
17 Latinos in the New world and the Uyghurs in Asia ([Xu et al. 2008](#); [Xu and Jin 2008](#); [Xu et](#)
18 [al. 2009](#)).

19 **Methods**

20 *Data assemble and quality control*

21 The genotypic data were obtained from International HapMap Project (HapMap;
22 <http://www.hapmap.org>), Human Genome Diversity Project (HGDP;

1 <http://www.cephb.fr/en/hgdp>), and Illumina iControlDB (<http://www.illumina.com>),
2 respectively. The combined data set processed with PEAS v1.0 ([Xu et al. 2010](#)) includes
3 the genotypic data of 588 HapMap samples (87 ASW, 167 YRI, 165 CEU, 85 CHD, 84
4 CHB) ([Frazer et al. 2007](#)), 300 HGDP samples (156 indigenous European, 102 indigenous
5 African and 42 Amerindian) ([Li et al. 2008](#)), 2,161 AfA and 3,294 Caucasians (referred to
6 as CAU-GWAS) from iControlDB genotyped by Illumina 550K Beadarray. The samples
7 from iControlDB have passed Illumina's rigorous quality control and have been used as
8 controls in five genome-wide association studies.

9 ASW from HapMap and AfA from iControlDB were merged into one AfA
10 population, and Yoruba from HGDP was collected from Nigeria and therefore was merged
11 with YRI from HapMap in subsequent analyses. The following samples were removed
12 from the data set: (1) the relatives based on sample information or PLINK ([Purcell et al.](#)
13 [2007](#)) result ($IBD > 0.2$), (2) individuals identified as outliers of each population (except
14 AfA) based on the top ten principal components of PCA analysis ($SD > 6$), (3) AfA
15 individuals with $> 2\%$ Native-American/East-Asian, (4) AfA individuals with $> 99\%$
16 European contribution, which are likely to be descendents of individuals of European
17 ancestry, (5) AfA individuals with $> 99\%$ African contribution, in which recent African
18 immigrants could not be practically identified ([Bryc et al. 2010](#)).

19 These filtered samples described above were merged and the SNPs sharing
20 reference SNP ID (rs) and vendor-specified strands were kept in combined data. Then, the
21 data set was further filtered for individuals with $> 10\%$ missing genotypes and SNPs with
22 $> 10\%$ missing data, as well as Hardy-Weinberg disequilibrium ($P < 2 \times 10^{-6}$) within each
23 population except AfA. The final data set comprised of 503,694 SNPs (491,526 autosomal

1 SNPs) shared by 5,210 individuals from 21 population groups, with total genotyping call
2 rate 99.74%.

3 *Populations and samples*

4 Overall, 1,890 unrelated AfA samples with ignorable ancestry outside of Africa and
5 Europe were studied, among which 1,838 individuals were downloaded from iControlDB,
6 and another 52 individuals were from the International HapMap Project ([Altshuler et al.](#)
7 [2010](#)). In addition, 113 YRI from HapMap ([Altshuler et al. 2010](#)) and 102 indigenous
8 Africans in 7 different groups collected from HGDP ([Li et al. 2008](#)) were merged,
9 representing the extant Africans, while 113 CEU from HapMap ([Altshuler et al. 2010](#))
10 were merged with 156 Europeans in 8 different groups collected from HGDP, representing
11 the extant Europeans. In total 2,648 CAU-GWAS samples from iControlDB were taken as
12 another representation of the extant Europeans. In addition, 84 CHB and 85 CHD
13 represented populations from East Asia, and 24 pure Amerindians from HGDP represented
14 Native American.

15 *Population genetic analysis*

16 In order to reduce the linkage disequilibrium (LD) between markers, those with r^2
17 >0.5 were removed, calculated in sliding window of 50 SNPs and shifted every five SNPs
18 ([see Supplemental Text](#)). This process reduced the original dataset to 341,672 autosomal
19 SNPs. Based on these thinned markers, principle component analysis (PCA) was
20 performed at the individual level using *smartpca*, from the package EIGENSOFT
21 ([Patterson et al. 2006](#)). Individual ancestry proportion was estimated using FRAPPE,
22 which implements an expectation-maximization (EM) algorithm ([Tang et al. 2005](#)).
23 FRAPPE was run on all 491,526 SNPs with 10,000 iterations by setting K from 2 to 4. We

1 also ran STRUCTURE (Falush et al. 2003) on SNPs with inter-marker distance more than
2 1M (see Supplemental Text). Genetic difference between populations was measured using
3 F_{ST} following Weir and Cockerham (Weir and Cockerham 1984), which accounts for
4 differences in the sample size in each population. The locus-specific F_{ST} between any two
5 populations was also calculated using the same formula.

6 *Locus-specific ancestry inference*

7 Various methods and software have been developed for inferring locus-specific
8 ancestry based on high density SNPs data, such as ANCESTRYMAP (Patterson et al.
9 2004), SABER (Tang et al. 2006), LAMP and LAMP-ANC (Sankararaman et al. 2008b),
10 uSWITCH and uSWITCH-ANC (Sankararaman et al. 2008a), HAPAA (Sundquist et al.
11 2008) and HAPMIX (Price et al. 2009). A simple analysis showed that HAPMIX
12 outperformed other methods with implemented software based on our simulated data (see
13 Supplemental Text). Therefore, HAPMIX was used to infer locus-specific ancestry in
14 African Americans in this study. We also used LAMP-ANC to do a similar analysis since it
15 also performs very well.

16 The phased data of HapMap 3 were downloaded from HapMap website (Altshuler et
17 al. 2010). We used haplotypes of 88 CEU and 88 YRI (all from trio samples) representing
18 the European and African ancestral populations, respectively in the subsequent analysis.
19 The mean European ancestry proportion in AfA (θ), which is a required input parameter
20 for HAPMIX, was based on the estimation of FRAPPE. Generation since admixture (λ)
21 with the largest likelihood was taken as its estimation. By running HAPMIX in diploid
22 mode, we obtained the haplotypes and ancestry segments for each AfA individual. Then
23 we reconstructed an ancestral African population (AAF) and an ancestral European

1 population (AEU) using inferred chromosomal segments of African ancestry and those of
2 European ancestry in AfA, respectively. In brief, we constructed AAF using only those
3 chromosomal segments with pure African ancestry. For each given SNP, allele frequency
4 of AAF was calculated based on the available genotypes with African ancestry across all
5 AfA. This procedure was also applied to the construction of AEU using chromosomal
6 segments with pure European ancestry.

7 *Simulation of AfA and its parental populations*

8 Under selective neutrality, genetic drift of the admixed population and its parental
9 populations contribute to the variation of ancestry proportion and population
10 differentiation among loci ([Weir and Cockerham 1984](#); [Long 1991](#)). Therefore, we
11 performed an extensive forward-time simulation to explore the potential impact of genetic
12 drift on the locus-specific population differentiations between AAF and YRI, as well as
13 that between AEU and CEU. In this simulation, recombination was introduced according
14 to the genetic map adapted from HapMap (release #22) ([Frazer et al. 2007](#)), and mutation
15 was ignored given the short history of AfA. The effective population sizes (N_e) of each
16 population was obtained from the HapMap website ([Altshuler et al. 2010](#)). In particular, N_e
17 for African, European and AfA were set to 17,094, 11,418 and 17,094, respectively.
18 Continuous-gene-flow (CGF) model ([Figure S15](#)) was used based on previous studies
19 ([Pfaff et al. 2001](#); [Price et al. 2009](#)).

20 The aforementioned phased data of 88 YRI and 88 CEU were taken respectively as
21 the genotypes of common African ancestry and common European ancestry before
22 population admixture. We simulated individuals of AfA by constructing their genomes
23 from a mosaic of haploid YRI and haploid CEU genomes. The generations since admixture

1 (λ) was set to 12 according to CGF model based on previous reports ([Pfaff et al. 2001](#);
2 [Price et al. 2009](#)), which was also supported by our observations ([Figure S6](#)). The gene
3 flow (α) that African ancestry received from European each generation was calculated by
4 $\alpha = 1 - (m_1)^{1/t}$, in which m_1 represents the mean proportion of European ancestry in AfA.
5 Both African and European parental populations evolved 12 generations simultaneously.
6 Finally, genotypes of 113 European, 113 African and 1,890 AfA, as well as their
7 SNP-specific ancestral status, were output to match the sample sizes of the data. Based on
8 the primary simulation, we performed extended simulations by setting a bottleneck event
9 in the first generation during which N_e for AfA were reduced to 8,000, 5,000, 3,000, 2,000
10 and 1,000, respectively.

11 *Function annotations and ingenuity pathway analysis (IPA)*

12 Genomic regions that deviated from genome-wide ancestral contributions or with
13 extremely high F_{ST} were annotated based on HapMap website ([Thorisson et al. 2005](#)). The
14 SNPs showing substantial population differentiation were interrogated for network and
15 functional interrelatedness using the Ingenuity Pathway Analysis (IPA) version 8.5
16 software tools. This software searches for information on genes in Ingenuity Pathways
17 Knowledge Base, a repository of molecular interactions, regulatory events,
18 gene-to-phenotype associations, and chemical knowledge, all collected from the full text of
19 the peer-reviewed life sciences literatures. With IPA, we can analyze data in the context of
20 molecular mechanisms, identify key mechanistic differences between subpopulations, and
21 further relate molecular events to higher-order cellular and disease processes.

22 *Statistical analysis*

1 All statistical computation and graphics were performed using R version 2.9 (Ihaka
2 and Gentleman 1996). Kolmogorov–Smirnov tests were performed to compare the
3 empirical distributions of locus-specific F_{ST} with those simulated. And Chi-squared (χ^2
4 test) tests were performed to test the over-representation of each SNP category compared
5 with nongenic SNPs among the high F_{ST} bin.

6 **Acknowledgements**

7 SX was supported by the National Science Foundation of China (30971577, 31171218),
8 Shanghai Rising-Star Program (11QA1407600), and Science Foundation of The Chinese
9 Academy of Sciences (KSCX2-EW-Q-1-11, KSCX2-EW-R-01-05,
10 KSCX2-EW-J-15-05). LJ was supported by the National Science Foundation of China
11 (30890034, 30625016) and the Science and Technology Commission of Shanghai
12 Municipality (09540704300). SX is a Max-Planck Independent Junior Research Group
13 Leader and a member of CAS Youth Innovation Promotion Association. SX also gratefully
14 acknowledges the support of K.C.Wong Education Foundation, Hong Kong. BW is a
15 senior author of Harvard group. This work was also supported by the MoST International
16 Cooperation Base of China.

17 **Authors' contributions**

18 **S.X.** and **L.J.** conceived and designed the study. **H.W.** collected genotype data from the
19 Illumina iControl Database (iControlDB). **W.J.** performed data analysis, with contribution
20 from **S.X.**, **Y.Y.**, **Y.S.** and **B.W.** performed IPA analysis. **S.X.** and **W.J.** interpreted the
21 data. **W.J.**, **S.X.** and **L.J.** wrote the paper. All authors read and approved the final
22 manuscript. All authors declare that no competing financial interests exist.

1

2

1 References

- 2 Adams, J. and R.H. Ward. 1973. Admixture studies and the detection of selection. *Science*
3 180: 1137-1143.
- 4 Aitman, T.J., L.D. Cooper, P.J. Norsworthy, F.N. Wahid, J.K. Gray, B.R. Curtis, P.M.
5 McKeigue, D. Kwiatkowski, B.M. Greenwood, R.W. Snow et al. 2000. Malaria
6 susceptibility and CD36 mutation. *Nature* 405: 1015-1016.
- 7 Akey, J.M. 2009. Constructing genomic maps of positive selection in humans: where do
8 we go from here? *Genome Res* 19: 711-722.
- 9 Akey, J.M., G. Zhang, K. Zhang, L. Jin, and M.D. Shriver. 2002. Interrogating a
10 high-density SNP map for signatures of natural selection. *Genome Res* 12:
11 1805-1814.
- 12 Al Olama, A.A., Z. Kote-Jarai, G.G. Giles, M. Guy, J. Morrison, G. Severi, D.A.
13 Leongamornlert, M. Tymrakiewicz, S. Jhavar, E. Saunders et al. 2009. Multiple
14 loci on 8q24 associated with prostate cancer susceptibility. *Nat Genet* 41:
15 1058-1060.
- 16 Altshuler, D.M., R.A. Gibbs, L. Peltonen, D.M. Altshuler, R.A. Gibbs, L. Peltonen, E.
17 Dermitzakis, S.F. Schaffner, F. Yu, L. Peltonen et al. 2010. Integrating common
18 and rare genetic variation in diverse human populations. *Nature* 467: 52-58.
- 19 Ashley-Koch, A., Q. Yang, and R.S. Olney. 2000. Sickle hemoglobin (HbS) allele and
20 sickle cell disease: a HuGE review. *Am J Epidemiol* 151: 839-845.
- 21 Balaesque, P.L., S.J. Ballereau, and M.A. Jobling. 2007. Challenges in human genetic
22 diversity: demographic history and adaptation. *Hum Mol Genet* 16 (R2): R134-139.
- 23 Barreiro, L.B., G. Laval, H. Quach, E. Patin, and L. Quintana-Murci. 2008. Natural
24 selection has driven population differentiation in modern humans. *Nat Genet* 40:
25 340-345.
- 26 Baruch, D.I., J.A. Gormely, C. Ma, R.J. Howard, and B.L. Pasloske. 1996. Plasmodium
27 falciparum erythrocyte membrane protein 1 is a parasitized erythrocyte receptor for
28 adherence to CD36, thrombospondin, and intercellular adhesion molecule 1. *Proc*
29 *Natl Acad Sci U S A* 93: 3497-3502.
- 30 Basu, A., H. Tang, X. Zhu, C.C. Gu, C. Hanis, E. Boerwinkle, and N. Risch. 2008.
31 Genome-wide distribution of ancestry in Mexican Americans. *Hum Genet* 124:
32 207-214.
- 33 Blumberg, B.S. and J.E. Hesser. 1971. Loci differentially affected by selection in two
34 American black populations. *Proc Natl Acad Sci U S A* 68: 2554-2558.
- 35 Bryc, K., A. Auton, M.R. Nelson, J.R. Oksenberg, S.L. Hauser, S. Williams, A. Froment,
36 J.-M. Bodo, C. Wambebe, S.A. Tishkoff et al. 2010. Genome-wide patterns of
37 population structure and admixture in West Africans and African Americans. *Proc*
38 *Natl Acad Sci U S A* 107: 786-791.
- 39 Cui, Y., M. Tian, M. Zong, M. Teng, Y. Chen, J. Lu, J. Jiang, X. Liu, and J. Han. 2009.
40 Proteomic analysis of pancreatic ductal adenocarcinoma compared with normal
41 adjacent pancreatic tissue and pancreatic benign cystadenoma. *Pancreatology* 9:
42 89-98.
- 43 Davila, S., F.E. Froeling, A. Tan, C. Bonnard, G.J. Boland, H. Snippe, M.L. Hibberd, and
44 M. Seielstad. 2010. New genetic associations detected in a host response study to
45 hepatitis B vaccine. *Genes Immun* 11: 232-238.

- 1 Erdman, L.K., G. Cosio, A.J. Helmers, D. Gowda, S. Grinstein, and K.C. Kain. 2009.
2 CD36 and TLR Interactions in Inflammation and Phagocytosis: Implications for
3 Malaria. *J Immunol* 183: 6452-6459.
- 4 Esselens, C., J. Malapeira, N. Colome, C. Casal, J.C. Rodriguez-Manzaneque, F. Canals,
5 and J. Arribas. 2010. The cleavage of semaphorin 3C induced by ADAMTS1
6 promotes cell migration. *J Biol Chem* 285: 2463-2473.
- 7 Falush, D., M. Stephens, and J.K. Pritchard. 2003. Inference of population structure using
8 multilocus genotype data: linked loci and correlated allele frequencies. *Genetics*
9 164: 1567-1587.
- 10 Frazer, K.A. D.G. Ballinger D.R. Cox D.A. Hinds L.L. Stuve R.A. Gibbs J.W. Belmont A.
11 Boudreau P. Hardenbol S.M. Leal et al. 2007. A second generation human
12 haplotype map of over 3.1 million SNPs. *Nature* 449: 851-861.
- 13 Freedman, M.L., C.A. Haiman, N. Patterson, G.J. McDonald, A. Tandon, A. Waliszewska,
14 K. Penney, R.G. Steen, K. Ardlie, E.M. John et al. 2006. Admixture mapping
15 identifies 8q24 as a prostate cancer risk locus in African-American men. *Proc Natl*
16 *Acad Sci U S A* 103: 14068-14073.
- 17 Fry, A.E., A. Ghansa, K.S. Small, A. Palma, S. Auburn, M. Diakite, A. Green, S. Campino,
18 Y.Y. Teo, T.G. Clark et al. 2009. Positive selection of a CD36 nonsense variant in
19 sub-Saharan Africa, but no association with severe malaria phenotypes. *Hum Mol*
20 *Genet* 18: 2683-2692.
- 21 Garrigan, D. and P.W. Hedrick. 2003. Perspective: detecting adaptive molecular
22 polymorphism: lessons from the MHC. *Evolution* 57: 1707-1722.
- 23 Goran, M.I. 2008. Ethnic-specific pathways to obesity-related disease: the Hispanic vs.
24 African-American paradox. *Obesity (Silver Spring)* 16: 2561-2565.
- 25 Hancock, A.M., D.B. Witonsky, A.S. Gordon, G. Eshel, J.K. Pritchard, G. Coop, and A. Di
26 Rienzo. 2008. Adaptations to climate in candidate genes for common metabolic
27 disorders. *PLoS Genet* 4: e32.
- 28 Ihaka, R. and R. Gentleman. 1996. R: A language for data analysis and graphics. *J Comput*
29 *Graph Statist* 5: 299-314.
- 30 Jemal, A., R. Siegel, E. Ward, T. Murray, J. Xu, C. Smigal, and M.J. Thun. 2006. Cancer
31 statistics, 2006. *CA Cancer J Clin* 56: 106-130.
- 32 Kao, W.H., M.J. Klag, L.A. Meoni, D. Reich, Y. Berthier-Schaad, M. Li, J. Coresh, N.
33 Patterson, A. Tandon, N.R. Powe et al. 2008. MYH9 is associated with nondiabetic
34 end-stage renal disease in African Americans. *Nat Genet* 40: 1185-1192.
- 35 Kimura, M. 2003. The neutral theory of molecular evolution. *Cambridge University Press,*
36 *Cambridge, United Kingdom.*
- 37 Kresse, S.H., H.O. Ohnstad, E.B. Paulsen, B. Bjerkehagen, K. Szuhai, M. Serra, K.L.
38 Schaefer, O. Myklebost, and L.A. Meza-Zepeda. 2009. LSAMP, a novel candidate
39 tumor suppressor gene in human osteosarcomas, identified by array comparative
40 genomic hybridization. *Genes Chromosomes Cancer* 48: 679-93.
- 41 Kwiatkowski, D. 1999. The molecular genetic approach to malarial pathogenesis and
42 immunity. *Parassitologia* 41: 233-240.
- 43 Kwiatkowski, D.P. 2005. How malaria has affected the human genome and what human
44 genetics can teach us about malaria. *Am J Hum Genet* 77: 171-192.
- 45 Li, J.Z., D.M. Absher, H. Tang, A.M. Southwick, A.M. Casto, S. Ramachandran, H.M.
46 Cann, G.S. Barsh, M. Feldman, L.L. Cavalli-Sforza et al. 2008. Worldwide human

- 1 relationships inferred from genome-wide patterns of variation. *Science* 319:
2 1100-1104.
- 3 Li, Y., D.L. Dinwiddie, K.S. Harrod, Y. Jiang, and K.C. Kim. 2010. Anti-inflammatory
4 effect of MUC1 during respiratory syncytial virus infection of lung epithelial cells
5 in vitro. *Am J Physiol Lung Cell Mol Physiol* 298: L558-563.
- 6 Long, J.C. 1991. The genetic structure of admixed populations. *Genetics* 127: 417-428.
- 7 Matsuo, K., K. Tajima, T. Suzuki, T. Kawase, M. Watanabe, K. Shitara, K. Misawa, S. Ito,
8 A. Sawaki, K. Muro et al. 2009. Association of prostate stem cell antigen gene
9 polymorphisms with the risk of stomach cancer in Japanese. *Int J Cancer* 125:
10 1961-1964.
- 11 McAllister, C.S., A.M. Toth, P. Zhang, P. Devaux, R. Cattaneo, and C.E. Samuel. 2010.
12 Mechanisms of protein kinase PKR-mediated amplification of beta interferon
13 induction by C protein-deficient measles virus. *J Virol* 84: 380-386.
- 14 Meltzer, M. 1993. *Slavery: A World History*. Da Capo Press, New York.
- 15 Nielsen, R. 2005. Molecular signatures of natural selection. *Annu Rev Genet* 39: 197-218.
- 16 Nielsen, R., I. Hellmann, M. Hubisz, C. Bustamante, and A.G. Clark. 2007. Recent and
17 ongoing selection in the human genome. *Nat Rev Genet* 8: 857-868.
- 18 Oleksyk, T.K., M.W. Smith, and S.J. O'Brien. 2010. Genome-wide scans for footprints of
19 natural selection. *Philos Trans R Soc Lond B Biol Sci* 365: 185-205.
- 20 Omi, K., J. Ohashi, J. Patarapotikul, H. Hananantachai, I. Naka, S. Looareesuwan, and K.
21 Tokunaga. 2002. Fcγ receptor IIA and IIIB polymorphisms are associated
22 with susceptibility to cerebral malaria. *Parasitol Int* 51: 361-366.
- 23 Omi, K., J. Ohashi, J. Patarapotikul, H. Hananantachai, I. Naka, S. Looareesuwan, and K.
24 Tokunaga. 2003. CD36 polymorphism is associated with protection from cerebral
25 malaria. *Am J Hum Genet* 72: 364-374.
- 26 Oquendo, P., E. Hundt, J. Lawler, and B. Seed. 1989. CD36 directly mediates
27 cytoadherence of Plasmodium falciparum parasitized erythrocytes. *Cell* 58:
28 95-101.
- 29 Patterson, N., N. Hattangadi, B. Lane, K.E. Lohmueller, D.A. Hafler, J.R. Oksenberg, S.L.
30 Hauser, M.W. Smith, S.J. O'Brien, D. Altshuler et al. 2004. Methods for
31 high-density admixture mapping of disease genes. *Am J Hum Genet* 74: 979-1000.
- 32 Patterson, N., A.L. Price, and D. Reich. 2006. Population structure and eigenanalysis.
33 *PLoS Genet* 2: e190.
- 34 Pereira, R.M., K.L. Teixeira, V. Barreto-de-Souza, T.C. Calegari-Silva, L.D. De-Melo,
35 D.C. Soares, D.C. Bou-Habib, A.M. Silva, E.M. Saraiva, and U.G. Lopes. 2010.
36 Novel role for the double-stranded RNA-activated protein kinase PKR: modulation
37 of macrophage infection by the protozoan parasite Leishmania. *Faseb J* 24:
38 617-626.
- 39 Pfaff, C.L., E.J. Parra, C. Bonilla, K. Hiester, P.M. McKeigue, M.I. Kamboh, R.G.
40 Hutchinson, R.E. Ferrell, E. Boerwinkle, and M.D. Shriver. 2001. Population
41 structure in admixed populations: effect of admixture dynamics on the pattern of
42 linkage disequilibrium. *Am J Hum Genet* 68: 198-207.
- 43 Platt, O.S., D.J. Brambilla, W.F. Rosse, P.F. Milner, O. Castro, M.H. Steinberg, and P.P.
44 Klug. 1994. Mortality in sickle cell disease. Life expectancy and risk factors for
45 early death. *N Engl J Med* 330: 1639-1644.

- 1 Praml, C., W. Schulz, A. Claas, J. Mollenhauer, A. Poustka, R. Ackermann, M. Schwab,
2 and K.O. Henrich. 2008. Genetic variation of Aflatoxin B1 aldehyde reductase
3 genes (AFAR) in human tumour cells. *Cancer Lett* 272: 160-166.
- 4 Price, A.L., A. Tandon, N. Patterson, K.C. Barnes, N. Rafaels, I. Ruczinski, T.H. Beaty, R.
5 Mathias, D. Reich, and S. Myers. 2009. Sensitive detection of chromosomal
6 segments of distinct ancestry in admixed populations. *PLoS Genet* 5: e1000519.
- 7 Price, A.L., M.E. Weale, N. Patterson, S.R. Myers, A.C. Need, K.V. Shianna, D. Ge, J.I.
8 Rotter, E. Torres, K.D. Taylor et al. 2008. Long-range LD can confound genome
9 scans in admixed populations. *Am J Hum Genet* 83: 132-135.
- 10 Pritchard, J.K., M. Stephens, and P. Donnelly. 2000. Inference of population structure
11 using multilocus genotype data. *Genetics* 155: 945-959.
- 12 Purcell, S., B. Neale, K. Todd-Brown, L. Thomas, M.A.R. Ferreira, D. Bender, J. Maller,
13 P. Sklar, P.I.W. de Bakker, M.J. Daly et al. 2007. PLINK: a tool set for
14 whole-genome association and population-based linkage analyses. *Am J Hum*
15 *Genet* 81: 559-575.
- 16 Reed, T.E. 1969. Caucasian genes in American Negroes. *Science* 165: 762-768.
- 17 Reich, D., N. Patterson, P.L. De Jager, G.J. McDonald, A. Waliszewska, A. Tandon, R.R.
18 Lincoln, C. DeLoa, S.A. Fruhan, and P. Cabre. 2005. A whole-genome admixture
19 scan finds a candidate locus for multiple sclerosis susceptibility. *Nature Genet* 37:
20 1113-1118.
- 21 Sabeti, P.C. P. Varilly B. Fry J. Lohmueller E. Hostetter C. Cotsapas X. Xie E.H. Byrne
22 S.A. McCarroll R. Gaudet et al. 2007. Genome-wide detection and characterization
23 of positive selection in human populations. *Nature* 449: 913-918.
- 24 Sakamoto, H., K. Yoshimura, N. Saeki, H. Katai, T. Shimoda, Y. Matsuno, D. Saito, H.
25 Sugimura, F. Tanioka, S. Kato et al. 2008. Genetic variation in PSCA is associated
26 with susceptibility to diffuse-type gastric cancer. *Nat Genet* 40: 730-740.
- 27 Sankararaman, S., G. Kimmel, E. Halperin, and M.I. Jordan. 2008. On the inference of
28 ancestries in admixed populations. *Genome Res* 18: 668-675.
- 29 Sankararaman, S., S. Sridhar, G. Kimmel, and E. Halperin. 2008. Estimating local ancestry
30 in admixed populations. *Am J Hum Genet* 82: 290-303.
- 31 Smith, M.W. and S.J. O'Brien. 2005. Mapping by admixture linkage disequilibrium:
32 advances, limitations and guidelines. *Nat Rev Genet* 6: 623-632.
- 33 Smith, M.W., N. Patterson, J.A. Lautenberger, A.L. Truelove, G.J. McDonald, A.
34 Waliszewska, B.D. Kessing, M.J. Malasky, C. Scafe, and E. Le. 2004. A
35 high-density admixture map for disease gene discovery in african americans. *The*
36 *Am J Hum Genet* 74: 1001-1013.
- 37 Stannard, D. 1993. American Holocaust. *Oxford University Press*, Oxford, United
38 Kingdom.
- 39 Sundquist, A., E. Fratkin, C.B. Do, and S. Batzoglou. 2008. Effect of genetic divergence in
40 identifying ancestral origin using HAPAA. *Genome Res* 18: 676-682.
- 41 Tang, H., S. Choudhry, R. Mei, M. Morgan, W. Rodriguez-Cintron, E.G. Burchard, and
42 N.J. Risch. 2007. Recent genetic selection in the ancestral admixture of Puerto
43 Ricans. *Am J Hum Genet* 81: 626-633.
- 44 Tang, H., M. Coram, P. Wang, X. Zhu, and N. Risch. 2006. Reconstructing genetic
45 ancestry blocks in admixed individuals. *Am J Hum Genet* 79: 1-12.

- 1 Tang, H., J. Peng, P. Wang, and N.J. Risch. 2005. Estimation of individual admixture:
2 analytical and study design considerations. *Genetic epidemiology* 28: 289-301.
- 3 Thomas, H. 1999. The Slave Trade: The Story of the Atlantic Slave Trade. *Simon &*
4 *Schuster* 1440-1870.
- 5 Thorisson, G.A., A.V. Smith, L. Krishnan, and L.D. Stein. 2005. The international
6 HapMap project web site. *Genome Res* 15: 1592-1593.
- 7 Weir, B.S., L.R. Cardon, A.D. Anderson, D.M. Nielsen, and W.G. Hill. 2005. Measures of
8 human population structure show heterogeneity among genomic regions. *Genome*
9 *Res* 15: 1468-1476.
- 10 Weir, B.S. and C.C. Cockerham. 1984. Estimating F-statistics for the analysis of
11 population structure. *Evolution* 38: 1358-1370.
- 12 Winichagoon, P., S. Fucharoen, P. Chen, and P. Wasi. 2000. Genetic factors affecting
13 clinical severity in beta-thalassemia syndromes. *J Pediatr Hematol Oncol* 22:
14 573-580.
- 15 Wood, E.T., D.A. Stover, M. Slatkin, M.W. Nachman, and M.F. Hammer. 2005. The beta
16 -globin recombinational hotspot reduces the effects of strong selection around
17 HbC, a recently arisen mutation providing resistance to malaria. *Am J Hum Genet*
18 77: 637-642.
- 19 WORKMAN, P.L., B.S. BLUMBERG, and A.J. COOPER. 1963. SELECTION, GENE
20 MIGRATION AND POLYMORPHIC STABILITY IN A U. S. WHITE AND
21 NEGRO POPULATION. *Am J Hum Genet* 15: 429-437.
- 22 Wu, X., Y. Ye, L.A. Kiemeny, P. Sulem, T. Rafnar, G. Matullo, D. Seminara, T. Yoshida,
23 N. Saeki, A.S. Andrew et al. 2009. Genetic variation in the prostate stem cell
24 antigen gene PSCA confers susceptibility to urinary bladder cancer. *Nat Genet* 41:
25 991-995.
- 26 Yeager, M., N. Chatterjee, J. Ciampa, K.B. Jacobs, J. Gonzalez-Bosquet, R.B. Hayes, P.
27 Kraft, S. Wacholder, N. Orr, S. Berndt et al. 2009. Identification of a new prostate
28 cancer susceptibility locus on chromosome 8q24. *Nat Genet* 41: 1055-1057.
- 29 Yen, C.C., W.M. Chen, T.H. Chen, W.Y. Chen, P.C. Chen, H.J. Chiou, G.Y. Hung, H.T.
30 Wu, C.J. Wei, and C.Y. Shiau. 2009. Identification of chromosomal aberrations associated
31 with disease progression and a novel 3q13. 31 deletion involving LSAMP gene in
32 osteosarcoma. *Int J Oncol* 35: 775-788.
- 33 Zakharia, F., A. Basu, D. Absher, T.L. Assimes, A.S. Go, M.A. Hlatky, C. Iribarren, J.W.
34 Knowles, J. Li, B. Narasimhan et al. 2009. Characterizing the admixed African
35 ancestry of African Americans. *Genome Biol* 10: R141.
- 36
37
38

1 **Tables**2 **Table 1. Regions Showing Excess of European or African Ancestry.**

Regions	Position	Excess ancestry	Size (bp)	SNPs	Highest deviation	Genes	Pathways	Related diseases
1p36	chr1:17409539..21604321	African	4194783	489	0.0253	<i>AKR7A2*</i> , <i>IGSF21</i> , <i>DDOST*</i> , <i>HTR6 et al.</i>	Diabetes pathways, signaling by GPCR, metabolism of amino acids	Diabetes, pancreatic cancer
2q37	chr2:241750403..242568618	African	818216	16	0.0231	<i>SEPT2*</i> , <i>HDLBP*</i> , <i>PDCD*1</i> , <i>FARP2 et al.</i>	Signaling in immune system, Axon guidance, metabolism of nucleotides	Bladder cancer, lung cancer, coronary atherosclerosis
2p22	chr2:37451925..37508581	European	56657	9	0.0230	<i>QPCT</i> , (<i>EIF2AK2*</i> 222kb)	Influenza infection	Influenza infection
3q13	chr3:116930811..118313302	European	1382492	216	0.0253	<i>LSAMP*</i>	Homophilic adhesion	Osteosarcoma
6q26	chr6:163653158..163653428	European	271	2	0.0225	<i>PACRG*</i>	Mediate proteasomal degradation	Juvenile Parkinson's disease
16q21	chr16:61214438..61242497	European	28060	9	0.0229	NA	NA	NA

3

4 NA: not available. * Denotes genes associated with diseases. Genes in parentheses are strong candidates

5 out of the chromosome location but closest.

6

1 **Table 2. Regions with highly differentiated allele frequency between AAF and YRI**
 2 ($F_{ST} > 0.0452$).

Regions or SNPs	Position	Size (bp)	SNPs	Highest F_{ST}	Genes	Pathways	Related disease
1p21	chr1:100125058..100183875	58817	2	0.0562	<i>AGL*</i>	Metabolism of carbohydrate	Glycogen storage disease
1q22	chr1:153401959..153464086	62127	4	0.0692	<i>THBS3*</i> , <i>MUC1*</i> , <i>MTX1</i> , <i>TRIM46</i> , <i>KRTCAP2</i>	Signaling by PDGF	Stomach cancer, breast cancer, osteosarcoma
rs12094201	chr1:236509336	1	1	0.0561	(<i>ZP4*</i> 389kb)	NA	Hypertension, Non-alcoholic fatty liver
rs7642575	chr3: 31400165	1	1	0.0453	(<i>STT3B</i> , <i>OSBPL10*</i> 149 kb)	NA	Peripheral arterial disease
6p21-p22	chr6:26554684..33961049	7406365	11	0.0711	<i>HLA-B*</i> , <i>HLA-C</i> , <i>EHMT2*</i> , <i>HLA-DPA1*</i> , <i>HLA-DRB5</i> , <i>EHM</i> , <i>BTN3A3</i> , <i>et al</i>	Signaling by GPCG, signaling in immune system, HIV infection, Diabetes pathway	HIV, Crohn's disease, rheumatoid arthritis, juvenile idiopathic arthritis, colorectal cancer, systemic sclerosis
6q25	chr6:151555551..151569258	13707	2	0.0545	(<i>AKAP12*</i> 40kb)	Cell growth	Hypertension, hemorrhagic stroke
rs10499542	chr7: 22235870		1	0.04606	<i>RAPGEF5*</i>	GTP/GDP-regulation	Thyroid stimulating hormone
7q21	chr7:79768487..80482597	714110	10	0.0946	<i>CD36*</i> , <i>SEMA3C</i>	Metabolism of lipids and lipoprotein	Metabolic syndrome, malaria
8q24	chr8:143754039..143758933	4894	2	0.04679	<i>PSCA*</i>	NA	Prostate cancer, bladder cancer, gastric cancer
11p15	chr11:5034229..5421456	387227	3	0.0617	<i>HBB*</i> , <i>HBD*</i> , <i>HBE1*</i> , <i>HBG2</i> , <i>OR5111</i> , <i>et al</i>	Signaling by GPCR	Sickle cell disease, beta-thalassemia, malaria
rs4883422	chr12:7189594	1	1	0.04721	<i>CLSTN3</i>	NA	NA
rs6491096	chr13:25488362	1	1	0.04716	<i>ATP8A2</i>	NA	NA
rs1075875	chr16: 47595721	1	1	0.0766	(<i>CBLN1</i> 277kb)	NA	NA
rs6015945	chr20:59319574	1	1	0.0627	<i>CDH4*</i>	Cell junction organization	Alzheimer's Disease

3

4 NA: not available. * Denotes the genes associated with diseases. Genes in parentheses are strong
 5 candidates out of the chromosome location but closest.

6

7

8

9

1 **FIGURE LEGENDS**

2

3 **Figure 1. Genome-wide Distribution of European Ancestral Contributions.** Mean
4 European ancestral contribution across 1,890 African American individuals at each SNP.
5 Green line is the estimated genome-wide mean European ancestral contribution (21.68%).
6 Blue bands indicate +2 and -2 SDs from the mean ancestral contribution and red Bands
7 indicate +3 and -3 SDs from the mean ancestral contribution.

8

9 **Figure 2. Genomic Distribution of F_{ST} between AAF and YRI (A) and Genomic**
10 **Distribution of F_{ST} between African American and rAfA (B).** The dashed red
11 horizontal line indicates the cutoff threshold (99.99th percentile). Locus-specific F_{ST}
12 between YRI and CEU were calculated when MAF >0.05 in both populations. The rAfA
13 was constituted according to the ancestry proportion of CEU and YRI under neutrality.

14

15 **Figure 3. Enrichment of high F_{ST} loci for different SNP categories.** Observed excess of
16 high F_{ST} loci in different SNP classes, with respect to non-genic class, among high F_{ST} bin
17 (99th percentile; F_{ST} >0.0164). The values on the bar are p-values of χ^2 tests. “NS” stands
18 for “not significant”.

19

20

21

Proportion of European Ancestry

0.24
0.23
0.22
0.21
0.20
0.19

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

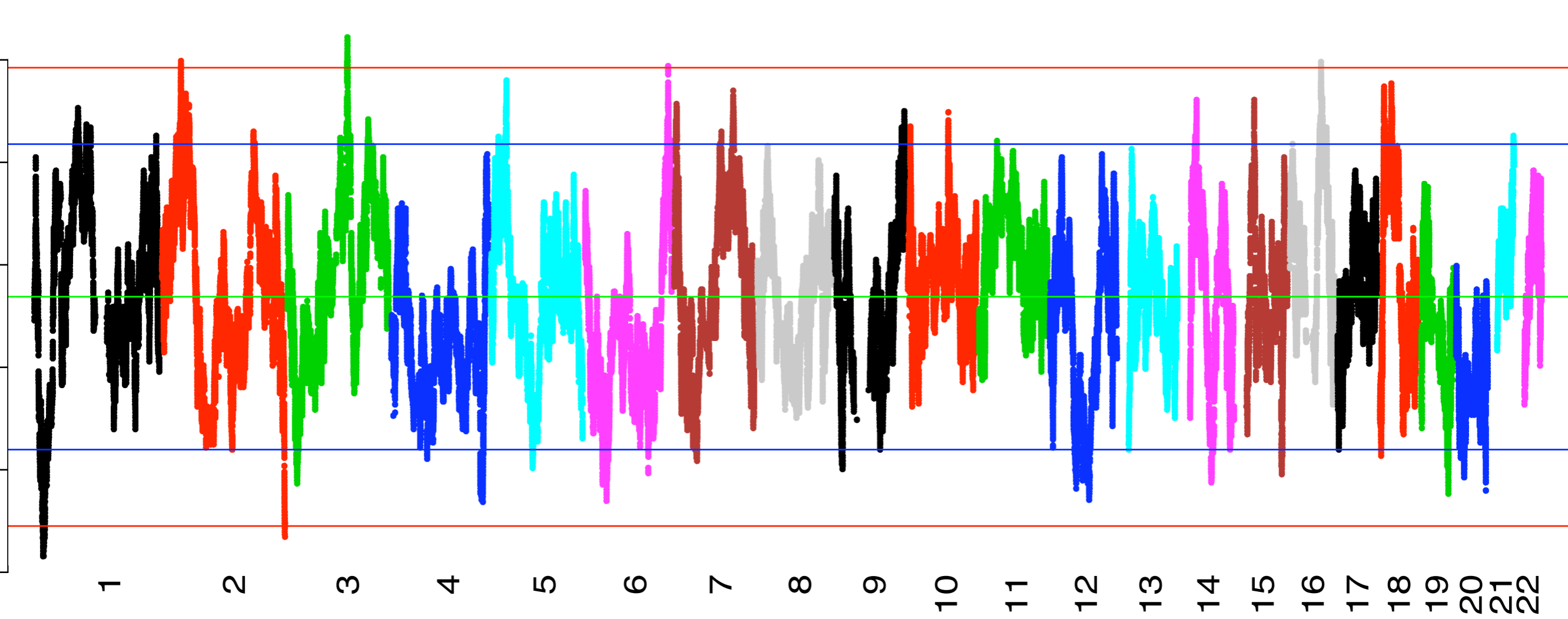
19

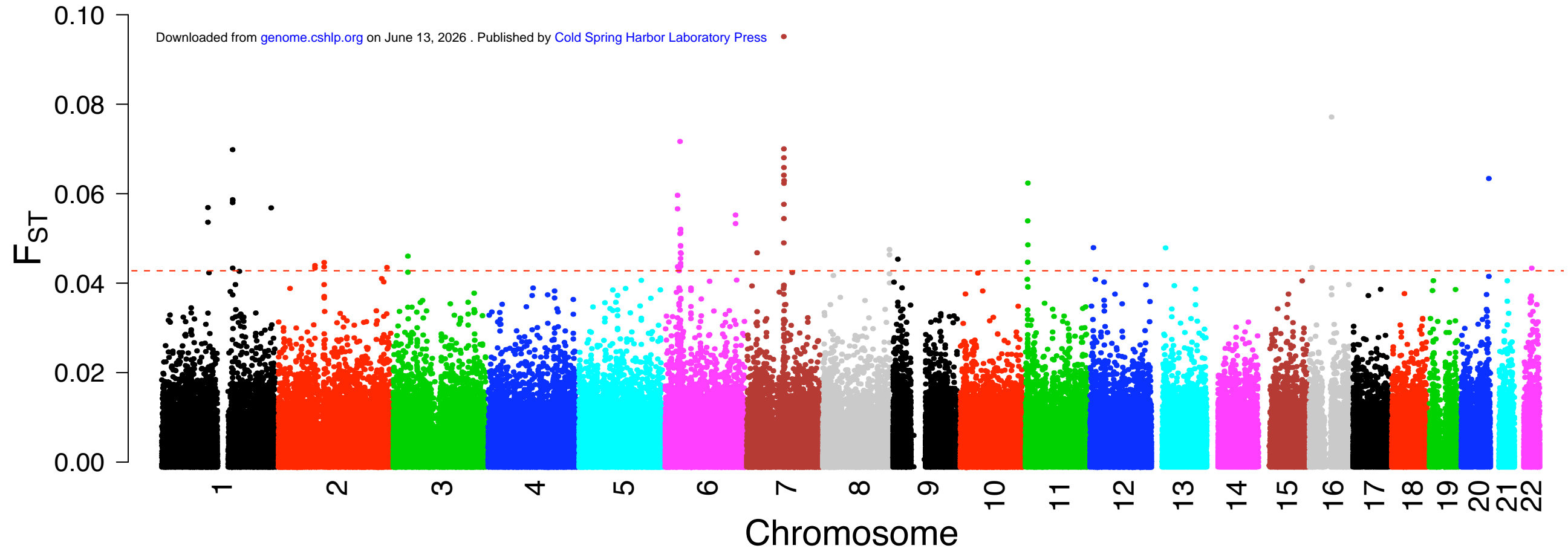
20

21

22

Chromosome



A**B**