

1 **Supplemental Text**

2 *Additional Information on populations and samples*

3 Overall 1,890 unrelated AfA samples with ignorable ancestry outside of Africa and
4 Europe were studied, from which 1,838 individual were download from iControlDB, and
5 52 African Americans were ASW (individuals of African ancestry from the southwest
6 U.S.) from the International HapMap Project ([Altshuler et al. 2010](#)). 113 YRI (Yoruba in
7 Ibadan, Nigeria) from the International HapMap Project and 102 HGDP indigenous
8 Africans from 7 different geographic or linguistic groups (including 21 Yoruba from
9 West Central Africa, 22 Mandenka from West Africa, 11 Bantu speakers from Kenya, 8
10 Bantu speakers from Southwestern Africa, 22 Biaka Pygmy and 13 Mbuti Pygmy from
11 Central Africa, and 5 San from Southern African) ([Li et al. 2008](#)) represent the potential
12 African ancestry of AfA. 113 CEU (Utah residents with northern and western European
13 ancestry from the CEPH collection) from the International HapMap Project and 156
14 HGDP indigenous Europeans from 8 different geographic and linguistic groups
15 (including 17 Adygei from Russia, 24 Basque from France, 28 French in France, 12
16 Italian from North Italy, 15 Orcadian from Orkney Islands, 25 Russian in European, 28
17 Sardinian in Sardinia island and 7 Tuscan from Italy) ([Li et al. 2008](#)) represent the
18 potential European ancestry of AfA. Overall 2,648 CAU-GWAS also represent the
19 European ancestry of AfA. The 84 CHB (Han Chinese from Beijing) and 85 CHD (Han
20 Chinese from metropolitan Denver, CO, USA) represent populations in East Asia. In
21 addition, 24 Amerindian samples showing no genetic admixture with European and
22 African from HGDP represent Native American.

23

24 *Correcting for Linkage disequilibrium and PCA analysis*

25 Principle component analysis (PCA) is a statistics method widely used to identify
26 population genetic variation across geographical location and ethnic background. It is
27 well known that the uncorrected linkage disequilibrium (LD) would seriously distort the
28 results of PCA and thus bias the population structure ([Patterson et al. 2006](#)). In order to
29 reduce the LD between markers, we used the parameter "--indep-pairwise 50 5 0.5" in
30 PLINK ([Purcell et al. 2007](#)) to remove any SNP that had a pairwise r^2 greater than 0.5 in

1 window of 50 SNPs, shifted and recalculated every 5 SNPs. The original 491,526
2 autosomal SNPs were reduced to 341,672 SNPs after this process.

3 Based on these thinned markers, PCA was performed using *smartpca* in the
4 EIGENSOFT package (Patterson et al. 2006) under default setting with no outlier
5 removal. In the PCA analysis, the first principal component (PC1) is the coordinate
6 drawn in the multidimensional space so that the projections of the points (each point
7 represents an individual) to the coordinate have the largest variance. The second principal
8 component (PC2) is the coordinate drawn in the multidimensional space so that the
9 projections of the points to the coordinate have the second largest variance, and so forth.
10 Intuitively, PC1 and PC2 could reflect the two primary genetic differences between
11 samples. In this study, only PC1 and PC2 were plotted in our PCA plot and have been
12 carefully analyzed.

13 The projections of individuals onto existing axes not only identified the population
14 structure, but also could identify admixture proportion of each admixed individual
15 (McVean 2009). So we performed PCA on AfA and its two ancestral parental
16 populations (YRI and CEU), and inferred the ancestral proportion of each AfA based on
17 projection of samples on PC1.

18

19 *YRI and CEU are better choices for ancestral parental populations than the others*

20 Using the 341,672 not high-linked SNPs, the 1,890 AfA samples were studied
21 together with 3320 samples from 19 putative ancestral parental populations. In the PCA
22 plot, populations from the same continents were clustered together, which reflected
23 genetic diversity of populations from the same continent is rather lower compared with
24 that of inter-continent populations. AfA lie on the direct line between European (except
25 Russian and Adygei, which may have some components from Asia) and African, which
26 reflected their varying levels of admixture in AfA between the two ancestral groups
27 (Figure S2A). When samples from Asia and America (as well as Russian and Adygei)
28 were removed, AfA almost distributed only between West African and European, which
29 reflected the West African, especially YRI, contributing to most of the African
30 components of ancestry (Figure S2B). Substructure in European populations showed up
31 when only YRI, AfA and European populations were analyzed. All European populations

1 except Sardinian contribute to European components of AfA. However, CEU contributed
2 to most of the European components of AfA as only a few samples were beyond the line
3 between YRI-CEU (Figure S2C). Then, we performed PCA on samples from YRI, CEU,
4 and AfA. The PC1 explained 80% of the total variance of top ten PCs while PC2
5 explained only 2.4%, which reflect that most of the variances were caused by the genetic
6 differentiation between West African and European ancestry, especially like YRI and
7 CEU (Figure S2D).

9 *Population structure analysis using FRAPPE and STRUCUTE*

10 A number of methods and software have been developed to infer the proportional
11 ancestry of each individual using genotype data without prior knowledge of individual
12 ancestral and geographic information. FRAPPE (Tang et al. 2005), which implements an
13 expectation-maximization (EM) algorithm, was used to infer the admixture proportion of
14 AfA in this study. We used FRAPPE not just because of its fastness, but also because we
15 could take full advantage of the extremely high-density SNPs data since it was almost
16 unaffected by Linkage disequilibrium (LD). The program was allowed to run for 10,000
17 iteration, with prior-specified cluster numbers. First, we run FRAPPE on the 1,890 AfA
18 and all putative ancestral parental populations (7 HGDP African populations, 8 HGDP
19 European populations, CEU and YRI) by taking K=2. Second, we only used CEU and
20 YRI as the ancestral populations of AfA, and ran FRAPPE on AfA, CEU and YRI by
21 setting K=2. The two independent runs yield almost identical results (Figure S3, S4).

22 STRUCTURE (Pritchard et al. 2000; Falush et al. 2003), implementing a model-
23 based clustering method, was widely used to infer population structure. However, the
24 model in the program is best suited to these markers not highly linked. Therefore, any
25 contiguous SNPs <1M were removed and only a set of 2,654 SNPs was kept for further
26 analysis. We used the admixture-model and assumed that allele frequencies were
27 correlated. The program was run with 100,000 iterations and 100,000 burn-ins.

28 *Comparisons of methods for inferring locus-specific ancestry*

29 Various methods have been developed to infer locus-specific ancestry based on
30 high density SNPs data, such as ANCESTRYMAP (Patterson et al. 2004), SABER (Tang

1 et al. 2006), LAMP and LAMP-ANC (Sankararaman et al. 2008b), uSWITCH and
2 uSWITCH-ANC (Sankararaman et al. 2008a), HAPAA (Sundquist et al. 2008) and
3 HAPMIX (Price et al. 2009). There is no doubt that these methods outperform older
4 methods that only analyzed a few unlinked markers. However, It is still difficult to
5 distinguish which method performs better than the other since the data used in different
6 studies may vary greatly, and sometimes maybe completely different. In particular, we
7 are interested in which method performs the best on our dataset.

8 First, the genotypes of the two parental populations were labeled in our simulations
9 so that the ancestral status of the simulated AfA was known. Second, by comparing the
10 inferred ancestral status of AfA with the real ancestral status, the accuracies of five
11 methods (SABER, LAMP-ANC, LAMP, HAPAA, HAPMIX) were calculated. The
12 accuracy was simply calculated by using the number of correct inferences divided the
13 total number or comparisons. The inferred accuracy of each method was: HAPMIX
14 (98.1%), LAMP-ANC (96.5%), LAMP (96.3%), HAPAA (95.7%), SABER (92.3%). It is
15 obvious that the performance of HAPMIX is much better than the other methods in this
16 case. These observations were essentially similar to those of previous studies
17 (Sankararaman et al. 2008a; Sundquist et al. 2008; Price et al. 2009). HAPMIX, which
18 applies an extension of haplotype-based method (Price et al. 2009), was used to infer the
19 locus specific ancestry in this study. However, This does not mean that HAPMIX
20 performs better in other situations (e.g., the performance of HAPMIX was not so good if
21 high-qualified haplotypes of ancestral parental populations were not available; data not
22 shown). Moreover, many other situations, such as multiply-wave admixture and closely
23 related parental populations, were not simulated in this study.

24

25 *Generation since admixture for each individual*

26 By running HAPMIX in diploid mode, we also detected chromosomal segments of
27 distinct ancestry in each admixed individual. Then, we computed the generation since
28 admixture for each individual (g) following Price et al. (Price et al. 2009):

29
$$g = \frac{N}{4u(1-u) \times 35.3}$$

1 where N is the number of ancestry transition observed, and μ is the ancestry proportion of
2 the admixed individual. A fraction of $2u(1-u)$ of all recombination events was
3 expected, and 35.3 Morgan was the overall genetic distance used in our genetic map,
4 which was adapted from HapMap.

5

6 *Locus-specific ancestry inferred by LAMP-ANC*

7 To investigate the repeatability of selection signals detected by HAPMIX, we also
8 inferred the locus-specific ancestry using LAMP-ANC, which also performs very well.
9 The input genotypic data for LAMP-ANC were the same as those used by HAPMIX
10 except their ancestral allele frequencies were provided directly. The generation since
11 admixture was set to 7 (hybrid isolation model) based on HAPMIX results. Genome-wide
12 distribution of European ancestral contribution in AfA based on CAU-GWAS and YRI
13 was calculated (Figure S6), and the results are essentially the same when other putative
14 ancestral populations were used. The locus-specific European ancestry proportion across
15 the genome of AfA estimated by LAMP-ANC was $21.69\% \pm 0.82\%$ (mean \pm SD). The
16 standard deviation (SD) estimated by LAMP-ANC was higher than that based on
17 HAPMIX, which might be caused by higher statistic noise as the result of higher error
18 rate of LAMP-ANC. We found that all the ancestry deviation regions can be replicated
19 using LAMP-ANC to some extent.

20

21 *The results are unlikely to be affected by ascertainment bias*

22 The data of this study displayed a deficit of rare variants and an excess of
23 intermediate frequency SNPs. This frequency bias is likely to affect statistics based on
24 allele frequencies including F_{ST} (Clark et al. 2005). However, our goal was not to define
25 precise levels of population differentiation, but to detect differentiations of F_{ST} among
26 different SNP categories. In this context, as the ascertainment of our SNPs was agnostic
27 to any functional annotation, all SNP classes should be equally influenced by

1 ascertainment bias. These biases should therefore have a minimal effect on our analysis
2 and conclusions. In order to clarify whether ascertained bias has an effect on the
3 enrichment analysis of functional SNPs, we randomly chose several sets of SNPs (1%
4 and 5%), and we found that essentially no functional SNPs classes were significantly
5 enriched in the selected SNPs bin.

6

7 **References**

- 8 Altshuler, DM, Gibbs, RA, Peltonen, L, Altshuler, DM, Gibbs, RA, Peltonen, L,
9 Dermitzakis, E, Schaffner, SF, Yu, F, Peltonen, L et al. 2010. Integrating common
10 and rare genetic variation in diverse human populations. *Nature* **467**: 52-58.
- 11 Clark, AG, Hubisz, MJ, Bustamante, CD, Williamson, SH, Nielsen, R. 2005.
12 Ascertainment bias in studies of human genome-wide polymorphism. *Genome*
13 *research* **15**: 1496-1502.
- 14 Falush, D, Stephens, M, Pritchard, JK. 2003. Inference of population structure using
15 multilocus genotype data: Linked loci and correlated allele frequencies. *Genetics*
16 **164**: 1567-1587.
- 17 Li, JZ, Absher, DM, Tang, H, Southwick, AM, Casto, AM, Ramachandran, S, Cann, HM,
18 Barsh, GS, Feldman, M, Cavalli-Sforza, LL et al. 2008. Worldwide human
19 relationships inferred from genome-wide patterns of variation. *Science (New*
20 *York, NY* **319**: 1100-1104.
- 21 Long, JC. 1991. The genetic structure of admixed populations. *Genetics* **127**: 417-428.
- 22 McVean, G. 2009. A genealogical interpretation of principal components analysis. *PLoS*
23 *genetics* **5**.
- 24 Patterson, N, Hattangadi, N, Lane, B, Lohmueller, KE, Hafler, DA, Oksenberg, JR,
25 Hauser, SL, Smith, MW, O'Brien, SJ, Altshuler, D et al. 2004. Methods for high-
26 density admixture mapping of disease genes. *Am J Hum Genet* **74**: 979-1000.
- 27 Patterson, N, Price, AL, Reich, D. 2006. Population structure and eigenanalysis. *PLoS*
28 *genetics* **2**: e190.
- 29 Pfaff, CL, Parra, EJ, Bonilla, C, Hiester, K, McKeigue, PM, Kamboh, MI, Hutchinson,
30 RG, Ferrell, RE, Boerwinkle, E, Shriver, MD. 2001. Population structure in
31 admixed populations: Effect of admixture dynamics on the pattern of linkage
32 disequilibrium. *American journal of human genetics* **68**: 198-207.
- 33 Price, AL, Tandon, A, Patterson, N, Barnes, KC, Rafaels, N, Ruczinski, I, Beaty, TH,
34 Mathias, R, Reich, D, Myers, S. 2009. Sensitive detection of chromosomal
35 segments of distinct ancestry in admixed populations. *PLoS Genet* **5**: e1000519.
- 36 Pritchard, JK, Stephens, M, Donnelly, P. 2000. Inference of population structure using
37 multilocus genotype data. *Genetics* **155**: 945-959.

1 Purcell, S, Neale, B, Todd-Brown, K, Thomas, L, Ferreira, MAR, Bender, D, Maller, J,
2 Sklar, P, de Bakker, PIW, Daly, MJ et al. 2007. Plink: A tool set for whole-
3 genome association and population-based linkage analyses. *American journal of*
4 *human genetics* **81**: 559-575.

5 Sankararaman, S, Kimmel, G, Halperin, E, Jordan, MI. 2008a. On the inference of
6 ancestries in admixed populations. *Genome Res* **18**: 668-675.

7 Sankararaman, S, Sridhar, S, Kimmel, G, Halperin, E. 2008b. Estimating local ancestry
8 in admixed populations. *Am J Hum Genet* **82**: 290-303.

9 Sundquist, A, Fratkin, E, Do, CB, Batzoglou, S. 2008. Effect of genetic divergence in
10 identifying ancestral origin using hapaa. *Genome research* **18**: 676-682.

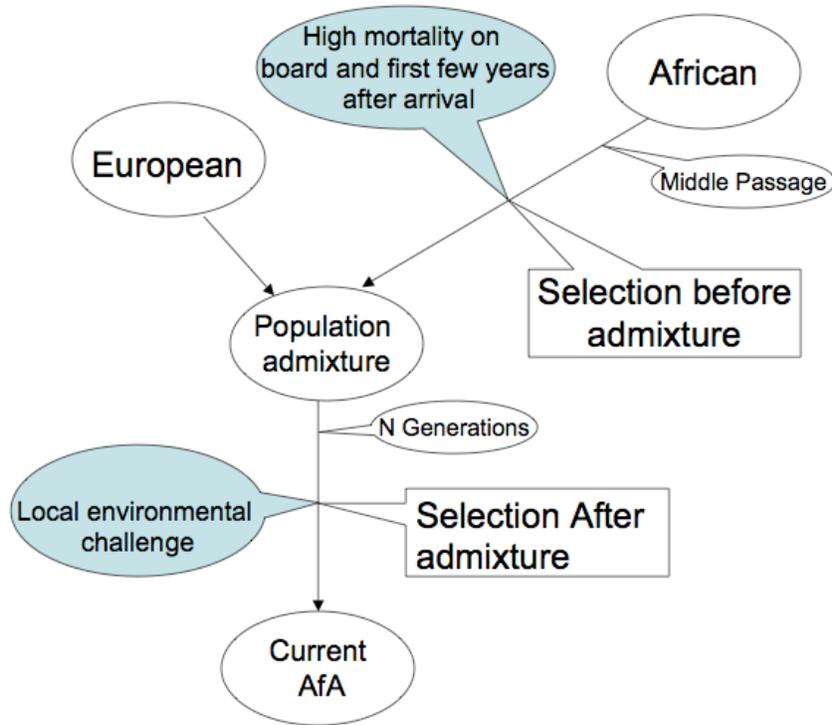
11 Tang, H, Coram, M, Wang, P, Zhu, X, Risch, N. 2006. Reconstructing genetic ancestry
12 blocks in admixed individuals. *American journal of human genetics* **79**: 1-12.

13 Tang, H, Peng, J, Wang, P, Risch, NJ. 2005. Estimation of individual admixture:
14 Analytical and study design considerations. *Genetic epidemiology* **28**: 289-301.
15
16
17
18
19
20
21
22

1 **Figures**

2 **Figure S1.** Schematic of the two natural selection events (pre-and post-admixture) in
3 African Americans.

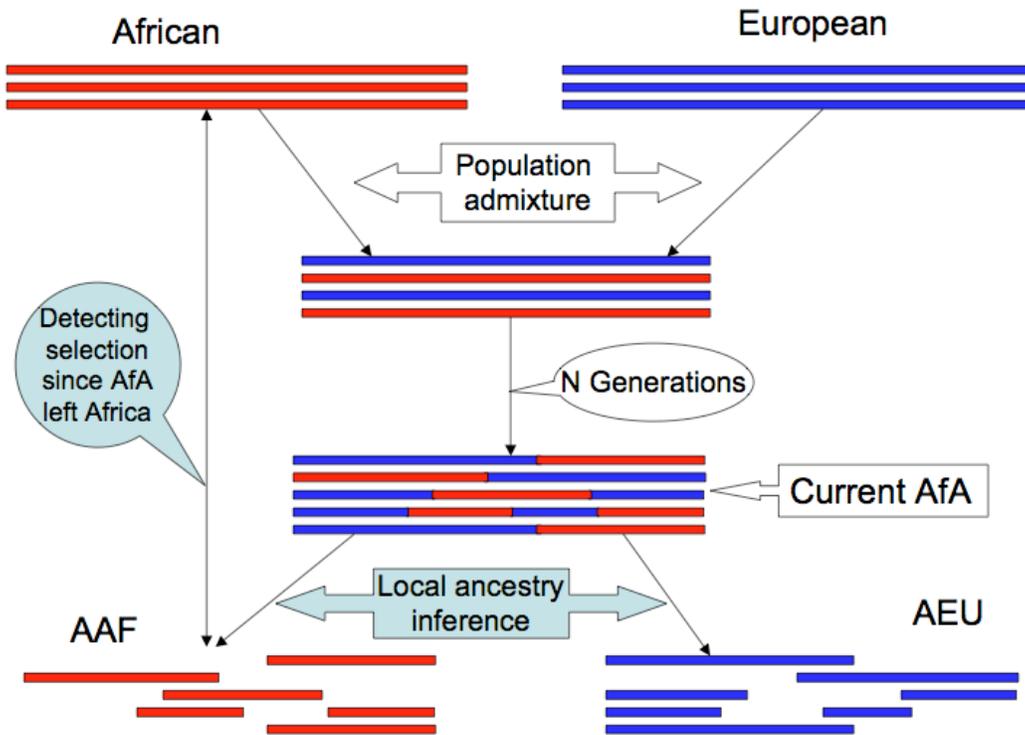
4



5

6

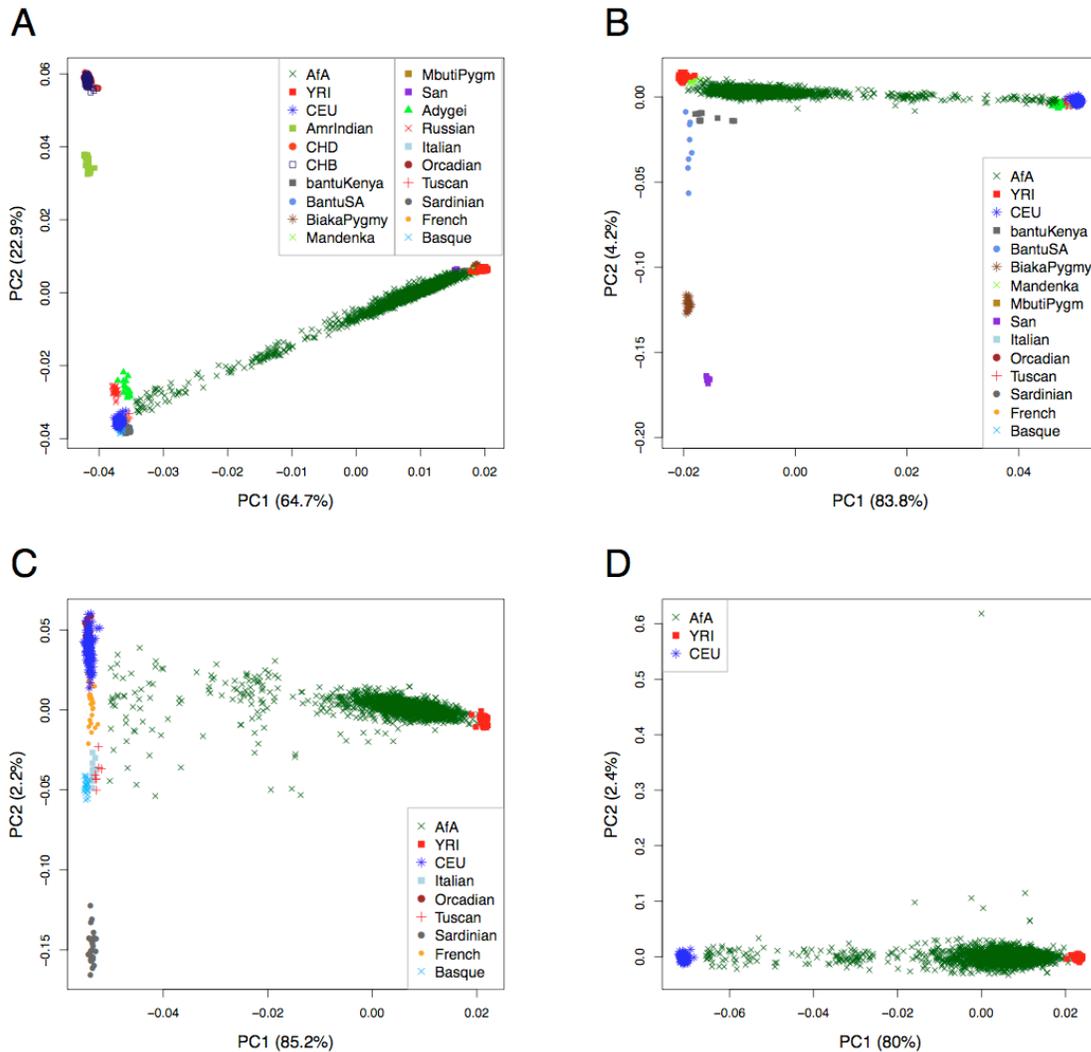
- 1 **Figure S2.** Schematic of the strategy to detect natural selection by comparing the
- 2 African components of ancestry in AfA with African populations.



- 3
- 4

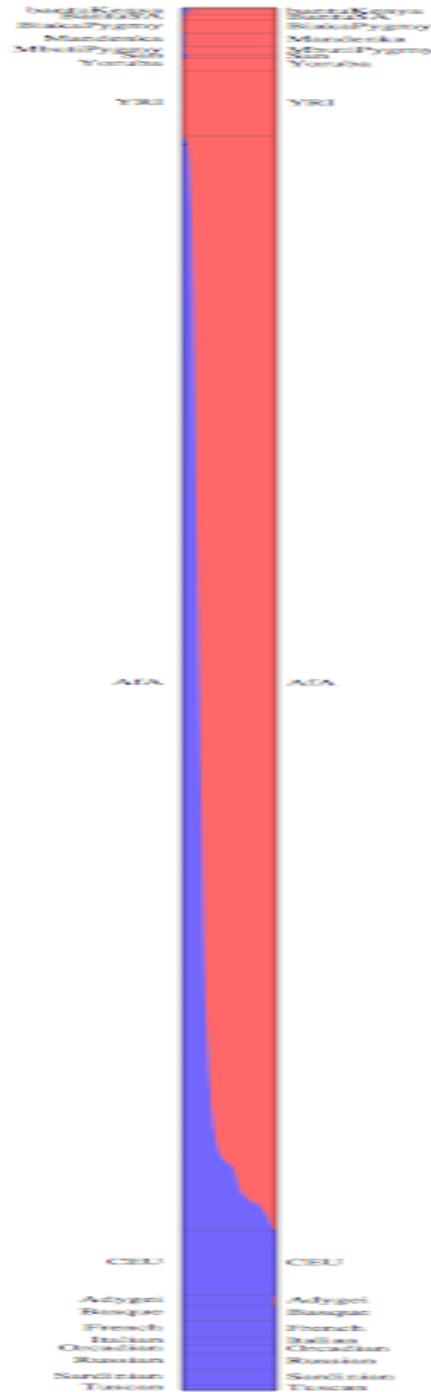
1 **Figure S3. Analyses of the first two principal components.** The % Eigenvalue is the
 2 percentage of the total variance in the Top ten PCs. (A) 1,890 AfA and 3320 samples
 3 from 19 other populations represent global-wide populations. (B) 1,890 AfA and all
 4 possible European or African ancestral populations. (C) 1,890 AfA, YRI and all possible
 5 European ancestral populations. (D) 1,890 AfA and two putative parental populations
 6 (YRI and CEU).

7
 8



9

- 1 **Figure S4. FRAPPE analysis of 8 African populations, 9 European populations and**
- 2 **AfA when K=2.** Each individual is represented by a vertical line, which is partitioned
- 3 into two segments corresponding to the inferred membership of the genetic clusters
- 4 indicated by the colors.

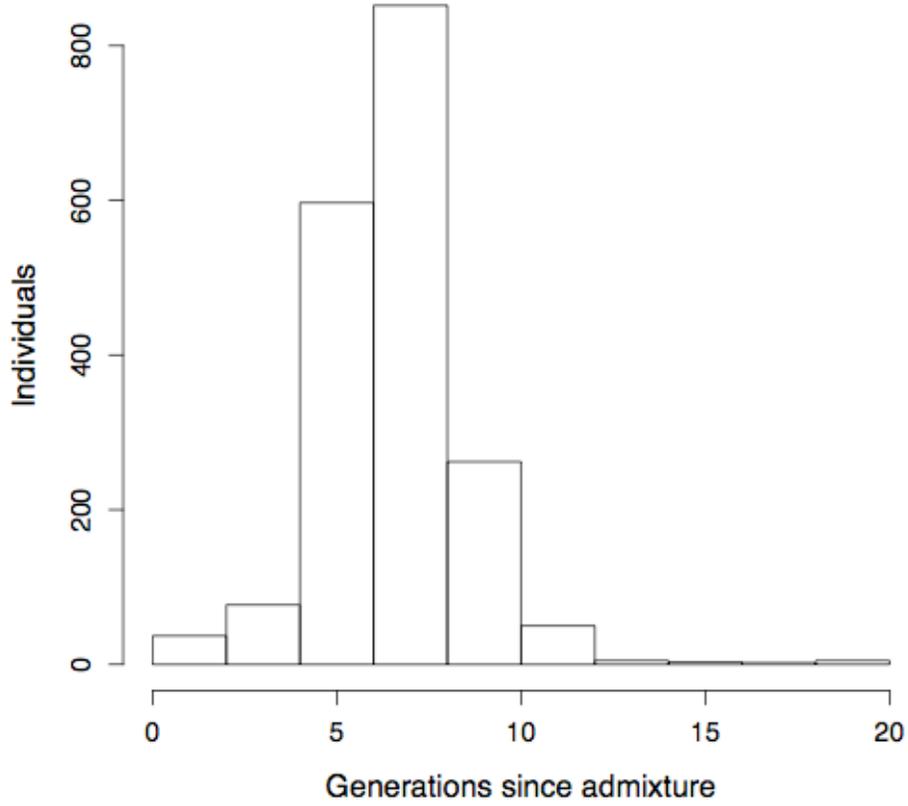


- 1 **Figure S5. FRAPPE analysis of YRI, AfA and CEU when K=2.** Each individual is
- 2 represented by a vertical line, which is partitioned into two segments corresponding to the
- 3 inferred membership of the genetic clusters indicated by the colors.



1

2 **Figure S6. Histogram of generations since admixture for the 1,890 African**
3 **Americans.**

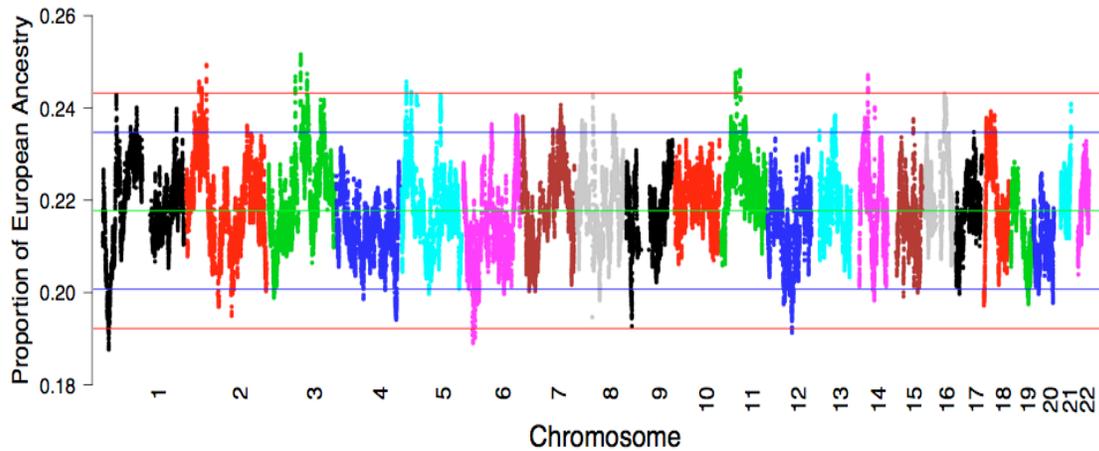


4

5

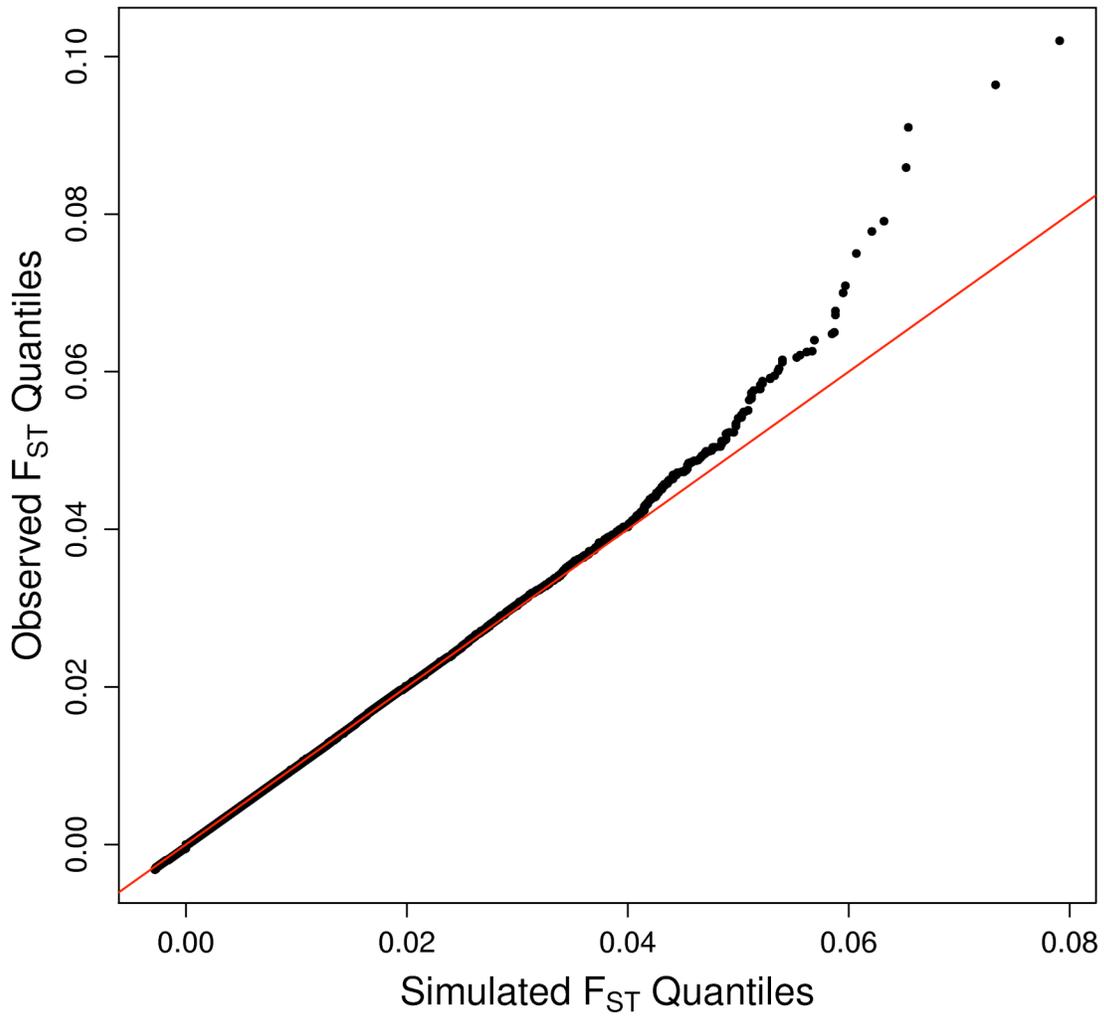
6

1 **Figure S7. Genome-wide distribution of European ancestral contributions estimated**
2 **by LAMP-ANC. Mean European ancestral contributions across 1,890 African American**
3 **individuals at each SNP. Green line is the estimated genome-wide mean European**
4 **ancestral contributions (21.69%). Blue bands indicate +2 and -2 SDs from the mean**
5 **ancestral contributions and red Bands indicate +3 and -3 SDs from the mean ancestral**
6 **contribution.**



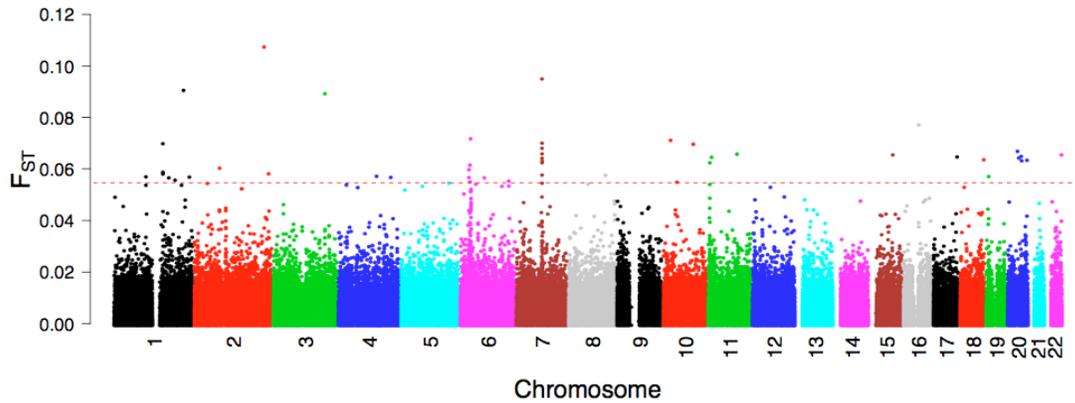
7
8

1 **Figure S8. Q-Q plot of locus-specific F_{ST} between AAF and YRI.** This Q-Q plot
2 compares empirical values on the vertical axis to simulated values on horizontal axis. Red
3 line shows null distribution that F_{ST} of empirical data are the same as that of simulated.
4



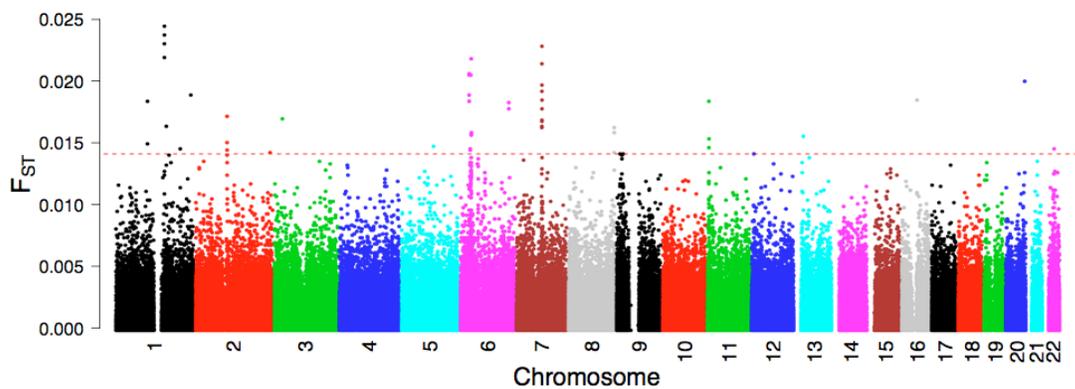
5
6
7

1 **Figure S9. Genome-wide distribution of F_{ST} between AAF and YRI for all**
2 **autosomal SNPs.** The dashed red horizontal line indicates the cutoff threshold (99.99th
3 percentile).



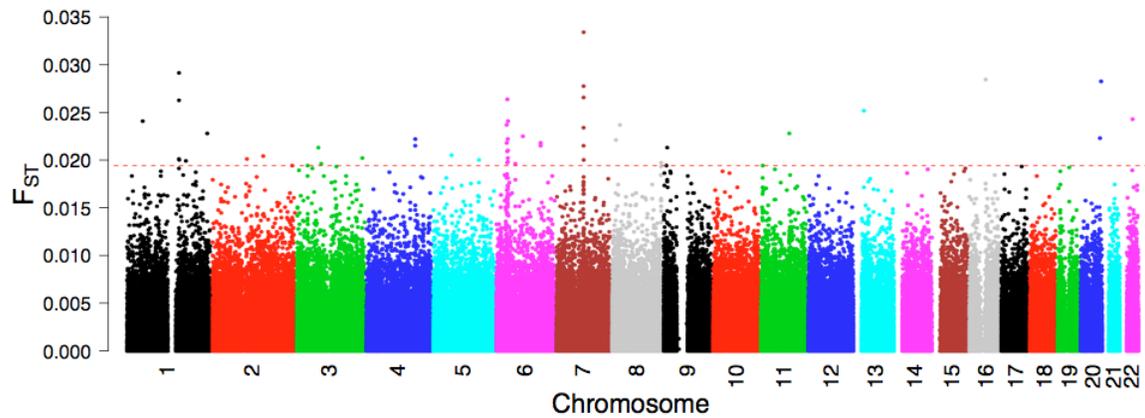
4
5
6
7

8 **Figure S10. Genome-wide distribution of F_{ST} between AfA and rAfA.** rAfA is
9 construct by CAU-GWAS and YRI according to estimated admixture proportion. The
10 dashed red horizontal line indicates the cutoff threshold (99.99th percentile).



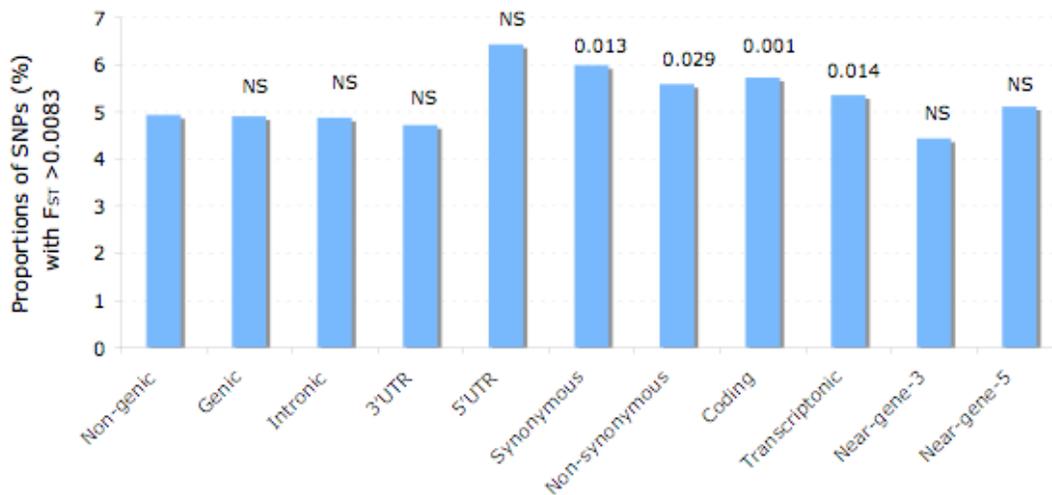
11
12
13
14

1 **Figure S11. Genome-wide distribution of F_{ST} between AfA and rAfA.** rAfA is
 2 construct by CAU-GWAS and APP (African parental population of AfA constructed by
 3 YRI, Mandenka and Bantu) according to estimated admixture proportion. The dashed red
 4 horizontal line indicates the cutoff threshold (99.99th percentile).



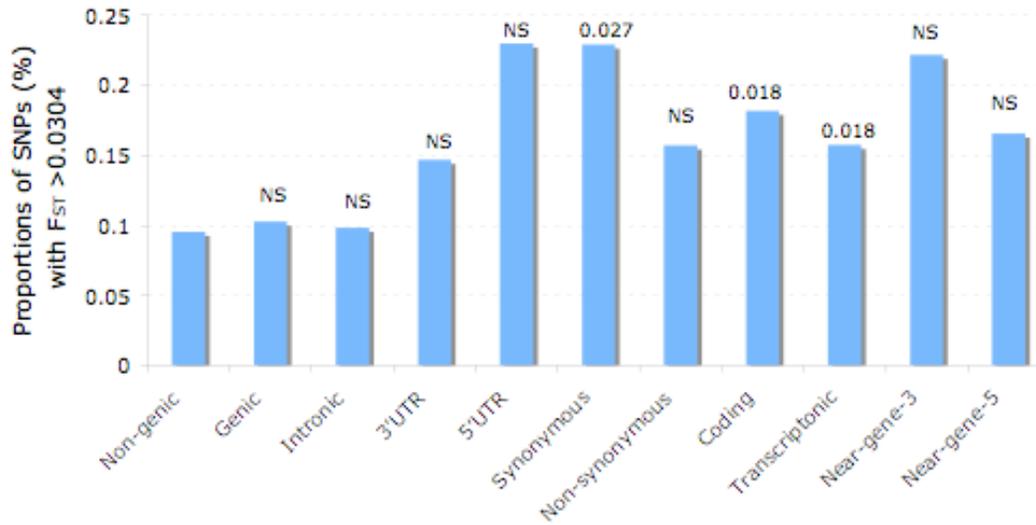
5
 6
 7

8 **Figure S12. Enrichment of high F_{ST} loci for different SNP categories (95th**
 9 **percentile; $F_{ST} > 0.0083$).** Observed excess of high F_{ST} loci for different SNP classes,
 10 with respect to nongenic class, among high F_{ST} bin (95th percentile; $F_{ST} > 0.0083$). The
 11 values on the bar are p-values of χ^2 tests. “NS” stands for “not significant”.



12

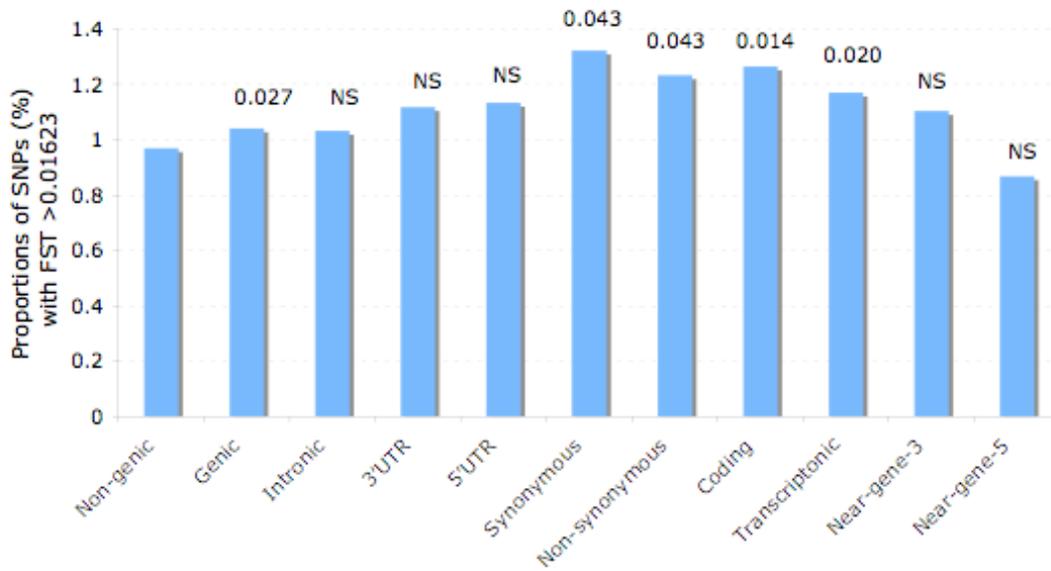
1 **Figure S13. Enrichment of high F_{ST} loci for different SNP categories (99.9th**
2 **percentile; $F_{ST} > 0.0304$).** Observed excess of high F_{ST} loci for different SNP classes,
3 with respect to nongenic class, among high F_{ST} bin (99.9th percentile; $F_{ST} > 0.0304$). The
4 values on the bar are p-values of χ^2 tests. “NS” stands for “not significant”.
5



6
7

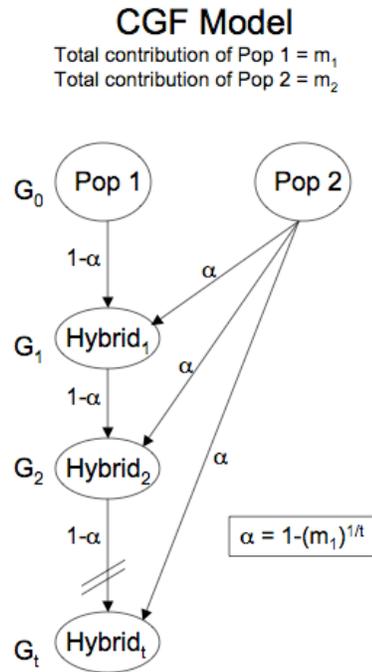
1 **Figure S14. Enrichment of high F_{ST} loci for different SNP categories (MAF>0.05;**
2 **99th percentile; F_{ST} >0.01623).** Observed excess of high F_{ST} loci for different SNP
3 classes, with respect to nongenic class, among high F_{ST} bin (MAF>0.05; 99th percentile;
4 F_{ST} >0.01623). The values on the bar are p-values of χ^2 tests. “NS” stands for “not
5 significant”.

6
7



8
9

- 1 **Figure S15.** Schematic of continuous-gene-flow (CGF) model adapted from Long
- 2 (Long 1991) and Pfaff *et al.* (Pfaff *et al.* 2001).
- 3



- 4
- 5
- 6

1 **Tables**

2

3 **Table S1. SNPs showing strong linkage with rs3211938(G) in CD36.**

dbSNP_id	Alleles	Allele linked with rs3211938(G)	Frequency (YRI)	Frequency (AAF)	F_{ST} (AAF-YRI)	r^2 (YRI)
rs10216027	C/T	T	0.29	0.16	0.0569	0.514
rs1404315	A/C	C	0.32	0.16	0.0693	0.489
rs1722504	C/T	C	0.38	0.22	0.0652	0.403

4

5