



A probabilistic method for testing and estimating selection differences between populations

YUNGANG HE, Minxian Wang, Xin Huang, et al.

Genome Res. published online October 13, 2015

Access the most recent version at doi:[10.1101/gr.192336.115](https://doi.org/10.1101/gr.192336.115)

| | |
|---------------------------------|---|
| P<P | Published online October 13, 2015 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. |
| Creative Commons License | This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see http://genome.cshlp.org/site/misc/terms.xhtml). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at http://creativecommons.org/licenses/by-nc/4.0/ . |
| 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>

Published by Cold Spring Harbor Laboratory Press

1

2 **A probabilistic method for testing and estimating selection differences**

3 **between populations**

4 Yungang He,¹ Minxian Wang,¹ Xin Huang,¹ Ran Li,¹

5 Hongyang Xu,¹ Shuhua Xu,¹ Li Jin^{1,2}

6

7 ¹ Chinese Academy of Sciences Key Laboratory of Computational Biology, Chinese
8 Academy of Sciences-Max Planck Society Partner Institute for Computational Biology,
9 Shanghai Institutes for Biological Sciences, Chinese Academy of Sciences, Shanghai
10 200031, China

11 ² State Key Laboratory of Genetic Engineering and Ministry of Education Key
12 Laboratory of Contemporary Anthropology, Collaborative Innovation Center for
13 Genetics and Development, School of Life Sciences, Fudan University, Shanghai
14 200433, China

15 **Corresponding author:** lijin@fudan.edu.cn (L.J.), yunganghe@picb.ac.cn (Y.H.)

16 **Running title:** Investigating selection differences of populations

17 **Key words:** natural selection, population divergence.

1 **Abstract**

2 Human populations around the world encounter various environmental challenges,
3 and consequently develop genetic adaptations to different selection forces. Identifying the
4 differences in natural selection between populations is critical for understanding the roles
5 of specific genetic variants in evolutionary adaptation. Although numerous methods have
6 been developed to detect genetic loci under recent directional selection, a probabilistic
7 solution for testing and quantifying selection differences between populations is lacking.
8 Here we report the development of a probabilistic method for testing and estimating
9 selection differences between populations. Using a probabilistic model of genetic drift
10 and selection, we showed that logarithm odds ratios of allele frequencies provide
11 estimates of the differences in selection coefficients between populations. The estimates
12 approximate a normal distribution and variance can be estimated using genome-wide
13 variants. This allows us to quantify differences in selection coefficients and to determine
14 the confidence intervals of the estimate. Our work also revealed the link between genetic
15 association testing and hypothesis testing of selection differences. It therefore supplies a
16 solution for hypothesis testing of selection differences. This method was applied to a
17 genome-wide data analysis of Han and Tibetan populations. The results confirmed that
18 both *EPAS1* and *EGLN1* genes are under statistically different selection in Han and
19 Tibetan populations. We further estimated differences in the selection coefficients for
20 genetic variants involved in melanin formation and determined their confidence intervals
21 between continental population groups. Application of the method to empirical data
22 demonstrated the outstanding capability of this novel approach for testing and quantifying
23 differences in natural selection.

1 **Introduction**

2 When anatomically modern humans emerged from Africa (Mellars 2006), and
3 subsequently colonized throughout the world (Mellars et al. 2013; Hellenthal et al. 2008),
4 they encountered many challenges, including essential environmental alterations, food
5 resource shifts, and infectious diseases (Hancock et al. 2010, 2011; Leffler et al. 2013).
6 The current huge size and wide distribution of modern human populations demonstrates
7 the evolutionary success of human beings, which intrigues and attracts geneticists to
8 investigate the natural selection and genetic adaptation of human populations. Studies of
9 natural selection, especially directional selection, focus mainly on beneficial heritable
10 traits and related genetic alterations (Fu and Akey 2013; Williams 2008). In recent years,
11 genetic alterations under directional selection have attracted more attention than ever
12 before. Consequently, some highly irregular genetic variants were discovered and further
13 explored using various approaches (Grossman et al. 2010; Xu et al. 2011; Vitti et al.
14 2013; Bhatia et al. 2011; Sabeti et al. 2007; Xiang et al. 2013; Kamberov et al. 2013;
15 Sabeti et al. 2006).

16 Directional selection usually involves genetic adaptation to local environments.
17 Comparison of selection differences between populations is therefore important in
18 genetic studies of directional selection. Differences in allele frequencies are indicators of
19 possible selection differences between populations. As a measure of frequency difference,
20 genetic distance, such as F_{ST} , is the most popular statistic in studies of natural selection
21 (Akey et al. 2002; Lewontin and Krakauer 1973). The two-dimensional site frequency
22 spectrum (2D-SFS) method was also designed to compare frequency differences between

1 populations and thus to identify selection differences (Nielsen et al. 2009). Selection
2 differences can also be detected by comparing selective sweeps in different populations,
3 such as cross-population extended haplotype homozygosity (XP-EHH) and cross-
4 population composite likelihood ratio methods (XP-CLR) (Chen et al. 2010; Sabeti et al.
5 2007). Unfortunately, these methods lack efficient strategies to identify statistical outliers
6 from the ‘background noise’ of genetic drift. Theoretical distributions of these statistics
7 are not known in closed-form expressions. All the methods determine confidence levels
8 based on empirical data distribution or computer simulation with limited prior knowledge
9 of the demographic history (Akey et al. 2002; Nielsen et al. 2009; Chen et al. 2010;
10 Sabeti et al. 2007). Although computer simulation can handle complicated genetic
11 scenarios, it is unlikely that the “real” population genetic history can be accurately
12 represented in computer simulations (Teshima et al. 2006). Furthermore, existing
13 approaches do not provide effective solutions to quantify selection differences between
14 populations.

15 In this report, we present a probabilistic method for estimating and testing selection
16 differences between populations. The theoretical distribution of the involved statistics is
17 well known and easy to compute. It enables us to conduct strict hypothesis testing
18 without tedious computer simulation. Our approach supplies estimates and their
19 confidence intervals for differences in selection coefficients. To demonstrate the
20 capability of our approach, we conducted statistical hypothesis testing on a whole-
21 genome dataset including samples of Han Chinese and Tibetan populations. The results
22 regarding the *EPASI* and *EGLNI* genes rejected the null hypothesis and confirmed their

1 significant differences in selection between the populations. We further estimated
 2 differences in the selection coefficients between continental populations for genetic
 3 variants involved in melanin formation.

4

5 **Result**

6 *Model*

7 In a scenario with two populations, we assumed that populations *A* and *B* have the
 8 same ancestral population *O*. For a given locus, we denoted the frequencies of mutated
 9 allele in the three populations as p_O^m , p_A^m , and p_B^m while frequencies of wild-type allele as
 10 p_O^w , p_A^w , and p_B^w , respectively. In a deterministic approximation with selection, the
 11 difference of logarithm ratio of frequencies was determined by divergence time t and
 12 selection coefficient s in the populations, say $\log(\frac{p_A^m}{p_A^w}) - \log(\frac{p_O^m}{p_O^w}) = s_A \times t$ and
 13 $\log(\frac{p_B^m}{p_B^w}) - \log(\frac{p_O^m}{p_O^w}) = s_B \times t$. Therefore, the difference of selection coefficients with
 14 uncertainty can be presented as

$$15 \quad \Phi = s_A - s_B = [\log(\frac{p_A^m}{p_A^w}) - \log(\frac{p_B^m}{p_B^w})] / t + \Omega ,$$

16 where $\Omega = \frac{1}{t} \sum_{i=1}^t [(\omega_{A,i}^w - \omega_{A,i}^m) - (\omega_{B,i}^w - \omega_{B,i}^m)]$ indicates uncertainty due to genetic drift (see
 17 supplementary for details)

1 ***Estimating***

2 Numbers of chromosomes sampling from populations *A* and *B* with mutated
 3 alleles are denoted C_A^m and C_B^m with those carrying wild-type alleles are denoted as C_A^w
 4 and C_B^w . When population divergence time t is large, the general effect of genetic drift
 5 Ω will approximate a normal distribution with mean zero following the central limit
 6 theorem (Feller 1968). The differences in the strength of natural selection between
 7 populations *A* and *B* can be estimated as

$$8 \quad \hat{\Phi} = E(s_B - s_A) = \frac{\log(Odds)}{t}, \quad (\text{Eq.1})$$

9 where $Odds = (C_A^m C_B^w) / (C_A^w C_B^m)$. Variance of the estimation could be calculated as

$$10 \quad \text{Var}(\hat{\Phi}) = \text{Var}[\log(Odds)] / t^2 + \text{Var}(\Omega). \quad (\text{Eq.2})$$

11 Consequently, 95% confidence interval of the estimation is determined as
 12 $\hat{\Phi} \pm 1.96 \cdot \text{std}(\hat{\Phi})$.

13 For a neural locus i , we have $\hat{\Phi}_i^2 = \text{Var}[\log(Odds_i)] / t^2 + \text{Var}(\Omega)$. Therefore, when
 14 a sample has n neural loci and the n is large, the general effect of genetic drift between
 15 population *A* and *B* can be estimated as

$$16 \quad \hat{\text{Var}}(\Omega) = \text{median}\{\hat{\Phi}_i^2 / 0.455 - \text{Var}[\log(Odds_i)] / t^2, n \geq i \geq 1\},$$

17 where the variance of the log-odds ratio could be effectively approximated as
 18 $\text{Var}[\log(Odds)] = 1 / C_A^m + 1 / C_B^w + 1 / C_A^w + 1 / C_B^m$.

1 **Testing**

2 It is straightforward to propose a statistic for natural selection of a candidate locus,
3 as follows:

$$4 \quad \delta = \hat{\Phi}^2 / \text{Var}(\hat{\Phi}) . \quad (\text{Eq.3})$$

5 Under the null hypothesis that differences in natural selection are absent, the statistic δ
6 follows a central chi-square distribution with a degree of freedom = 1. Under the
7 alternative hypothesis with a selection difference, the statistic δ has a noncentral chi-
8 square distribution with non-centrality parameter $\hat{\Phi}^2$ and a degree of freedom = 1.

9 The aforementioned statistical test for a single candidate locus could be
10 generalized for a scenario with multiple linked loci to boost its power for detecting
11 differences. We can rewrite the statistic as

$$12 \quad \delta = X' \Sigma^{-1} X$$

13 where X is a vector with elements $\{\hat{\Phi}_1, \hat{\Phi}_2, \dots, \hat{\Phi}_m\}$ and Σ is the covariance matrix of the
14 vector with elements

$$15 \quad \text{Cov}(\hat{\Phi}_i, \hat{\Phi}_j) = \begin{cases} \text{Cov}[\log(\text{odd}_i), \log(\text{odd}_j)] + \text{Var}(\Omega), & i = j \\ \text{Cov}[\log(\text{odd}_i), \log(\text{odd}_j)], & i \neq j \end{cases} .$$

16 The covariance of two correlated log-odds ratios is given as (Bagos 2012)

$$17 \quad \text{Cov}[\log(\text{odd}_i), \log(\text{odd}_j)] = \sum_k \sum_l \sum_m (-1)^{l-m} \left(\frac{C_{klm}}{C_{kl+} C_{k+m}} \right) .$$

1 The notations for the covariance calculation are defined in Table 1. When testing for
2 multiple linked loci, the statistic δ approximates a central chi-square distribution under
3 the null hypothesis and the degree of freedom is the same as the number of involved loci
4 (De Maesschalck et al. 2000).

5 *Connection with case-control studies and its statistical power*

6 The theoretical framework presented above bears an intrinsic conceptual and
7 statistical connection with population-based association studies, as presented in Figure 1.
8 The left panel illustrates the conceptual framework of the null hypothesis and alternative
9 hypothesis in a genetic association study with the population stratification described by
10 Devlin et al. (2001) (Devlin et al. 2001). Genetic association studies detect indirect
11 associations between genetic markers (G) and phenotype (Y) that are mediated by
12 correlations between the genetic cause (X) and phenotype (Y). A special approach, such
13 as genomic control (GC), is capable of eliminating sporadic associations due to genetic
14 confounding effects (C). The right panel presents the framework of the null hypothesis
15 and alternative hypothesis in our method. The difference between the panels is that our
16 method focuses on the difference in selection (D) but not phenotype (Y). Our method can
17 distinguish selection differences from genetic background noise (B). GC controls type I
18 errors using an inflation factor λ , while our method considers differences in the genetic
19 background in the variance calculation by introducing $\text{Var}(\Omega)$ (Eq.2). The GC method

1 remains the same as a regular association test if $\lambda \approx 1$; our method also degenerates to a
2 regular association test if $\text{Var}(\Omega)$ approximates zero.

3 As the statistic of our method follows a chi-square distribution, the statistical power
4 of the research design can be conveniently calculated. In this study, with a given
5 difference in selection coefficients of 5.0×10^{-3} per generation, we show examples to
6 demonstrate how sample size, genetic drift, and divergence time contribute to the
7 statistical power of our method. Given a population divergence time of 300 generations,
8 our calculation indicates that the statistical power effectively increases with an increase in
9 the sample size (Figure 2 A). With a sample size of 500 chromosomes for each of the
10 paired populations and genetic drift per generation $\text{Var}(\Omega) = 1.0 \times 10^{-6}$, the statistical
11 power of our method is as high as 0.98 (Figure 2 A). The power increase, however, is
12 limited with an increase in the sample size when genetic drift is large. With genetic drift
13 per generation $\text{Var}(\Omega) = 5.0 \times 10^{-5}$, power is only ~ 0.20 , even if we have a sample size as
14 large as 500 chromosomes for each population (Figure 2A).

15 We also investigated the relationship between the population divergence time and
16 statistical power. The increase in power with an increase in the sample size is prominent
17 when the divergence time of involved populations is small (power curve marked by ‘o’ or
18 ‘*’, Figure 2B). When the divergence time is large, however, an increase in the sample
19 size has only a minor effect on statistical power (power curve marked by ‘◇’ or ‘x’,
20 Figure 2B). This could be due to the fact that accumulated genetic drift contributes

1 significantly to the statistic's variance in this scenario. Because our method is based on
2 allele frequencies of individual loci, but not a strict selective sweep, it is especially
3 helpful for studying a 'soft sweep'. Other selection-sweep based methods, such as XP-
4 EHH and XP-CLR, cannot work as well without significant linkage disequilibrium. We
5 therefore strongly suggest that both our method and selection-sweep based methods
6 should be applied as complementary methods to selection identification. Furthermore, our
7 method supplies an estimate for differences in selection coefficients, whereas the others
8 do not.

9 *Testing selection differences between Tibetan and Han genomes*

10 We applied our method of hypothesis testing on genotype data of Tibetan and Han
11 Chinese. Because several genetic loci are reported to be involved in adaptation to high-
12 altitude, most of these were not further verified by strict hypothesis testing but solely by
13 inspection in simulation-based inference. A QQ plot of our single-variant testing showed
14 that the obtained p-value was well fitted to the expectation (Figure 3), suggesting that our
15 theoretical model handled genetic divergence of populations well, as least for this
16 example. In particular, population divergence between Han and Tibetan populations did
17 not lead to inflation of the type I error in our hypothesis testing.

18 The criterion to declare a genome-wide statistical significance is given by p-value
19 $\leq 1.0 \times 10^{-8}$ in this study. Nineteen variants of the *EPAS1* gene have p-values that fit the
20 criterion (Figure 4A). This observation agrees with previous reports suggesting that the
21 *EPAS1* gene plays a major role in the high-altitude adaptation of Tibetan people (Xu et al.
22 2011; Simonson et al. 2010; Peng et al. 2011). We also conducted our aforementioned

1 multi-variant analysis on single nucleotide polymorphism (SNP) bins with different sizes.
2 With a bin size of 5, 10, or 15 SNPs, SNP bins in the *EGLNI* gene region showed
3 significant selection differences between the populations in our genome-wide hypothesis
4 testing (Figure 4B). The power increase was consistent with previous reports that a multi-
5 variant analysis could be more powerful than a single-variant approach in statistical tests
6 of genetic data (Akey et al. 2001; He et al. 2011). These results support previous findings
7 that both *EPASI* and *EGLNI* genes are critical to high-altitude adaptation of the Tibetan
8 population (Lorenzo et al. 2014). We obtained no positive findings in other gene regions,
9 except for *EPASI* and *EGLNI*. Other reported candidate genes should be further verified
10 when more genetic data becomes available.

11 ***Estimating differences in selection***

12 We estimated differences in the selection coefficients for several genetic variants of
13 melanin formation between continental populations (Table 2). In this study, we assumed
14 a simplified stepping stone model with five worldwide population groups (Figure 5).
15 Selection differences were compared between neighboring groups while mutated alleles
16 of the ancestral group served as a reference to determine the direction of the selection
17 differences (Eq. 1).

18 Our estimations and their 95% confidence intervals suggested that most of the
19 involved variants had similar selection coefficients in south and north Eurasian groups,
20 except variant rs12913832 of the *OCA2* gene had an obvious difference $\hat{\Phi} = 4.87 \times 10^{-3}$
21 (Figure 6). For south and north Eurasian groups, 95% confidence intervals of the
22 estimations were larger than those of population-group pairs involving both African and

1 non-African groups (Figure 6). This finding indicated that sampling variance contributed
2 to the variance of the estimations of the south and north Eurasian groups. Therefore, the
3 estimations could be further improved by increasing the sample sizes. Selection
4 coefficients had only minor differences between Asians and Africans (Figure 6),
5 suggesting that there are other genetic variants having a critical role in melanin formation
6 in Asians (Edwards et al. 2010). The observed directions of the selection differences
7 suggest that mutated alleles of the variants involved in melanin formation were more
8 favorably selected in non-African populations (Wilde et al. 2014).

9

10 **Discussion**

11 We measured the differences in allele frequencies between populations using their
12 logarithm odds ratios. Because genetic association studies usually present the effect size
13 of risk alleles in odds ratios with estimated confidence intervals, this study revealed a
14 statistical connection between our approach and classical genetic association studies. The
15 close connection further allowed us the opportunity to explore natural selection in a
16 genome-wide statistical test. There are other statistics with statistical properties better
17 than logarithm odds ratio, especially when sample size is limited. As we present in this
18 report, however, the logarithm odds ratio is an estimate of differences in the selection
19 coefficient while the other statistics lack a direct connection with selection difference.
20 Further, performance of logarithm odds ratio was acceptable in the presented case of the
21 Han-Tibetan comparison, demonstrating the merits of logarithm odds ratio. When
22 population divergence is small, variance of our estimate is due mainly to sampling

1 variance but not genetic drift (Eq.2, Figure 2). It is therefore possible to significantly
2 improve the power of the statistical test by increasing the sample sizes. In this scenario,
3 the benefit introduced by the large sample size is similar to that in genetic association
4 studies. The statistical power of the hypothesis test using our method can be calculated
5 for a specified study design. This provides a great advantage for determining the
6 technical details of a research design, especially for determining sample sizes. When the
7 evaluated locus is neutral in one of the two involved populations, our method provides
8 estimation for selection coefficient in the rest population.

9 In our genetic model, the overall effects of demographic impact are summarized by
10 variance in genetic drift (Eq.2). It is therefore unnecessary to separately consider the scale
11 and duration of each demographic event in the analysis. The scales and durations of
12 demographic events of the populations are often unknown, although some consensus has
13 been reached in the research community. Tedious computer simulation is unnecessary in
14 our approach while simulation is the only way to determine the confidence level in most
15 previous reports. This simulation-free feature is a significant advantage for selection
16 studies because “real” population history is unlikely to be accurately represented by
17 computer simulation. It should be noted that our method of modeling genetic drift differs
18 from the Wright-Fisher process. We use total variance $Var(\Omega)$ to capture the overall
19 effect of genetic drift but not effective sample sizes.

20 There are other statistics that measure the differences in allele frequencies between
21 populations, such as F_{ST} and ΔDAF . Both F_{ST} and ΔDAF have been applied to studies of
22 natural selection (Akey et al. 2002). There is a close relationship between our logarithm

1 odds ratio with F_{ST} and ΔDAF . When F_{ST} or the absolute value of ΔDAF is larger, we
2 generally have a larger positive or smaller negative logarithm odds ratio. Theoretical
3 distributions of the F_{ST} and ΔDAF statistics, however, are not available in straightforward
4 approaches. It therefore hinders their application in testing of selection difference.
5 Furthermore, in the presence of population stratification, there is no convenient approach
6 for quantifying the contribution of natural selection to F_{ST} and ΔDAF of individual
7 variants. There lack a perfect quantitative correlation between the statistics.

8 In our genetic model, we considered only the mutations that occurred before the
9 population stratification. This assumption holds for most genetic variants of the human
10 genome, given its short evolutionary history. Our method is therefore applicable to
11 populations with limited genetic divergence (Figure 2). When the frequency of an allele
12 is low and the sample size is small, minor alleles may be missing from the samples. In
13 these cases, we suggested a continuity correction in the calculation of the logarithm odds
14 ratios and the variance (Friedrich et al. 2007). Consequently, differences in selection
15 coefficients may be underestimated in this scenario. This potentially biased estimation
16 could be partially improved in two ways. First, a larger sample size may be helpful for
17 counting the minor allele; second, Bayesian estimation may be helpful for determining
18 the frequency of the missing allele.

19 To summarize, we developed a probabilistic method for testing and estimating
20 selection differences between populations. This method offers a statistical solution to
21 study directional selection without tedious computer simulation. It is very powerful when
22 the populations under investigation have close genetic connection. This method can be

1 used to quantify differences in selection coefficients, but not genotype fitness. Efficient
2 estimation of genotype fitness remains a difficult task when no time-serial data is
3 available.

4

5 **Material and Methods**

6 *Data*

7 Genotype data for 137 Han Chinese and 123 Tibetan unrelated individuals from
8 three previous studies of human high-altitude adaptation were analyzed in this report
9 (Wuren et al. 2014; Xu et al. 2011; Xing et al. 2013). All involved individuals were
10 genotyped using Affymetrix Genome-Wide Human SNP Array 6.0. To investigate
11 differences in the selection of genetic variants involved in melanin formation, genotype
12 data of worldwide populations were downloaded from the website of the 1000 Genomes
13 Project (1KG) (The 1000 Genomes Project Consortium 2010).

14 *Computing*

15 Haplotypes of the individuals were reconstructed using BEAGLE (Version 4.0)
16 (Browning and Browning 2007). Other computing works of this report were conducted in
17 R (version 2.14.2) (R Core Team 2015), a free software environment for statistical
18 computing and graphics.

19

20 **Acknowledgements**

1 We thank three anonymous reviewers for their comments to improve this work. This
2 work was supported by grants from National Natural Science Foundation of China
3 (91331109 and 31171279 to Y.H.; 31271338 and 31330038 to L.J.; 91331204 and
4 31171218 to S.X.). S.X. was also supported by the Strategic Priority Research Program
5 of the Chinese Academy of Sciences (XDB13040100). L.J. was also supported by
6 Shanghai Leading Academic Discipline Project (B111) and the Center for Evolutionary
7 Biology at Fudan University. Y.H. was also grateful for the supports of SA-SIBS
8 scholarship program and the Youth Innovation Promotion Association of Chinese
9 Academy of Science.

10 **Competing interest statement**

11 The authors have declared that no competing interests exist.

12

13

14

15

1 **Figure Legends**

2 **Figure 1 Conceptual framework of statistical tests for an association study and our**
3 **method.** Left panel shows conceptual framework of an association study in the presence
4 of population stratification. Right panel shows the conceptual framework of our method
5 in the same manner. Upper figures show conceptual frameworks for the null (H0)
6 hypothesis; lower figures show conceptual frameworks of the alternative (H1) hypothesis.

7 **Figure 2 Statistical power of our single-variant method increasing with an increase**
8 **in the sample size.** Statistical power is represented by the y-axis and sizes of the
9 involved haplotypes are represented by the x-axis. Allele frequency of one population
10 was given to be constant at 0.9, and frequency of the other population was determined by
11 differences in selection coefficients of 5.0×10^{-3} per generation and divergence time. A.
12 Power curve with a constant divergence time of 300 generations is marked by different
13 symbols for different drift variances ('o' for $\text{Var}(\Omega) = 1.0 \times 10^{-6}$, '*' for $\text{Var}(\Omega) = 5.0 \times 10^{-6}$,
14 '◇' for $\text{Var}(\Omega) = 1.0 \times 10^{-5}$, 'x' for $\text{Var}(\Omega) = 2.0 \times 10^{-5}$) B. Power curve with constant drift
15 variance $\text{Var}(\Omega) = 5.0 \times 10^{-6}$ is marked by different symbols for different divergence times
16 ('o' for $t = 100$ generations, '*' for $t = 300$ generations, '◇' for $t = 600$ generations, 'x' for t
17 $= 1000$ generations)

18 **Figure 3 QQ plot of single-variant analysis of Han-Tibetan data.** Observed
19 significance levels are represented by the y-axis on a scale of $-\log_{10}(\text{p-value})$. Expected
20 quartile is represented by the x axis on the same scale.

1 **Figure 4 Manhattan plots of significance levels for analysis of Han-Tibetan data.**

2 Chromosomes are shown on the x-axis; y-axis shows significance levels in $-\log_{10}(P)$. A.
3 Manhattan plot of single-variant analysis of all autosomes. B. Manhattan plots of single-
4 and multi-variant analysis of chromosome 1 and 2. Bin sizes are shown under the x axis.

5 **Figure 5 Simplified stepping-stone model of five population groups.** Details of genetic

6 demographic history were ignored in the model, such as backward gene flows and genetic
7 admixture, etc. The five continental population groups are North Eurasian (NEU), South
8 Eurasian (SEU), African (AFR), and Asian (ASN). Divergence of African and non-
9 African groups was assumed to be 5000 generations. We further assumed that NEU and
10 SEU have a divergence time of 400 generations.

11 **Figure 6 Differences in selection coefficients between population groups.** Estimated

12 differences in selection coefficients are represented by the y-axis. Error bars indicate 95%
13 confidence interval. Estimation for each neighboring group pair is marked by a group
14 name and allele-frequency pie chart of the corresponding descendant group. Frequency of
15 the derived allele is represented by the light color in the pie chart. North Eurasian
16 population (NEU) is a combination of the 1000 genome population CEU, FIN, and GBR;
17 South Eurasian population (SEU) is a combination of population IBS and TSI; African
18 population (AFR) is combination of population YRI and LWK; Asian population (ASN)
19 is a combination of population CHB, CHS, and JPT.

20

21

22

1 **Table 1 Notations for the covariance calculation.** C_{klm} is the haplotype count from
 2 population ' k ' which carries alleles in states ' l ' and ' m ' at locus 1 and 2, respectively.

3

| | | Locus 1 | | | |
|--------------|-------|----------------|-----------|-----------|-----------|
| | | $l=1$ | $l=0$ | | |
| | | Locus 2 | | | |
| | | $m=1$ | $m=0$ | $m=1$ | $m=0$ |
| Pop A | $k=1$ | C_{111} | C_{110} | C_{101} | C_{100} |
| Pop B | $k=0$ | C_{011} | C_{010} | C_{001} | C_{000} |

4

5

1 **Table 2 Candidate sites involved in our estimation.** ‘AA’ indicates ancestral allele and
 2 ‘DA’ indicates derived allele.

3

| Gene | Chr | rs# | Coordinate | AA | DA | Reference |
|-----------------------|-----|------------|------------|----|----|---|
| <i>OCA2</i> | 15 | rs12913832 | 28365618 | A | G | (Eiberg et al. 2008; Sturm et al. 2008) |
| <i>TYRP1</i> | 9 | rs1408799 | 12672097 | T | C | (Nan et al. 2009; Pośpiech et al. 2014) |
| <i>TYR</i> | 11 | rs1042602 | 88911696 | C | A | (Durso et al. 2014; Pośpiech et al. 2014) |
| <i>DCT</i> | 13 | rs1407995 | 95096013 | T | C | (Edwards et al. 2010; Zhu et al. 2007) |
| <i>SLC24A5</i> | 15 | rs1426654 | 48426484 | G | A | (Tekola-Ayele et al. 2014; Basu Mallick et al. 2013; Durso et al. 2014) |
| <i>SLC45A2</i> | 5 | rs16891982 | 33951693 | C | G | (Durso et al. 2014; Fernandez et al. 2008; Branicki et al. 2008) |

4

5

1 **References**

- 2 Akey J, Jin L, Xiong M. 2001. Haplotypes vs single marker linkage disequilibrium tests:
3 what do we gain? *Eur J Hum Genet EJHG* **9**: 291–300.
- 4 Akey JM, Zhang G, Zhang K, Jin L, Shriver MD. 2002. Interrogating a high-density SNP
5 map for signatures of natural selection. *Genome Res* **12**: 1805–1814.
- 6 Bagos PG. 2012. On the covariance of two correlated log-odds ratios. *Stat Med* **31**: 1418–
7 1431.
- 8 Basu Mallick C, Iliescu FM, Möls M, Hill S, Tamang R, Chaubey G, Goto R, Ho SYW,
9 Gallego Romero I, Crivellaro F, et al. 2013. The light skin allele of SLC24A5 in
10 South Asians and Europeans shares identity by descent. *PLoS Genet* **9**: e1003912.
- 11 Bhatia G, Patterson N, Pasaniuc B, Zaitlen N, Genovese G, Pollack S, Mallick S, Myers
12 S, Tandon A, Spencer C, et al. 2011. Genome-wide comparison of African-
13 ancestry populations from CARE and other cohorts reveals signals of natural
14 selection. *Am J Hum Genet* **89**: 368–381.
- 15 Branicki W, Brudnik U, Draus-Barini J, Kupiec T, Wojas-Pelc A. 2008. Association of
16 the SLC45A2 gene with physiological human hair colour variation. *J Hum Genet*
17 **53**: 966–971.
- 18 Browning SR, Browning BL. 2007. Rapid and accurate haplotype phasing and missing-
19 data inference for whole-genome association studies by use of localized haplotype
20 clustering. *Am J Hum Genet* **81**: 1084–1097.
- 21 Chen H, Patterson N, Reich D. 2010. Population differentiation as a test for selective
22 sweeps. *Genome Res* **20**: 393–402.
- 23 De Maesschalck R, Jouan-Rimbaud D, Massart DL. 2000. The Mahalanobis distance.
24 *Chemom Intell Lab Syst* **50**: 1–18.
- 25 Devlin B, Roeder K, Wasserman L. 2001. Genomic control, a new approach to genetic-
26 based association studies. *Theor Popul Biol* **60**: 155–166.
- 27 Durso DF, Bydlowski SP, Hutz MH, Suarez-Kurtz G, Magalhães TR, Pena SDJ. 2014.
28 Association of genetic variants with self-assessed color categories in Brazilians.
29 *PloS One* **9**: e83926.
- 30 Edwards M, Bigham A, Tan J, Li S, Gozdzik A, Ross K, Jin L, Parra EJ. 2010.
31 Association of the OCA2 polymorphism His615Arg with melanin content in east
32 Asian populations: further evidence of convergent evolution of skin pigmentation.
33 *PLoS Genet* **6**: e1000867.

- 1 Eiberg H, Troelsen J, Nielsen M, Mikkelsen A, Mengel-From J, Kjaer KW, Hansen L.
2 2008. Blue eye color in humans may be caused by a perfectly associated founder
3 mutation in a regulatory element located within the *HERC2* gene inhibiting *OCA2*
4 expression. *Hum Genet* **123**: 177–187.
- 5 Feller W. 1968. *An Introduction to Probability Theory and Its Applications, Vol. 1*. 3rd
6 edition, pp244. Wiley, New York.
- 7 Fernandez LP, Milne RL, Pita G, Avilés JA, Lázaro P, Benítez J, Ribas G. 2008.
8 *SLC45A2*: a novel malignant melanoma-associated gene. *Hum Mutat* **29**: 1161–
9 1167.
- 10 Friedrich JO, Adhikari NKJ, Beyene J. 2007. Inclusion of zero total event trials in meta-
11 analyses maintains analytic consistency and incorporates all available data. *BMC*
12 *Med Res Methodol* **7**: 5.
- 13 Fu W, Akey JM. 2013. Selection and adaptation in the human genome. *Annu Rev*
14 *Genomics Hum Genet* **14**: 467–489.
- 15 Grossman SR, Shlyakhter I, Shlyakhter I, Karlsson EK, Byrne EH, Morales S, Frieden G,
16 Hostetter E, Angelino E, Garber M, et al. 2010. A composite of multiple signals
17 distinguishes causal variants in regions of positive selection. *Science* **327**: 883–
18 886.
- 19 Hancock AM, Alkorta-Aranburu G, Witonsky DB, Di Rienzo A. 2010. Adaptations to
20 new environments in humans: the role of subtle allele frequency shifts. *Philos*
21 *Trans R Soc Lond B Biol Sci* **365**: 2459–2468.
- 22 Hancock AM, Witonsky DB, Alkorta-Aranburu G, Beall CM, Gebremedhin A, Sukernik
23 R, Utermann G, Pritchard JK, Coop G, Di Rienzo A. 2011. Adaptations to
24 climate-mediated selective pressures in humans. *PLoS Genet* **7**: e1001375.
- 25 Hellenthal G, Auton A, Falush D. 2008. Inferring human colonization history using a
26 copying model. *PLoS Genet* **4**: e1000078.
- 27 He Y, Li C, Amos CI, Xiong M, Ling H, Jin L. 2011. Accelerating haplotype-based
28 genome-wide association study using perfect phylogeny and phase-known
29 reference data. *PloS One* **6**: e22097.
- 30 Kamberov YG, Wang S, Tan J, Gerbault P, Wark A, Tan L, Yang Y, Li S, Tang K, Chen
31 H, et al. 2013. Modeling recent human evolution in mice by expression of a
32 selected *EDAR* variant. *Cell* **152**: 691–702.
- 33 Leffler EM, Gao Z, Pfeifer S, Ségurel L, Auton A, Venn O, Bowden R, Bontrop R, Wall
34 JD, Sella G, et al. 2013. Multiple instances of ancient balancing selection shared
35 between humans and chimpanzees. *Science* **339**: 1578–1582.

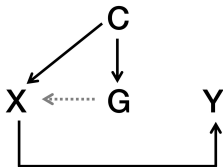
- 1 Lewontin RC, Krakauer J. 1973. Distribution of gene frequency as a test of the theory of
2 the selective neutrality of polymorphisms. *Genetics* **74**: 175–195.
- 3 Lorenzo FR, Huff C, Myllymäki M, Olenchock B, Swierczek S, Tashi T, Gordeuk V,
4 Wuren T, Ri-Li G, McClain DA, et al. 2014. A genetic mechanism for Tibetan
5 high-altitude adaptation. *Nat Genet* **46**: 951–956.
- 6 Mellars P. 2006. Why did modern human populations disperse from Africa ca. 60,000
7 years ago? A new model. *Proc Natl Acad Sci U S A* **103**: 9381–9386.
- 8 Mellars P, Gori KC, Carr M, Soares PA, Richards MB. 2013. Genetic and archaeological
9 perspectives on the initial modern human colonization of southern Asia. *Proc Natl*
10 *Acad Sci U S A* **110**: 10699–10704.
- 11 Nan H, Kraft P, Hunter DJ, Han J. 2009. Genetic variants in pigmentation genes,
12 pigmentary phenotypes, and risk of skin cancer in Caucasians. *Int J Cancer J Int*
13 *Cancer* **125**: 909–917.
- 14 Nielsen R, Hubisz MJ, Hellmann I, Torgerson D, Andrés AM, Albrechtsen A,
15 Gutenkunst R, Adams MD, Cargill M, Boyko A, et al. 2009. Darwinian and
16 demographic forces affecting human protein coding genes. *Genome Res* **19**: 838–
17 849.
- 18 Peng Y, Yang Z, Zhang H, Cui C, Qi X, Luo X, Tao X, Wu T, Ouzhuluobu null, Basang
19 null, et al. 2011. Genetic variations in Tibetan populations and high-altitude
20 adaptation at the Himalayas. *Mol Biol Evol* **28**: 1075–1081.
- 21 Pośpiech E, Wojas-Pelc A, Walsh S, Liu F, Maeda H, Ishikawa T, Skowron M, Kayser
22 M, Branicki W. 2014. The common occurrence of epistasis in the determination
23 of human pigmentation and its impact on DNA-based pigmentation phenotype
24 prediction. *Forensic Sci Int Genet* **11**: 64–72.
- 25 R Core Team. 2015. *R: A language and environment for statistical computing*. R
26 Foundation for Statistical Computing, Vienna, Austria.
- 27 Sabeti PC, Schaffner SF, Fry B, Lohmueller J, Varilly P, Shamovsky O, Palma A,
28 Mikkelsen TS, Altshuler D, Lander ES. 2006. Positive natural selection in the
29 human lineage. *Science* **312**: 1614–1620.
- 30 Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, Cotsapas C, Xie X, Byrne EH,
31 McCarroll SA, Gaudet R, et al. 2007. Genome-wide detection and
32 characterization of positive selection in human populations. *Nature* **449**: 913–918.
- 33 Simonson TS, Yang Y, Huff CD, Yun H, Qin G, Witherspoon DJ, Bai Z, Lorenzo FR,
34 Xing J, Jorde LB, et al. 2010. Genetic evidence for high-altitude adaptation in
35 Tibet. *Science* **329**: 72–75.

- 1 Sturm RA, Duffy DL, Zhao ZZ, Leite FPN, Stark MS, Hayward NK, Martin NG,
2 Montgomery GW. 2008. A single SNP in an evolutionary conserved region within
3 intron 86 of the HERC2 gene determines human blue-brown eye color. *Am J Hum*
4 *Genet* **82**: 424–431.
- 5 Tekola-Ayele F, Adeyemo A, Chen G, Hailu E, Aseffa A, Davey G, Newport MJ, Rotimi
6 CN. 2014. Novel genomic signals of recent selection in an Ethiopian population.
7 *Eur J Hum Genet EJHG*.
- 8 Teshima KM, Coop G, Przeworski M. 2006. How reliable are empirical genomic scans
9 for selective sweeps? *Genome Res* **16**: 702–712.
- 10 The 1000 Genomes Project Consortium. 2010. A map of human genome variation from
11 population-scale sequencing. *Nature* **467**: 1061–1073.
- 12 Vitti JJ, Grossman SR, Sabeti PC. 2013. Detecting natural selection in genomic data.
13 *Annu Rev Genet* **47**: 97–120.
- 14 Wilde S, Timpson A, Kirsanow K, Kaiser E, Kayser M, Unterländer M, Hollfelder N,
15 Potekhina ID, Schier W, Thomas MG, et al. 2014. Direct evidence for positive
16 selection of skin, hair, and eye pigmentation in Europeans during the last 5,000 y.
17 *Proc Natl Acad Sci U S A* **111**: 4832–4837.
- 18 Williams GC. 2008. *Adaptation and natural selection: a critique of some current*
19 *evolutionary thought*. Princeton University Press.
- 20 Wuren T, Simonson TS, Qin G, Xing J, Huff CD, Witherspoon DJ, Jorde LB, Ge R-L.
21 2014. Shared and unique signals of high-altitude adaptation in geographically
22 distinct Tibetan populations. *PLoS One* **9**: e88252.
- 23 Xiang K, Ouzhuluobu null, Peng Y, Yang Z, Zhang X, Cui C, Zhang H, Li M, Zhang Y,
24 Bianba null, et al. 2013. Identification of a Tibetan-specific mutation in the
25 hypoxic gene EGLN1 and its contribution to high-altitude adaptation. *Mol Biol*
26 *Evol* **30**: 1889–1898.
- 27 Xing J, Wuren T, Simonson TS, Watkins WS, Witherspoon DJ, Wu W, Qin G, Huff CD,
28 Jorde LB, Ge R-L. 2013. Genomic analysis of natural selection and phenotypic
29 variation in high-altitude mongolians. *PLoS Genet* **9**: e1003634.
- 30 Xu S, Li S, Yang Y, Tan J, Lou H, Jin W, Yang L, Pan X, Wang J, Shen Y, et al. 2011. A
31 genome-wide search for signals of high-altitude adaptation in Tibetans. *Mol Biol*
32 *Evol* **28**: 1003–1011.
- 33 Zhu G, Montgomery GW, James MR, Trent JM, Hayward NK, Martin NG, Duffy DL.
34 2007. A genome-wide scan for naevus count: linkage to CDKN2A and to other
35 chromosome regions. *Eur J Hum Genet EJHG* **15**: 94–102.

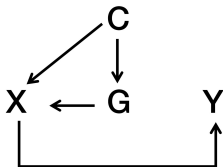
1

Test of genetic association

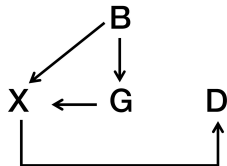
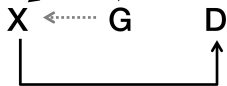
H0

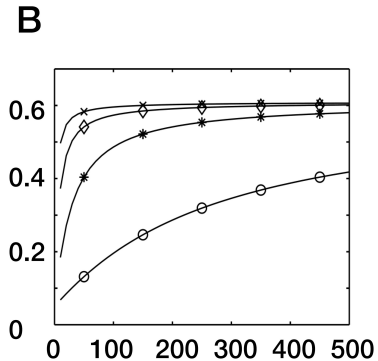
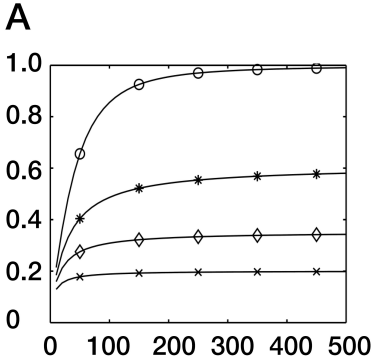


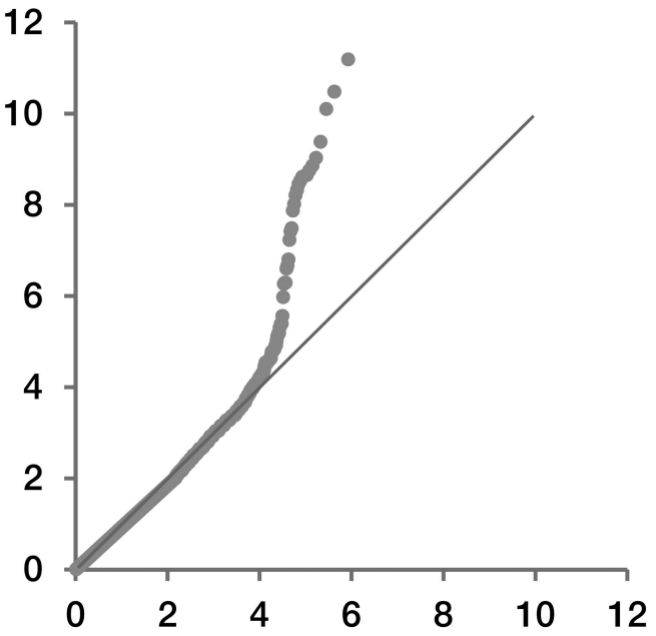
H1

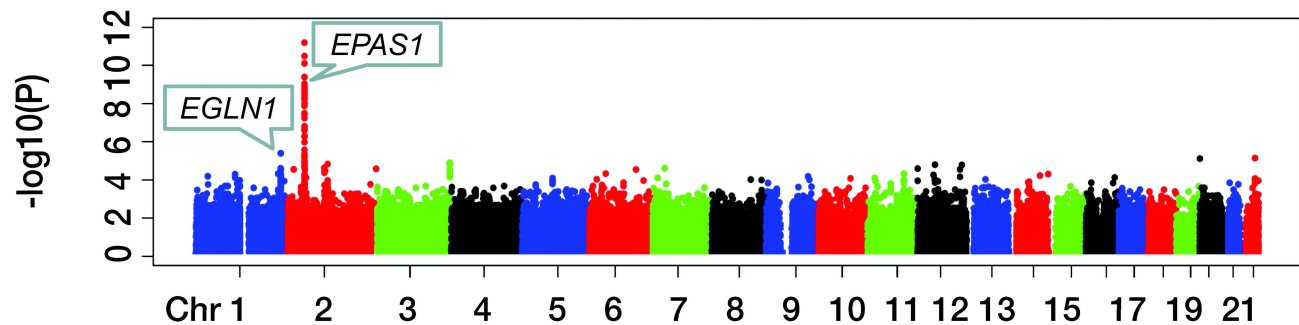
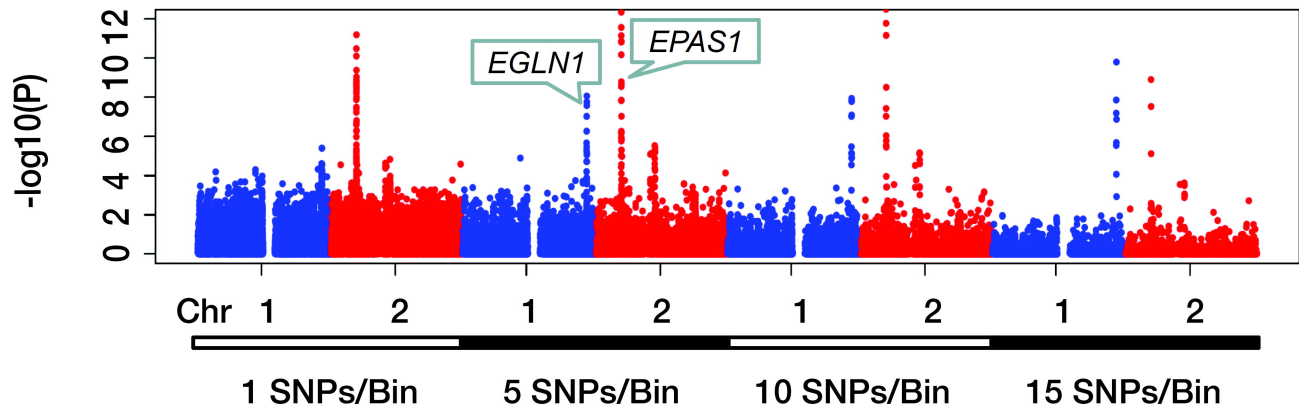


Test of selection difference

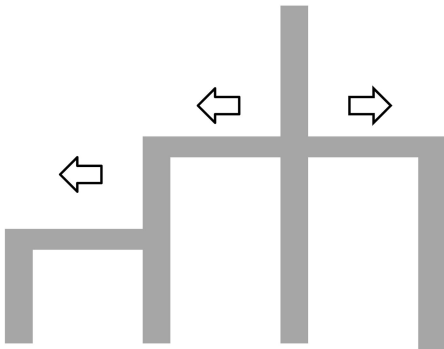






A**B**

MRCA

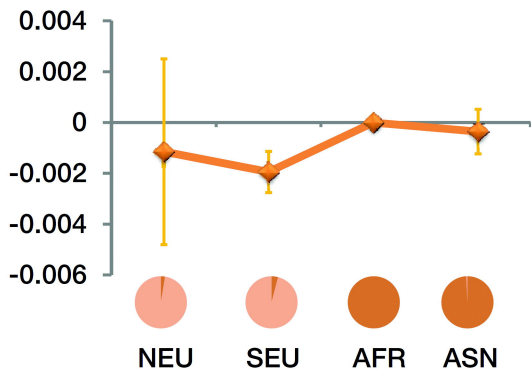
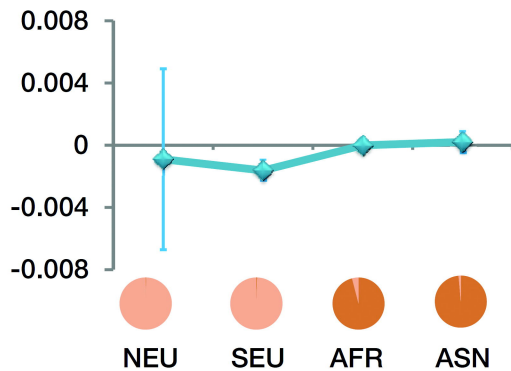
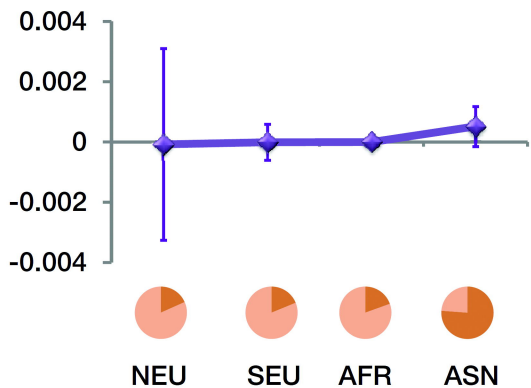
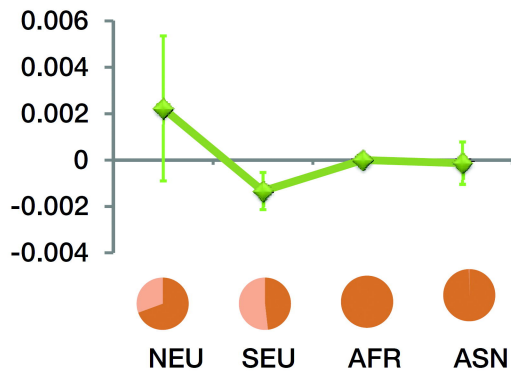
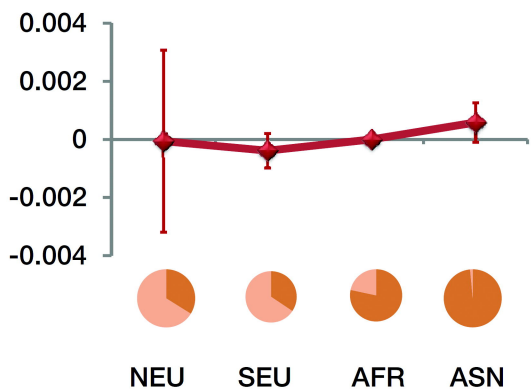


NEU

SEU

AFR

ASN

SLC45A2**SLC24A5****DCT****TYR****TYRP1****OCA2**