

RESEARCH ARTICLE

On the Apportionment of Population Structure

Yaron Granot^{1*}, Omri Tal², Saharon Rosset³, Karl Skorecki¹

1 Rappaport Faculty of Medicine and Research Institute, Technion–Israel Institute of Technology, and Rambam Medical Center, Haifa, Israel, **2** Max Planck Institute for Mathematics in the Sciences, Inselstr. 22-26, 04103, Leipzig, Germany, **3** School of Mathematical Sciences Tel Aviv University, Tel Aviv, Israel

* yarongranot@hotmail.com

Abstract

Measures of population differentiation, such as F_{ST} , are traditionally derived from the partition of diversity within and between populations. However, the emergence of population clusters from multilocus analysis is a function of genetic *structure* (departures from panmixia) rather than of diversity. If the populations are close to panmixia, slight differences between the mean pairwise distance within and between populations (low F_{ST}) can manifest as strong separation between the populations, thus population clusters are often evident even when the vast majority of diversity is partitioned within populations rather than between them. For any given F_{ST} value, clusters can be tighter (more panmictic) or looser (more stratified), and in this respect higher F_{ST} does not always imply stronger differentiation. In this study we propose a measure for the partition of structure, denoted E_{ST} , which is more consistent with results from clustering schemes. Crucially, our measure is based on a statistic of the data that is a good measure of internal structure, mimicking the information extracted by unsupervised clustering or dimensionality reduction schemes. To assess the utility of our metric, we ranked various human (HGDP) population pairs based on F_{ST} and E_{ST} and found substantial differences in ranking order. E_{ST} ranking seems more consistent with population clustering and classification and possibly with geographic distance between populations. Thus, E_{ST} may at times outperform F_{ST} in identifying evolutionary significant differentiation.



OPEN ACCESS

Citation: Granot Y, Tal O, Rosset S, Skorecki K (2016) On the Apportionment of Population Structure. PLoS ONE 11(8): e0160413. doi:10.1371/journal.pone.0160413

Editor: Francesc Calafell, Universitat Pompeu Fabra, SPAIN

Received: December 2, 2015

Accepted: July 19, 2016

Published: August 9, 2016

Copyright: © 2016 Granot et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: The SNP data used in this study have been accessed from the HGDP dataset: <http://www.hagsc.org/hgdp/files.html>.

Funding: The authors have no support or funding to report.

Competing Interests: The authors have declared that no competing interests exist.

Introduction

Genetic differentiation among populations is typically derived from the ratio of within- to between-population diversity. The most commonly used metric, F_{ST} , was originally introduced as a fixation index at a single biallelic locus [1], and subsequently adapted as a measure of population subdivision by averaging over multiple loci [2–3]. F_{ST} can be expressed mathematically in terms of population diversities as $F_{ST} = 1 - S/T$, where S and T represent the heterozygosity in subpopulations and in the total population, respectively [4–5]. The validity of F_{ST} as a measure of differentiation has been brought into question, especially when gene diversity is high (e.g., in

microsatellites), and various metrics, including G_{ST}^2 [6] and Jost's D [7], have been proposed to address this inadequacy (though see [8] for a counter-perspective).

Although these metrics vary considerably in their formulation, they all follow the same basic framework of partitioning genetic diversity into within- vs. between-group components. It has long been noted, however, that the apportionment of diversity [9] does not directly reflect the strength of separation between populations, and the emergence of population clusters has been demonstrated both empirically [10] and mathematically [5, 11–12] even when the vast majority of diversity is within rather than between populations. For example, humans sampled from across Europe [13] form identifiable clusters with pairwise F_{ST} as low as 0.001, even though 99.9% of the variation is contained within populations and only 0.1% is between them. Clearly, these clusters reflect an aspect of population differentiation that is not directly captured by F_{ST} , yet there is currently no commonly used metric for partitioning structure into within- and between-population components in the same way that F_{ST} partitions diversity. Dimensionality reduction schemes such as principal component analysis (PCA) [14] and clustering algorithms such as the widely used STRUCTURE [15] are highly popular, however such programs are primarily used in the population-genetics literature for visualization purposes, and there is still value in summary statistics for quantifying complex datasets on a simple 0–1 scale.

Here we propose a novel measure, denoted E_{ST} , which is based on the variation in pairwise genetic distance (which we show to be a measure of internal structure), thus exposing the excess structure inherent in the total population compared to subpopulations. Conceptually, E_{ST} is formulated in three steps: 1. Population structure is defined in terms of departures from *panmixia*. 2. Panmixia is defined in terms of pairwise genetic *equidistance* between individuals (we show that a population is panmictic if all individuals are equally distant from each other). 3. Departures from equidistance are defined in terms of the *standard deviation* of pairwise distances. E_{ST} reflects the decrease in panmixia when subpopulations are pooled. The basic formula is $E_{ST} = 1 - SD_S / SD_T$, where SD_S and SD_T represent the standard deviations of pairwise distances in subpopulations and in the total population respectively. While F_{ST} is weighed down (diminished) by high *diversity* within populations, E_{ST} is weighed down by high *structure* within populations.

A core insight of this proposal is that the asymptotic (in terms of number of genetic markers considered) standard deviation of pairwise genetic distances is a good unsupervised measure of internal structure, a statistic that mimics the information extracted by dimensionality reduction and clustering schemes, thus justifiable as a basis for the definition of E_{ST} . In particular, we show (in Appendix A) that a population is panmictic if and only if all individuals are *asymptotically* equally distant from each other, and that the standard deviation of pairwise genetic distances is highly associated with the deviation from panmixia and faithfully reflects the structure extracted by PCA.

Results and Discussion

Partitioning Diversity vs. Partitioning Structure

The difference between partitioning diversity and structure within and between two populations from the human genome diversity project (HGDP) [16] is illustrated in Fig 1. By zooming into a Russian and Chinese neighbor-joining tree of individual similarities we observe three layers of diversity and structure. Distances between individuals (black) account for most of the diversity, followed by the between-population component (red) and lastly, structure within populations (blue). The striking symmetry in the full-sized tree (1A) suggests high levels of panmixia in these two populations. Even at 100x magnification, most intra population branches (blue) are shorter than the 1x individual branches (black), indicating that these two

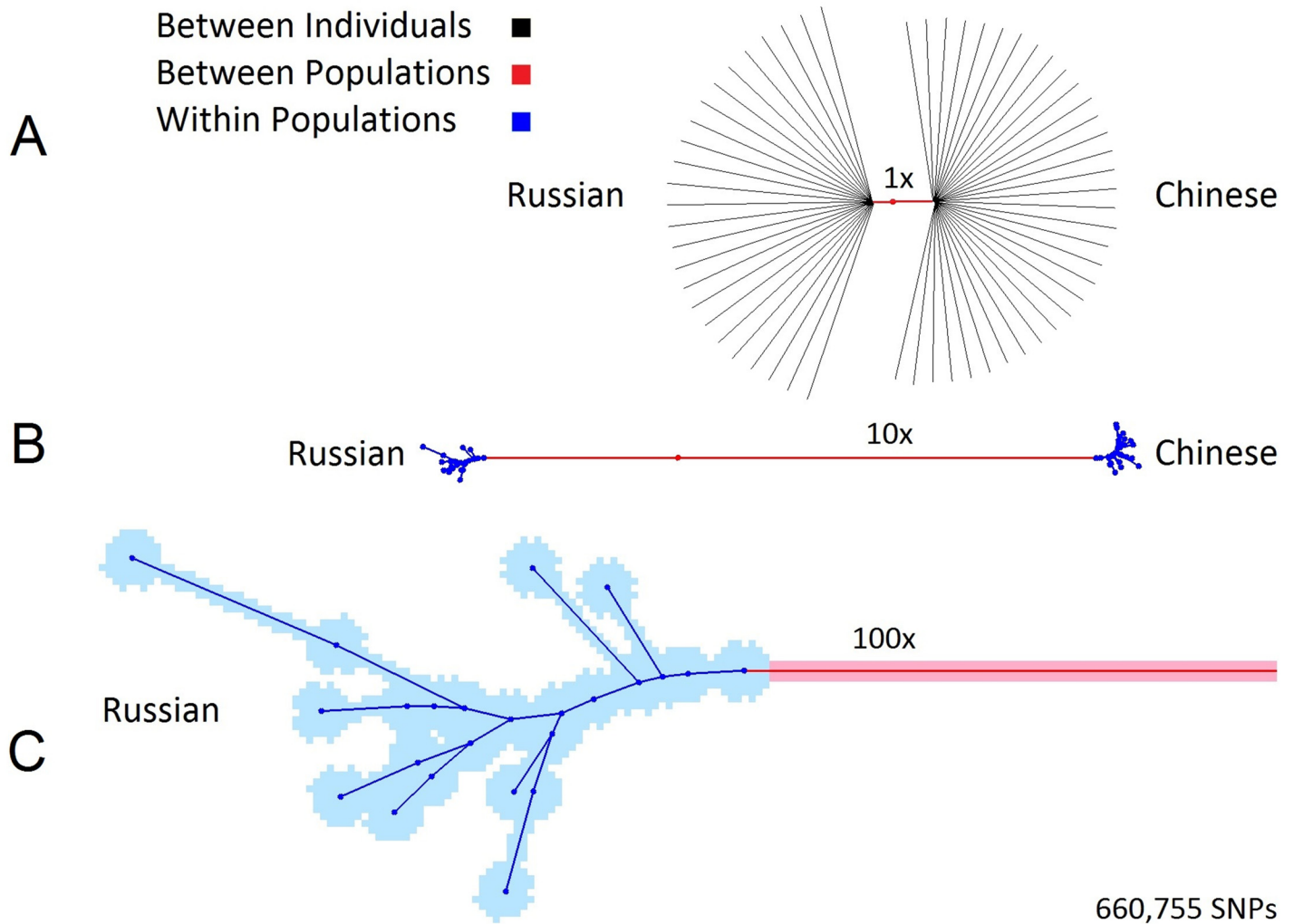


Fig 1. Zooming into a Russian (n = 25) and Chinese (n = 34) Neighbor Joining tree of individual similarities. (A) The length of the red branch compared to the overall tree length is a rough proxy to F_{ST} . (B) 10x magnification highlights the structure within and between populations. The blue branches inversely associate with E_{ST} values. (C) 100x magnification reveals fine substructure within Russian samples. Individual branches (black) were removed in B and C for clarity.

doi:10.1371/journal.pone.0160413.g001

sample sets are nearly panmictic. The red/black branch length ratio can be perceived as a rough proxy to the fixation index F_{ST} and (one minus) the blue/red branch length ratio can be perceived as a rough proxy to the equidistance index E_{ST} .

We compared F_{ST} , E_{ST} , and clustering among Russian and Chinese samples, with an increasing amount of single nucleotide polymorphisms (SNPs) ranging from 10 to 660,755 (Fig 2). Using multidimensional scaling (MDS), the two population clusters gradually diverge as SNP count increases, with no corresponding increase in F_{ST} . At the same time, we observe a steady increase in E_{ST} directly corresponding to the emerging clusters, indicating that the Russian and Chinese HGDP samples are close to panmixia. With few SNPs this is obfuscated by the variance of the genetic distance measure, hence E_{ST} is relatively small. The actual levels of panmixia become increasingly evident as more SNPs are added, thus revealing the population clusters [11]. However, this process does not proceed indefinitely; the finite number of pairwise differences among humans (~3 million SNPs) sets an upper limit to the number of available

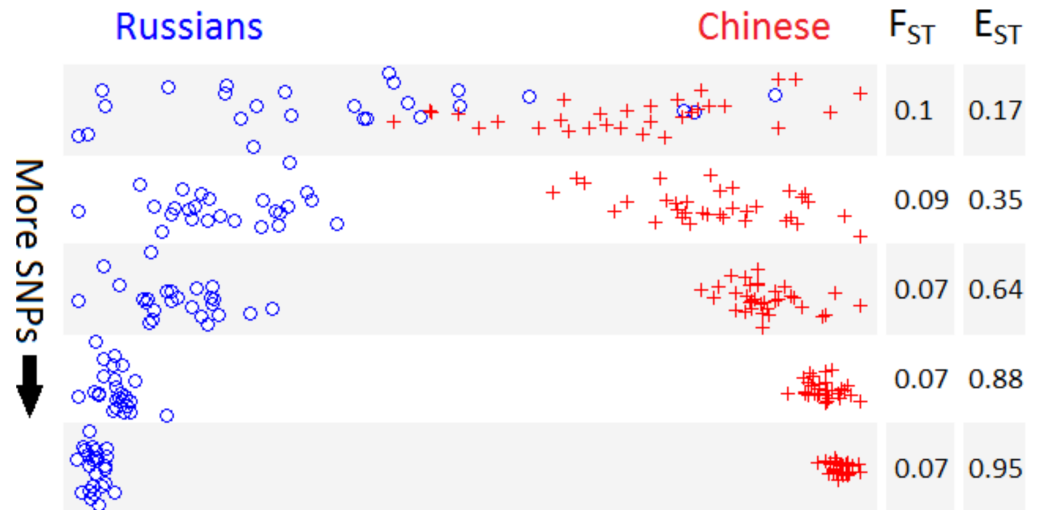


Fig 2. F_{ST} and E_{ST} vs. Clustering with increasing SNP count. Multidimensional scaling (MDS) plots with Russian ($n = 25$) and Chinese ($n = 34$) samples with increasing SNP count from top to bottom (10, 100, 1000, 10,000, and 660,755 SNPs). Two clusters gradually emerge as SNP count increases, along with an increase in E_{ST} , while F_{ST} remains relatively constant.

doi:10.1371/journal.pone.0160413.g002

markers, and the amount of extractable information is further reduced by physical linkage. In our HGDP data the increase in E_{ST} as a function of marker count reaches a plateau approximately above 100,000 SNPs (Fig 3). Although this upper bound can vary across different datasets and types of markers, it suggests that resolution may not improve substantially with further increases in marker count. Thus, these clusters can be considered close approximations of the “true” strength of separation among these populations. For this reason, when estimating E_{ST} one should include as many markers as possible, although at a certain point additional markers provide greatly diminishing returns.

In order to determine whether or not E_{ST} adds insight to the analysis of population structure, we sought to compare the rank order of population differentiation using F_{ST} and E_{ST} . Pairwise F_{ST} and E_{ST} values from various HGDP populations are given in Table 1 (see S1 Table and S1, S5 Figs for additional comparisons). It is noted that for almost all population pairs E_{ST} is larger than F_{ST} , and only the Colombian-Maya pair entails a slightly lower E_{ST} than F_{ST} , presumably due to a combination of relatively low differentiation and high levels of intra-population structure. According to the HGDP browser (http://spsmart.cesga.es/search.php?dataSet=ceph_stanford), the Colombians ($n = 7$) are the only HGDP population sample where two different tribes (Piapoco and Curripaco) were combined, which can help explain the high level of structure observed in this particular population (see S1 Table, S9, S5 Figs, and Materials and Methods for further analysis of E_{ST} range). There is a moderate positive correlation ($r = 0.61$) between F_{ST} and E_{ST} among all 60 population pairs included in our analysis (Fig 4), which is consistent with the two measures capturing somewhat different aspects of population structure. This seemingly high correlation can be attributed to the low (and therefore relatively similar) level of structure in many HGDP populations; the correlation between E_{ST} and F_{ST} tends to be higher when populations are more panmictic (see appendix).

Amazonians vs. Global Populations

A prominent example of the divergent behavior of E_{ST} and F_{ST} is the Surui and Karitiana, which have an unusually high pairwise F_{ST} . In fact, the Karitiana are as diverged from the

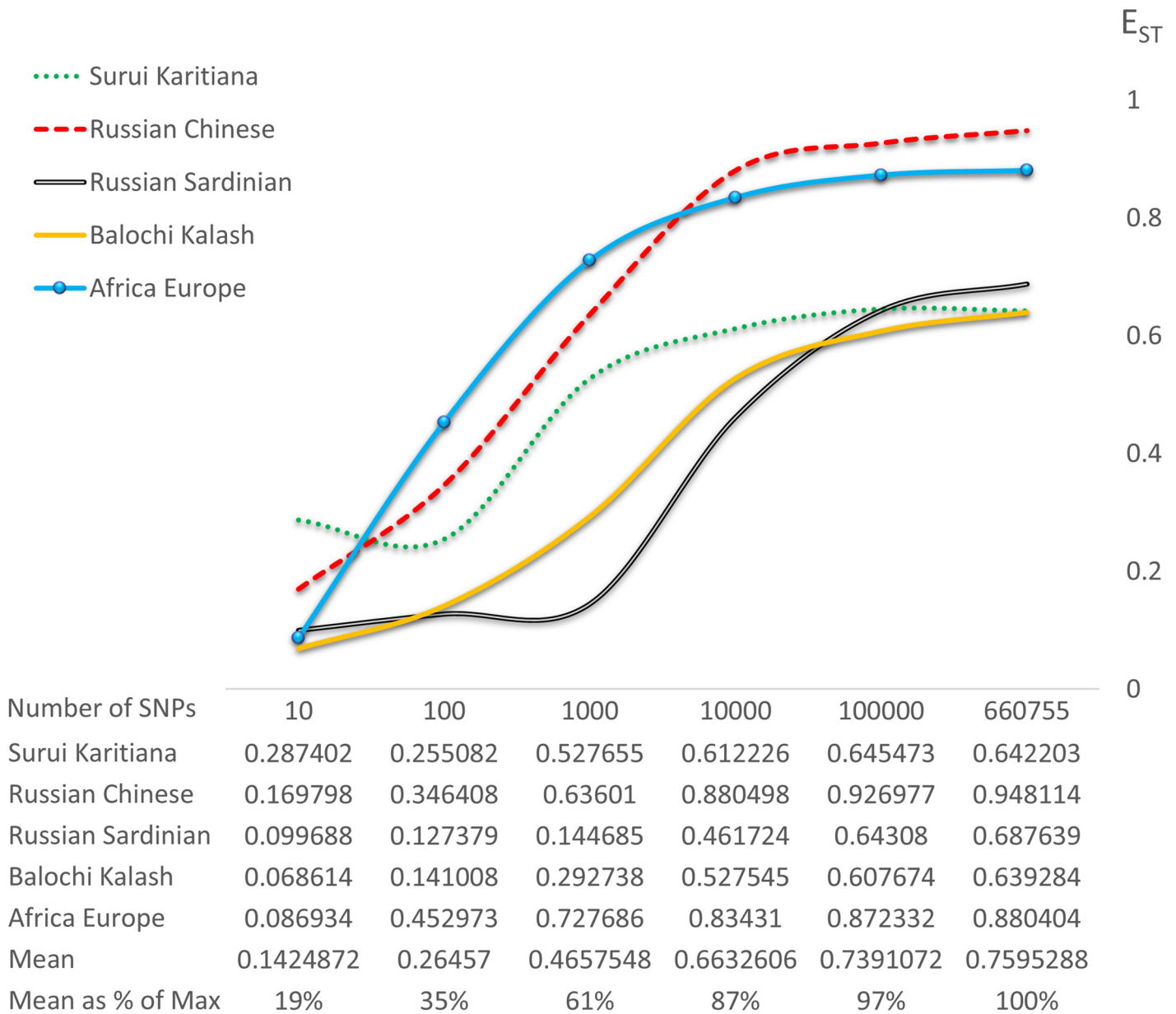


Fig 3. E_{ST} as a function of SNP sample size. E_{ST} was estimated in various population pairs with gradually increasing SNP sample size from 10 to 660,755. As expected, E_{ST} initially rises rather steeply, but tends to plateau before reaching the 660,755 SNP point. This suggests that we are approaching the maximal resolving power of genetic markers in this dataset, and adding markers beyond this point should not have a significant effect on ω -cluster separation and E_{ST} .

doi:10.1371/journal.pone.0160413.g003

neighboring Surui in terms of F_{ST} as they are from the Mongola on the other side of the world (Table 1, Fig 5, and S5 Fig). F_{ST} decreases initially with distance from the Amazon, from 0.13 between the two Amazonian tribes, to 0.08–0.1 between Amazonians and Colombians, and further decreases to 0.07–0.09 between Amazonians and the more distant Maya. Remarkably, the highest F_{ST} among all HGDP Native American populations is between the two geographically closest populations, the Surui and Karitiana. These apparent anomalies can be explained by the inflation of F_{ST} among genetic isolates. F_{ST} between pairs of isolates can be nearly twice as high as between either one of the isolates and a more cosmopolitan population, as pairwise

Table 1. Pairwise F_{ST} (above diagonal) and E_{ST} (below diagonal) in 5 New World and 5 Old World HGDP populations.

	Surui	Karitiana	Colombian	Maya	Pima	Yakut	Mongola	Russian	Bantu	San
Surui		0.13	0.1	0.09	0.12	0.15	0.15	0.17	0.23	0.3
Karitiana	0.58		0.08	0.07	0.11	0.13	0.13	0.16	0.22	0.29
Colombian	0.51	0.57		0.03	0.06	0.09	0.09	0.12	0.18	0.25
Maya	0.52	0.63	0.02		0.04	0.07	0.06	0.08	0.15	0.21
Pima	0.57	0.63	0.37	0.43		0.1	0.09	0.12	0.19	0.25
Yakut	0.74	0.8	0.6	0.69	0.74		0.01	0.06	0.13	0.19
Mongola	0.81	0.87	0.69	0.8	0.83	0.46		0.06	0.12	0.19
Russian	0.82	0.87	0.74	0.83	0.84	0.86	0.9		0.11	0.17
Bantu	0.88	0.92	0.85	0.91	0.91	0.93	0.94	0.95		0.07
San	0.92	0.95	0.89	0.95	0.94	0.96	0.97	0.98	0.89	

doi:10.1371/journal.pone.0160413.t001

F_{ST} reflects the *combined* isolation of both populations. Since the Surui and Karitiana are both isolated, their pairwise F_{ST} is nearly double that between any one of them and a larger, less isolated population such as the Maya.

Differentiation based on E_{ST} (Surui-Karitiana = 0.58, Karitiana-Mongola = 0.87, and Mongola-Bantu = 0.94) seems more consistent with the geographic distances among these populations (Fig 5). It should be noted that the Surui-Karitiana E_{ST} might be somewhat underestimated due to cryptic sampling of close relatives [17], however the wide range of

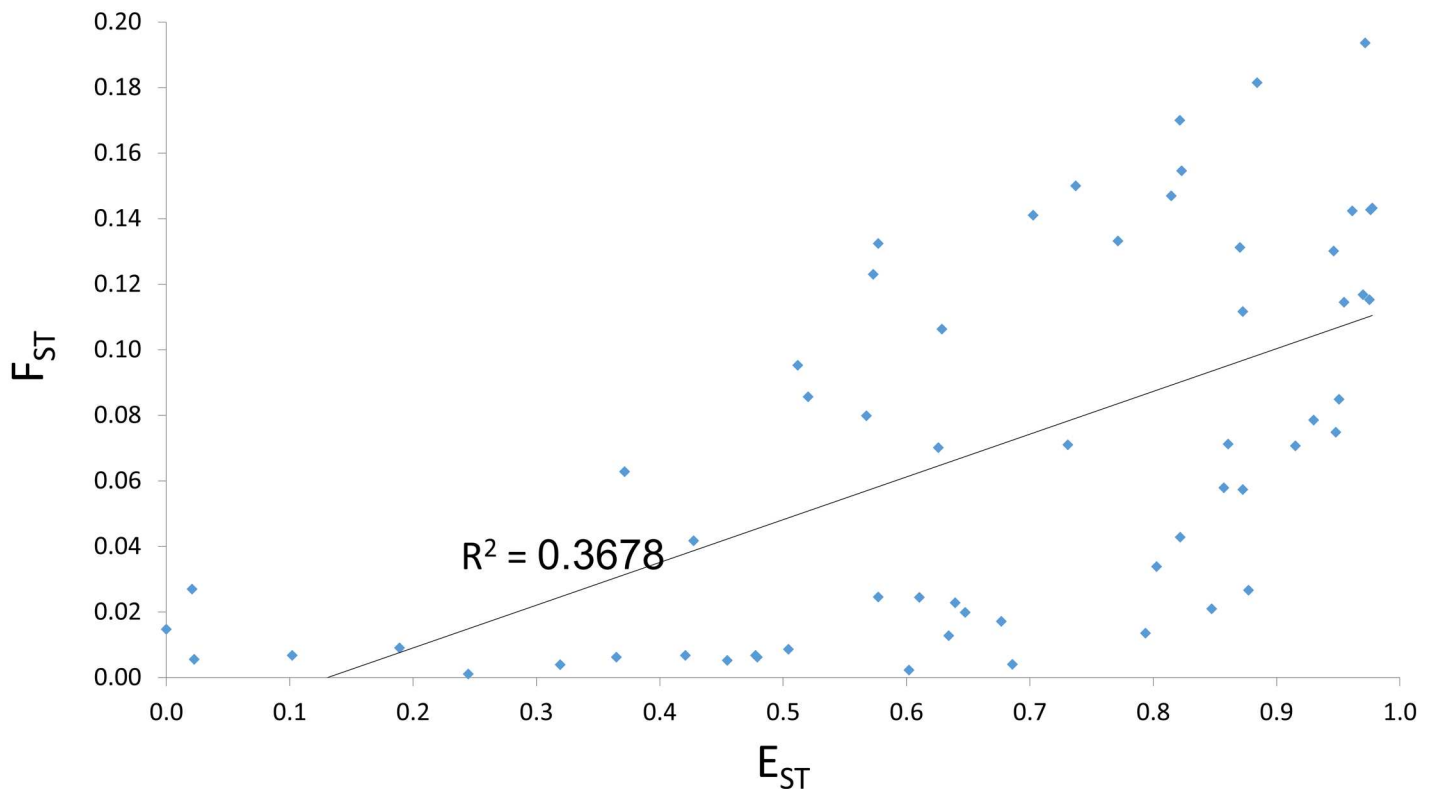


Fig 4. Positive correlation between F_{ST} and E_{ST} . $R = 0.61$. Note that E_{ST} has a much broader range, spanning nearly the entire 0–1 interval while F_{ST} only goes as high as 0.2 in these HGDP populations.

doi:10.1371/journal.pone.0160413.g004

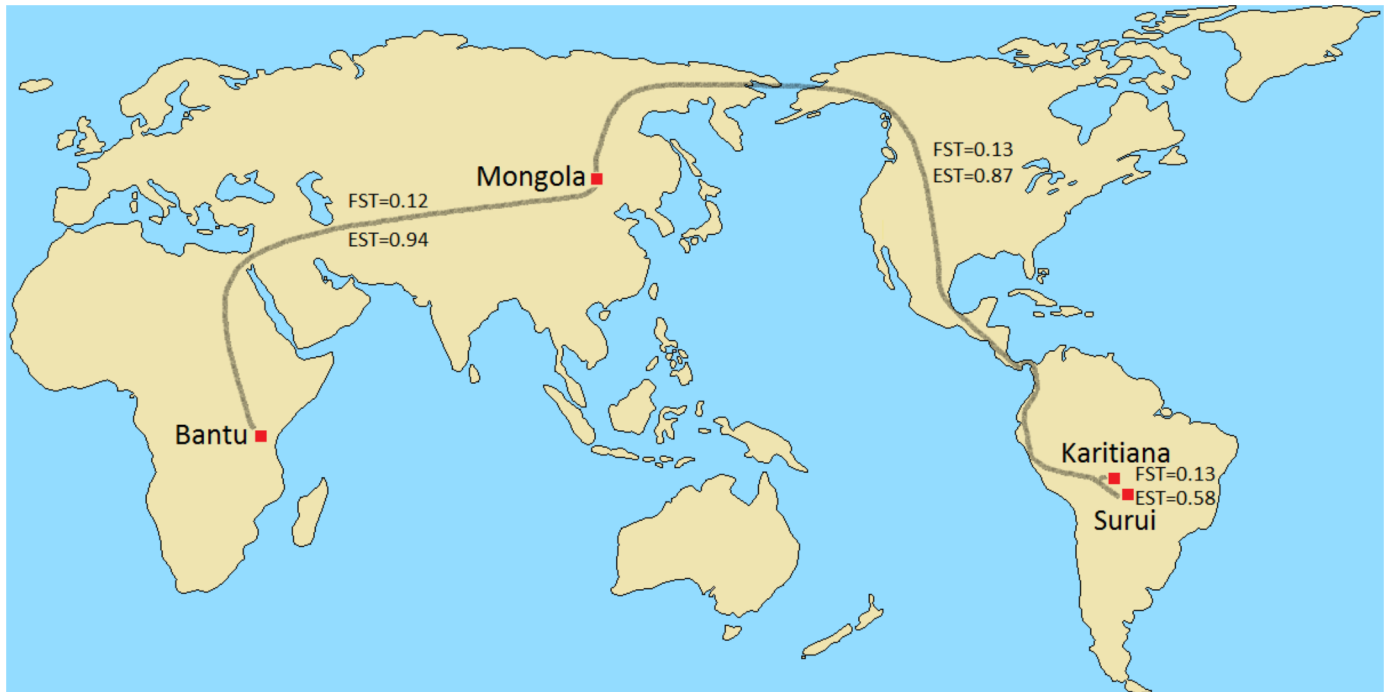


Fig 5. Geographic distance vs. F_{ST} and E_{ST} in various populations. In terms of F_{ST} , the Karitiana are roughly as diverged from the nearby Surui ($F_{ST} = 0.13$) as they are from the Mongola on the other side of the world ($F_{ST} = 0.13$) or as the Bantu are from the Mongola ($F_{ST} = 0.12$). In terms of E_{ST} , differentiation is far greater among these global populations ($E_{ST} \approx 0.9$) than between the neighboring Amazonian tribes ($E_{ST} \approx 0.6$).

doi:10.1371/journal.pone.0160413.g005

heterozygosity values (which are less sensitive to the sampling of close relatives) and the elevated structure across all Native American HGDP populations (S3–S5 Figs) suggest that this is not merely a sampling artifact. In some cases E_{ST} also decreases with distance from the Amazon (Table 1), however this decrease is more moderate than the decrease in F_{ST} (S5 Fig).

Neighbor-joining trees of individual similarities [18] are a convenient tool for representing multidimensional genetic data on a two-dimensional plane, while simultaneously displaying distances within and between populations. Two pairs of such trees, for Surui-Karitiana and Yoruba-Russians, are given in Fig 6, and we can see that in both cases distances are greater between individuals (black branches) than between populations (red branches) (Fig 6A).

The ratio of within- to between-population distance is roughly equivalent in the two population pairs, however the Yoruba-Russian tree is significantly *flatter*, indicating greater panmixia within these two populations (also see S6–S7 Figs). Adding a third dimension of intra-population structure (blue branches) highlights this discrepancy (Fig 6B), which is further accentuated by removing the inter-individual component (Fig 6C) and stretching the Yoruba-Russian tree to match the level of structure observed in the Surui-Karitiana tree (Fig 6D). At first glance the Amazonian tribes, with their long population branches, appear to be as differentiated as the Yoruba are from the Russians. Upon closer inspection, however, the Yoruba and Russians appear more strongly diverged. The Amazonian tribes are highly structured not only between them, but also within them, resulting in distant, but loosely separated clusters. This aspect of population structure is not captured by F_{ST} , which is actually slightly higher between the Surui and Karitiana (0.13) than between Yoruba and Russians (0.12), but is revealed by the higher E_{ST} between Yoruba and Russians (0.97) compared to the Surui and Karitiana (0.58).

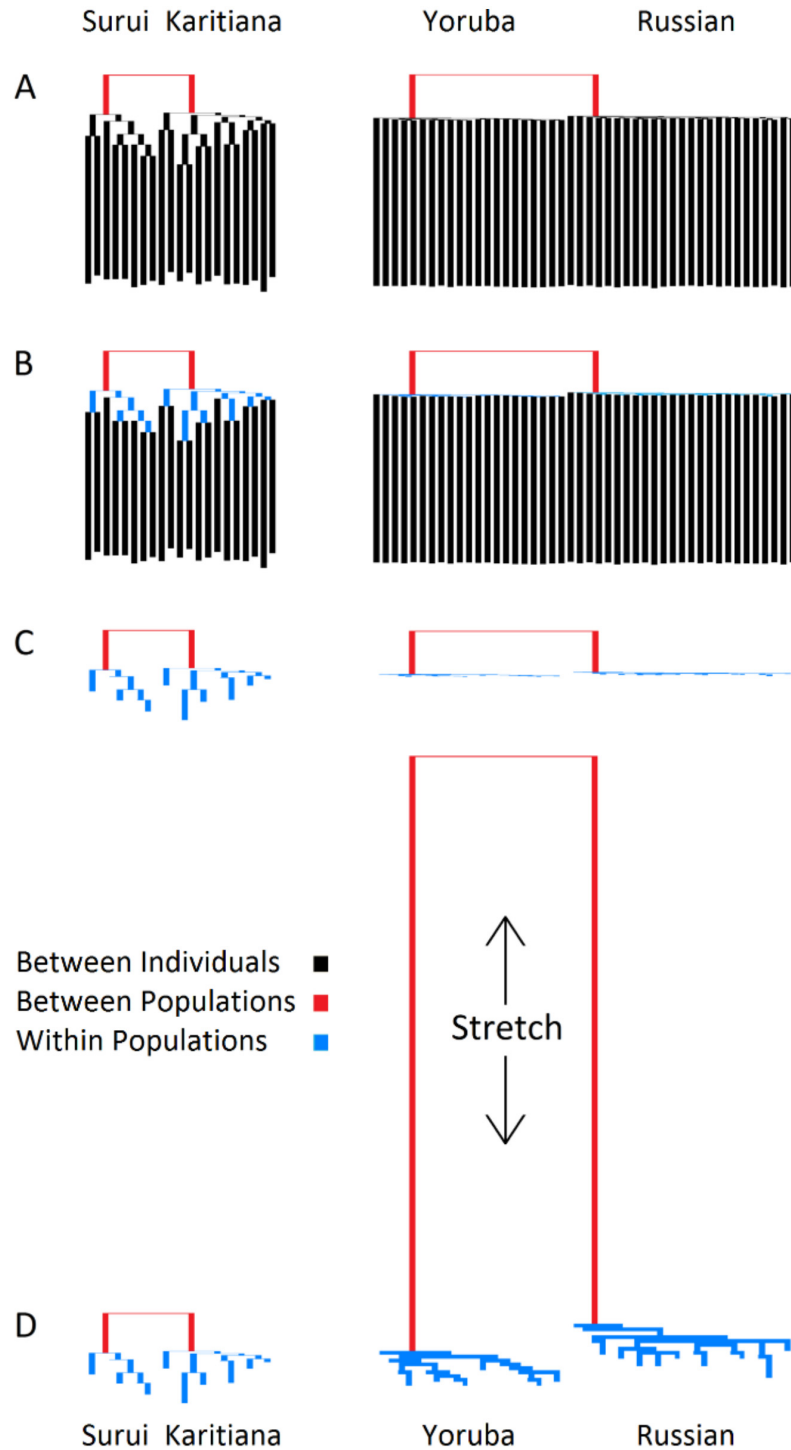


Fig 6. Surui-Karitiana vs. Yoruba-Russian NJ trees of individual similarities. (A) Diversity is apportioned into individual (black) and population (red) components. (B) A third component, structure within populations (blue), is added. (C) The individual component is removed. (D) The Yoruba-Russian tree is stretched to roughly match the level of structure within the Surui-Karitiana tree.

doi:10.1371/journal.pone.0160413.g006

E_{ST} and the Dissimilarity Fraction

Witherspoon et al. [10] have also examined population structure through the lens of pairwise genetic similarities and dissimilarities. They have defined the dissimilarity fraction, ω , as the probability that individuals are genetically more similar to members of a different population than to members of their own population. An intuitive proxy for ω is (half) the overlap of the within and between pairwise distance distributions. For population pairs, this probability has a 0–0.5 range, with the extremities $\omega = 0$ indicating that individuals are always more similar to members of their own population and $\omega = 0.5$ indicating that individuals are just as likely to be more similar to members of the other population as to members of their own population (see [5] for a formal analysis of such a metric and its relation to classification accuracy). Witherspoon et al. reported that when many thousands of loci are analyzed, individuals from “geographically separated populations” are never more similar to each other than they are to members of their own populations. The definition of “geographically separated” is, of course, open to interpretation. We found no overlap ($\omega = 0$) between the Adygei and Uygur HGDP samples, but some overlap ($\omega > 0$) between Mayans and Surui, despite a 4x higher F_{ST} (Fig 7). Thus, F_{ST} and the dissimilarity fraction (ω) are not necessarily congruent. The E_{ST} values for these two population pairs are more consistent with ω , showing strong separation between the Adygei and Uygur (0.79) and more moderate separation between Colombians and Maya (0.52) (see S8 Fig for a more detailed plot). The stronger association between ω and E_{ST} is not surprising since both measures are sensitive to the presence of within-population structure (e.g., in the Maya and Surui), a property that is not captured by F_{ST} .

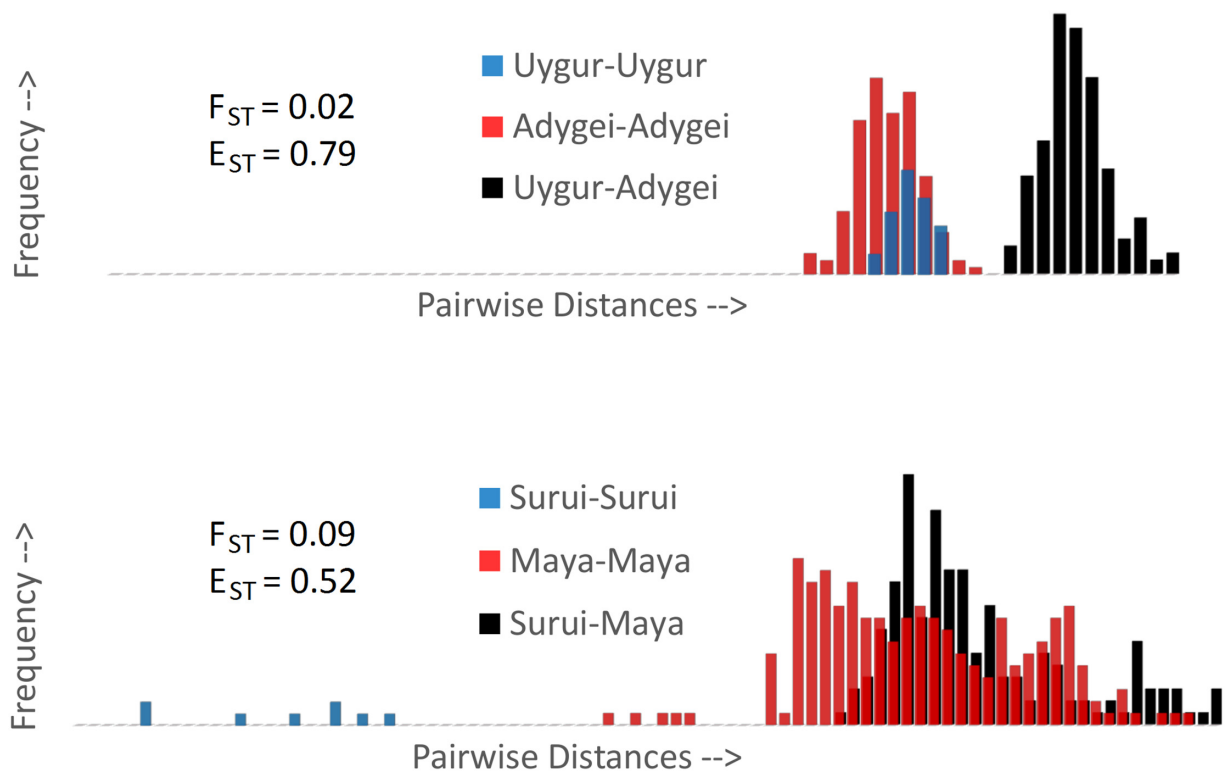


Fig 7. F_{ST} vs. genetic similarity in various population pairs. Pairwise distances are colored red or blue within populations and black between populations. (A) Even at a relatively low F_{ST} of 0.02 all within-population pairs among the Uygur and Adygei samples are genetically more similar than all the between-population pairs. (B) Separation is more ambiguous among Native Americans. Despite a relatively high F_{ST} of 0.09, there is substantial overlap between Maya-Maya (red) and Maya-Surui (black) samples. E_{ST} values are more consistent with the within- vs.-between population overlap and the dissimilarity fraction (ω).

doi:10.1371/journal.pone.0160413.g007

Summary and Conclusions

The core distinction between F_{ST} and E_{ST} is that F_{ST} partitions genetic diversity, whereas E_{ST} partitions genetic structure within and between populations. While F_{ST} is more sensitive to differences in within-population diversity, E_{ST} is more sensitive to outliers (though this is largely mitigated by using E_{ST} -median rather than E_{ST} -mean; see [Materials and Methods](#)). Since F_{ST} is weighed down by high levels of intrapopulation diversity, it can be close to zero even when population clusters are highly separated, however because it does not account for intrapopulation structure, high F_{ST} does not necessarily reflect highly separated population clusters. This is not necessarily a flaw in F_{ST} , but it does demonstrate a conceptual discrepancy between F_{ST} and strength of clustering.

Sewall Wright proposed a series of arbitrary F_{ST} thresholds ranging from 0.05 to 0.25, denoting little to very great differentiation [1]. Notably, the highest ranking of “very great differentiation” leaves most of the range (0.25–1) undefined. Given the wider empirical range of E_{ST} and its correspondence with results from clustering schemes (Fig 2), phylogeography (Fig 5), and the dissimilarity fraction (Fig 7), such arbitrary thresholds may not be necessary for E_{ST} . A value of E_{ST} larger than 0.5 simply indicates that most of the structure is between populations rather than within, corresponding to moderately separated populations such as Russians and Adygei ($E_{ST} = 0.5$), Bantu from South Africa and Kenya ($E_{ST} = 0.48$), or French and Sardinians ($E_{ST} = 0.48$) (S1 Table).

Differentiation metrics are judged by their ability to quantify meaningful evolutionary divergence, and can be indispensable in identifying *Evolutionarily Significant Units* (ESU) and *Distinct Population Segments* (DPS) for conservation [19]. For example, given several subpopulations within a species, it is reasonable to prioritize the most highly differentiated subpopulation for conservation in order to maximize biodiversity. However, higher F_{ST} does not necessarily reflect stronger separation and lower misclassification, as with the Uygur and Adygei, whose clusters are better defined than those of the Surui and Maya despite a fourfold lower F_{ST} (Fig 7). In this context humans can be a useful model species simply because we know so much about human populations due to our “long habit of observing ourselves” [20]. This allows us to make educated inferences about human populations that might otherwise be overlooked, e.g., we can be skeptical of the high Surui-Karitiana F_{ST} , and realize that this is most likely due to the relatively recent isolation of two small tribes. This is a luxury that we do not usually have with other species, in which case high F_{ST} can be misinterpreted as a deep phylogenetic divide, potentially leading to misguided conservation strategies. Our hope is that by combining information from both *fixation* (F_{ST}) and *equidistance* (E_{ST}) indices, researchers could make more informed decisions.

Unlike F_{ST} , which is typically averaged across any number of markers, E_{ST} is an asymptotic measure in the sense that it requires large datasets with many thousands of markers, which have only recently become widely accessible. With the latest SNP chips containing well over 100,000 markers, accurate estimates of departures from panmixia are finally within reach, and there is no longer a need for the simplifying assumption that subpopulations are effectively panmictic. By deriving an F_{ST} -type statistic for apportioning structure within and between populations, namely E_{ST} , we hope to add a new useful metric to the 21st century population genetics toolkit.

Materials and Methods

The HGDP data used in our analysis were accessed at: <http://www.hagsc.org/hgdp/files.html>. After removing the 163 mitochondrial SNPs and 105 samples previously inferred to be close relatives [18], the final file included 660,755 SNPs from 938 samples in 53 populations. Strings

of SNPs were treated as sequences, with mismatches summed and divided by the sequence length. Pairwise distances, based on Allele Sharing Distance (ASD) [21], were calculated as one minus half the average number of shared alleles per locus. The theoretical model, mathematical proofs and numerical simulations (using Mathematica v.8.0) of SD_T and SD_S appear in Appendix A.

In the empirical analysis we used Hudson’s pairwise-distance based F_{ST} estimator [4] adapted to diploid genotypes:

$$F_{ST} = 1 - \frac{S}{T} \tag{1}$$

where S and T are mean pairwise distances within subpopulations and in the total pooled population.

E_{ST} was formulated in terms of standard deviations as:

$$E_{ST} = 1 - \frac{SD_S}{SD_T} \tag{2}$$

where SD_S and SD_T are the standard deviations (SD) of pairwise distances within subpopulations and in the total population. This E_{ST} estimator is referred to as E_{ST} mean. We used three additional E_{ST} estimators: E_{ST} min, E_{ST} median, and E_{ST} max (S9 Fig). All four estimators use the same basic formula in Eq (2), with only the type of SD_S differing among estimators. In E_{ST} min, E_{ST} median, and E_{ST} max, SD_S is respectively replaced with the smallest, median, and largest individual SD, where the individual SD is the standard deviation of pairwise distances between a single sample and all other samples in the population. E_{ST} min uses the smallest individual SD_S from each population, i.e., the SD of the most panmictic sample, E_{ST} median uses the median individual SD_S , and E_{ST} max uses the highest individual SD_S . Each of these metrics has different sensitivities to various sampling biases. Due to E_{ST} mean’s sensitivity to the sampling of close relatives, we used E_{ST} median (which is unaffected by the inclusion of relatives as long as at least 50% of the samples are unrelated) as the primary measure of E_{ST} in this study. In the rare event that >50% of the samples are closely related, E_{ST} max may be preferable, as long as at least one individual has no close relatives among the samples. E_{ST} values, especially E_{ST} min and E_{ST} mean, can be negative if structure is high and differentiation is low (S9 Fig). Small sample sizes were often sufficient for estimating heterozygosity (S10 Fig) and F_{ST} and E_{ST} (S11 Fig) using all the SNPs in the HGDP dataset. Nevertheless, systematically developing estimators for E_{ST} is beyond the scope of the current treatment.

We derived an additional equidistance index, denoted E_{BT} , which is less sensitive to intra-population structure and the inclusion of relatives. Recall that E_{ST} reflects equidistance (E) within subpopulations (S) compared to the total (T) population. Similarly, E_{BT} reflects equidistance (E) between subpopulations (B) compared to the total (T) population:

$$E_{BT} = 1 - \frac{SD_B}{SD_T} \tag{3}$$

Where SD_B and SD_T are the standard deviations of pairwise distances between individuals from different subpopulations, and in the total pooled population respectively. In most cases $SD_T \geq SD_B$, because SD_T includes pairs of individuals from the same population as well as pairs from different populations, whereas SD_B only includes pairs of individuals from different populations. Pairs of individuals from the same population are likely to have a higher SD due to relatives in the samples, which disrupt the panmixia (e.g., in the Naxi population, see S2–S4 Figs). Panmictic populations are not just equidistant among themselves; they are also equidistant towards each other. Such populations should have similar SD_S and SD_B , and thus similar

E_{ST} and E_{BT} . Interestingly, some East Asians populations have relatively low E_{BT} , such as Cambodians vs. Mongola ($E_{BT} = 0.13$) and Japanese vs. Chinese ($E_{BT} = 0.16$). All F_{ST} , E_{ST} and E_{BT} estimates in this study are based on pairwise comparisons between two populations or population groups. Each of the two paired populations was given equal weight, as were the within- and between-population pairs. Thus, 25% of the total weight was given to each population, and 50% to between-population pairs.

We developed a custom MATLAB code for extracting genetic distances from SNP data and estimating heterozygosity, pairwise distances, F_{ST} , E_{ST} , and E_{BT} . The code corrects for missing data and small sample sizes, and identifies outliers, but includes no further assumptions or corrections. Phylogenetic trees and MDS plots were also generated with MATLAB. Equal angle and square neighbor-joining trees of individual similarities were generated from matrices of pairwise distances with the *seqneighjoin* command. An alternative script, based on the internal MATLAB *seqpdist* command for sequence distance, yielded similar results.

Appendix A

The standard deviation of pairwise distances as a measure of population structure

Our goal in this appendix is to substantiate the *asymptotic* (in terms of number of genetic markers) *standard deviation* of pairwise genetic distances as a good unsupervised measure of internal structure, thus justifiable as a basis for the definition of E_{ST} . In particular, we prove that this asymptotic standard deviation is zero if and only if there is no internal structure (i.e., the population is panmictic).

A model of pairwise genetic distances for genotypes from two diploid populations

Let p_i denote the frequency at locus i of allele 'A' in population 1, and let q_i denote the frequency of the same allele in population 2 and assume that both populations are effectively very large and have the same contribution to the total population. The commonly-used *allele sharing distance* (ASD) measures the dissimilarity of two individual genotypes. For *diploid* genotypes, it is commonly defined as 2 minus the number of shared alleles at each locus, averaged across loci [21–22]. For multiple loci genotypes we use a normalized (by the number of considered loci) version of ASD to simplify the analysis of means and variances of the ASD distribution, as in Tal (2013). Under the assumption of Hardy-Weinberg Equilibrium, allele frequencies fully determine per-locus genotype frequencies.

Let a categorical random variable X_i represent the ASD at diploid locus i , and let D_n represent the normalized ASD across n loci for pairs of genotypes sampled from the *total* population,

$$D_n = \frac{1}{n} \sum_{i=1}^n X_i$$

We are interested in arriving at an expression for the variance (and ultimately the asymptotic standard deviation) of D_n . Given *linkage equilibrium* (LE) within each of our two subpopulations (required in order to allow modeling LE in the total population), the X_i for the *total-population* pairs are *not* statistically independent, and therefore the formulation for the variance of D_n requires a partition into conditional expectations. From basic principles,

$$\text{Var}[D_n] = E[D_n^2] - E[D_n]^2$$

Now, to evaluate $E[D_n^2]$ we need to condition it upon classification of pairs of genotypes as within- or between-population since the X_i from the total population are not statistically independent,

$$\begin{aligned}
 E[D_n^2] &= E\left[\left(\frac{1}{n}\sum_{i=1}^n X_i\right)^2\right] = \frac{1}{n^2}E\left[\sum_{i=1}^n X_i^2 + 2\sum_{i<j}^n (X_i \cdot X_j)\right] = \frac{1}{n^2}\left(\sum_{i=1}^n E[X_i^2] + 2\sum_{i<j}^n E[X_i \cdot X_j]\right) \\
 &= \frac{1}{4} \cdot \frac{1}{n^2}\left(\sum_{i=1}^n E[X_i^2] + 2\sum_{i<j}^n E[X_i] \cdot E[X_j]\right) \Bigg| \text{both genotypes from pop 1} \\
 &+ \frac{1}{4} \cdot \frac{1}{n^2}\left(\sum_{i=1}^n E[X_i^2] + 2\sum_{i<j}^n E[X_i] \cdot E[X_j]\right) \Bigg| \text{both genotypes from pop 2} \\
 &+ \frac{1}{2} \cdot \frac{1}{n^2}\left(\sum_{i=1}^n E[X_i^2] + 2\sum_{i<j}^n E[X_i] \cdot E[X_j]\right) \Bigg| \text{one genotype from pop 1 and one from pop 2}
 \end{aligned} \tag{4}$$

where $E[X_i X_j] = E[x_i] \cdot E[X_j]$ since there is independence across any two loci for the within pairs and between pairs, and where the probabilities (assuming equal population sizes) for within-population 1 pairs, within-population 2 pairs, and between-population pairs are *at infinite population size* $\frac{1}{4}$, $\frac{1}{4}$, $\frac{1}{2}$ respectively (otherwise, for finite population sizes m we have probabilities $\frac{m-1}{4m-2}$, $\frac{m-1}{4m-2}$, $\frac{m}{2m-1}$ respectively).

Now, from per-locus probabilities in Tal (2013, [Eq 6](#) and [Table 1](#)) we derive the expected values,

$$\begin{aligned}
 E[X_i] \Big| \text{both genotypes from pop 1} &= 4(-p_i^4 + 2p_i^3 - 2p_i^2 + p_i) \\
 E[X_i^2] \Big| \text{both genotypes from pop 1} &= 4p_i(1 - p_i) \\
 E[X_i] \Big| \text{both genotypes from pop 2} &= 4(-q_i^4 + 2q_i^3 - 2q_i^2 + q_i) \\
 E[X_i^2] \Big| \text{both genotypes from pop 2} &= 4q_i(1 - q_i) \\
 E[X_i] \Big| \text{one genotype from pop 1 and one from pop 2} &= 2(2p_i q_i^2 + 2q_i p_i^2 - 2p_i^2 q_i^2 - 4p_i q_i + p_i + q_i) \\
 E[X_i^2] \Big| \text{one genotype from pop 1 and one from pop 2} &= 2(p_i^2 + q_i^2 - 4p_i q_i + p_i + q_i)
 \end{aligned} \tag{5}$$

So that,

$$\begin{aligned}
 E[D_n^2] &= \frac{1}{n^2}\sum_{i=1}^n p_i(1 - p_i) + 8\frac{1}{n^2}\sum_{i<j}^n (-p_i^4 + 2p_i^3 - 2p_i^2 + p_i)(-p_j^4 + 2p_j^3 - 2p_j^2 + p_j) \\
 &+ \frac{1}{n^2}\sum_{i=1}^n q_i(1 - q_i) + 8\frac{1}{n^2}\sum_{i<j}^n (-q_i^4 + 2q_i^3 - 2q_i^2 + q_i)(-q_j^4 + 2q_j^3 - 2q_j^2 + q_j) \\
 &+ \frac{1}{n^2}\sum_{i=1}^n (p_i^2 + q_i^2 - 4p_i q_i + p_i + q_i) \\
 &+ 4\frac{1}{n^2}\sum_{i<j}^n (2p_i q_i^2 + 2q_i p_i^2 - 2p_i^2 q_i^2 - 4p_i q_i + p_i + q_i)(2p_j q_j^2 + 2q_j p_j^2 - 2p_j^2 q_j^2 - 4p_j q_j + p_j + q_j)
 \end{aligned} \tag{6}$$

Also, since the expectation of a sum of dependent random variables is the sum of their expectations we have for the ‘total population’ X_i ,

$$E[D_n]^2 = \left(E \left[\frac{1}{n} \sum_{i=1}^n X_i \right] \right)^2 = \frac{1}{n^2} \left(\sum_{i=1}^n E[X_i] \right)^2$$

From Tal (2013, section 3.2) we have the expression for the pmf of X_i and thus can derive $E[X_i]$,

$$\begin{aligned} E[X_i] &= 0 \cdot Pr(X_i = 0) + 1 \cdot Pr(X_i = 1) + 2 \cdot Pr(X_i = 2) \\ &= \frac{1}{4}((p_i + q_i)^3(2 - p_i - q_i) + (2 - p_i - q_i)^3(p_i + q_i) + (p_i + q_i)^2(2 - p_i - q_i)^2) \\ &= \frac{1}{4}(p_i + q_i)(2 - p_i - q_i)((p_i + q_i)^2 + (2 - p_i - q_i)^2 + (p_i + q_i)(2 - p_i - q_i)) \\ &= \frac{1}{4}(p_i + q_i)(2 - p_i - q_i)(4 - (p_i + q_i)(2 - p_i - q_i)) \end{aligned}$$

Such that,

$$E[D_n]^2 = \frac{1}{16n^2} \left(\sum_{i=1}^n (p_i + q_i)(2 - p_i - q_i)(4 - (p_i + q_i)(2 - p_i - q_i)) \right)^2$$

So that finally,

$$\begin{aligned} Var[D_n] = E[D_n^2] - E[D_n]^2 &= \frac{1}{n^2} \left[\sum_{i=1}^n p_i(1 - p_i) + 8 \sum_{i<j}^n (-p_i^4 + 2p_i^3 - 2p_i^2 + p_i)(-p_j^4 + 2p_j^3 - 2p_j^2 + p_j) \right. \\ &+ \sum_{i=1}^n q_i(1 - q_i) + 8 \sum_{i<j}^n (-q_i^4 + 2q_i^3 - 2q_i^2 + q_i)(-q_j^4 + 2q_j^3 - 2q_j^2 + q_j) \\ &+ \sum_{i=1}^n (p_i^2 + q_i^2 - 4p_iq_i + p_i + q_i) \\ &+ 4 \sum_{i<j}^n (2p_iq_i^2 + 2q_i p_i^2 - 2p_i^2q_i^2 - 4p_iq_i + p_i + q_i)(2p_jq_j^2 + 2q_j p_j^2 - 2p_j^2q_j^2 - 4p_jq_j + p_j + q_j) \\ &\left. - \frac{1}{16} \left(\sum_{i=1}^n (p_i + q_i)(2 - p_i - q_i)(4 - (p_i + q_i)(2 - p_i - q_i)) \right)^2 \right] \end{aligned} \tag{7}$$

Thus we have an explicit formulation for the variance of the pairwise distance distribution of genotypes from two panmictic populations in terms of the allele frequencies across a given number of loci, n .

Crucially, we would like to prove that at the limit, the pairwise distance variance is asymptotically above zero *if and only if* the population has internal structure; i.e., if for any $F_{ST} > 0$,

$$S = \lim_{n \rightarrow \infty} Var[D_n] > 0$$

We will proceed by deriving an explicit expression for S . Consider an *equivalent setting* comprised of three random variables W, Y and Z , which represent the pairwise distances of genotypes within population 1, within population 2 and between populations 1 and 2, respectively. We sample n values X_i from just *one* of these distributions, by first flipping a 3-sided coin to decide from which: with a probability α for W , a probability β for Y and a probability γ for Z . Once the distribution has been selected, the sampling of X_i is done *i.i.d.* Note that due to

the randomized choice of the distribution from which to sample *all* the X_i , they are identically distributed but *not* independent. Now we set,

$$D_n = \frac{1}{n} \sum_{i=1}^n X_i$$

We would like to get an expression for S in terms of the expectations of W, Y, Z and α, β, γ , where,

$$S = \lim_{n \rightarrow \infty} \text{Var}(D_n)$$

From the law of total variance,

$$\text{Var}(D_n) = \text{Var}(E[D_n|B]) + E[\text{Var}(D_n|B)] = \text{Var}(U_{XYZ}) + \frac{1}{n} C \cdot \text{Var}(D_n)$$

where B is a categorical random variable that describes which of the distributions W, Y, Z we are sampling from, with probabilities α, β, γ respectively, where U_{XYZ} is a discrete random variable taking the values of $\mu_W = E[W], \mu_Y = E[Y], \mu_Z = E[Z]$ with corresponding probabilities α, β, γ respectively, and where C is some constant. Hence at the limit $n \rightarrow \infty$ we have,

$$S = \text{Var}(U_{XYZ}) = \alpha(\mu_W - \mu)^2 + \beta(\mu_Y - \mu)^2 + \gamma(\mu_Z - \mu)^2 \tag{8}$$

$$\mu = \alpha\mu_W + \beta\mu_Y + \gamma\mu_Z$$

and $S = 0$ if and only if the three means are equal, i.e., $\mu_W = \mu_Y = \mu_Z$.

Now consider three *sequences* of random variables $W_i, Y_i, Z_i, i:1 \dots n$, instead of the three single random variables, and sample n values from one of these sequences (again according to the prior probabilities α, β, γ). Once the sequence is selected, these samples are independent but now *not identically distributed*. We would again like to find S , and more importantly, the condition for which it is zero, this time in terms of $E[W_i], E[Y_i], E[Z_i]$ (and the prior probabilities). Sampling from a sequence with fixed probabilities just defines a new mixture distribution—so the problem gets reduced to the one already solved. Therefore U_{XYZ} is now defined by the three limits (since we have derived S in Eq (8) at the limit $n \rightarrow \infty$),

$$\mu_W = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n E[W_i], \mu_Y = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n E[Y_i], \mu_Z = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n E[Z_i] \tag{9}$$

with probabilities α, β, γ .

Crucially, this sampling scenario corresponds to our original setting of formulating the variance of the genetic distance of genotypes sampled from the total population, given the sequencing of an infinite number of loci, where μ_W and μ_Y in Eq (9) represent the two within-population pairwise distance means and μ_Z the total-population mean (derived below), and where the respective probabilities are as in Eq (4), $\alpha = 1/4, \beta = 1/4, \gamma = 1/2$, assuming *infinite population size*. Again, $S = 0$ if and only if these means are equal, i.e., $\mu_W = \mu_Y = \mu_Z$. Let us analyze the conditions for these equalities, given the corresponding formulations of the pairwise distance means. First, using the additivity of expectations,

$$E[D_n] = \frac{1}{n} \sum_{i=1}^n E[X_i]$$

we get from Eq (5) the expressions for any finite n ,

$$\begin{aligned}
 E[D_n] \mid \text{both genotypes from pop 1} &= \frac{1}{n} \sum_{i=1}^n 4(-p_i^4 + 2p_i^3 - 2p_i^2 + p_i) \\
 E[D_n] \mid \text{both genotypes from pop 2} &= \frac{1}{n} \sum_{i=1}^n 4(-q_i^4 + 2q_i^3 - 2q_i^2 + q_i) \\
 E[D_n] \mid \text{one genotype from pop 1 and one from pop 2} &= \frac{1}{n} \sum_{i=1}^n 2(2p_i q_i^2 + 2q_i p_i^2 - 2p_i^2 q_i^2 - 4p_i q_i + p_i + q_i)
 \end{aligned} \tag{10}$$

bearing in mind the analysis pertains to $E[D_n]$ as $n \rightarrow \infty$. We proceed to examine what can be concluded from the equalities $\mu_W = \mu_Y = \mu_Z$ (the only case where $S = 0$), given the means in Eq (10), about the allele frequencies p_i and q_i for any finite n (and this also holds at $n \rightarrow \infty$). Thus we start by explicitly writing the equalities (where the $1/n$ cancels out),

$$\begin{aligned}
 \sum_{i=1}^n (-p_i^4 + 2p_i^3 - 2p_i^2 + p_i) &= \sum_{i=1}^n (-q_i^4 + 2q_i^3 - 2q_i^2 + q_i) \\
 \sum_{i=1}^n (-p_i^4 + 2p_i^3 - 2p_i^2 + p_i) &= \sum_{i=1}^n (2p_i q_i^2 + 2q_i p_i^2 - 2p_i^2 q_i^2 - 4p_i q_i + p_i + q_i)
 \end{aligned} \tag{11}$$

To proceed we substitute new variables,

$$x_i = p_i(1 - p_i)$$

$$y_i = q_i(1 - q_i)$$

Then, the 1st equation in (8) becomes,

$$\sum_{i=1}^n x_i(1 - x_i) = \sum_{i=1}^n y_i(1 - y_i)$$

such that,

$$\sum_{i=1}^n [x_i(1 - x_i) + y_i(1 - y_i)] = \sum_{i=1}^n 2x_i(1 - x_i)$$

and using the 2nd equation in (8),

$$\begin{aligned}
 &= \sum_{i=1}^n [2p_i q_i (q_i + p_i - p_i q_i - 1) + p_i + q_i - 2p_i q_i] \\
 &= \sum_{i=1}^n [-2p_i q_i (1 - p_i)(1 - q_i) + (p_i - p_i^2) + (q_i - q_i^2) + (p_i^2 + q_i^2 - 2p_i q_i)]
 \end{aligned}$$

And again in terms of the new variables,

$$\sum_{i=1}^n [-2x_i y_i + x_i + y_i + (p_i - q_i)^2]$$

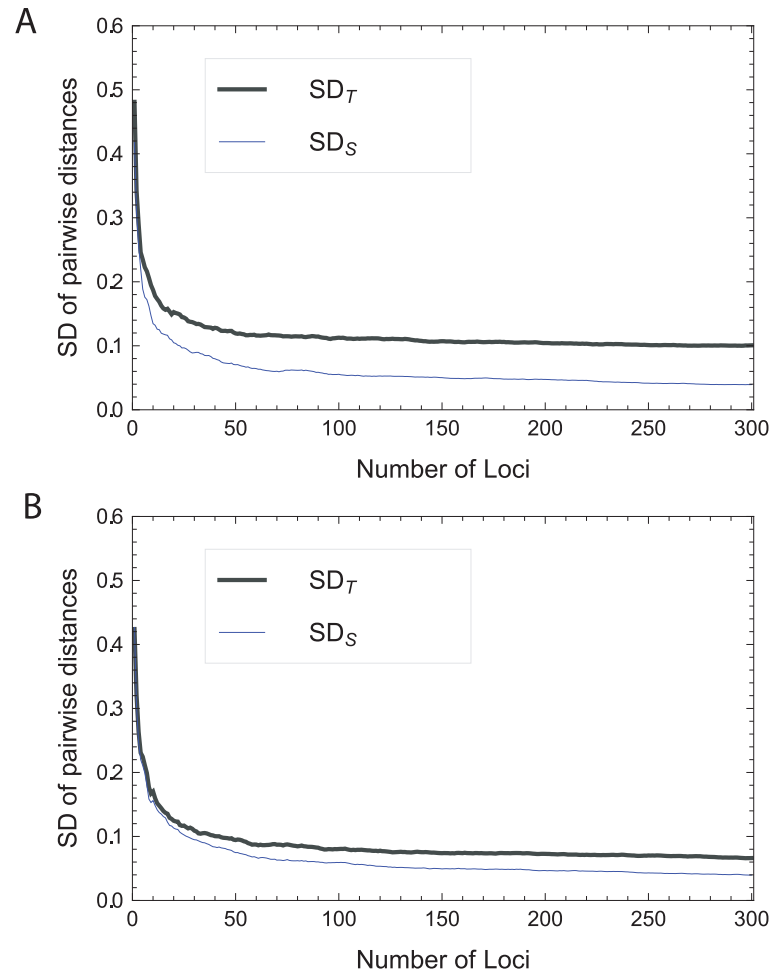


Fig 8. A simulation of SD_T and SD_S under a two panmictic population model demonstrating the divergent behavior of these two statistics with an increasing number of SNP loci. SNP frequencies are modeled on *Beta* distributions (as in [5]). (A) with $F_{ST} = 0.10$. (B) with $F_{ST} = 0.03$.

doi:10.1371/journal.pone.0160413.g008

This implies that,

$$\sum_{i=1}^n [(x_i - y_i)^2 + (p_i - q_i)^2] = 0$$

which occurs only if $p_i = q_i$ for all $i = 1, \dots, n$.

Therefore the *asymptotic* variance of the pairwise genetic distances (normalized by number of loci) of genotypes sampled from the combined population, comprising two subpopulations, is zero if this combined population is essentially a single panmictic population (i.e., $p_i = q_i$ for all i , or $F_{ST} = 0$). Since we have defined E_{ST} in terms of standard deviations rather than variances, we will subsequently consider the *asymptotic standard deviation* SD_T , which is simply defined as the square root of the *asymptotic variance*, S for the ‘total’ population. Fig 8 depicts numerical simulations of both SD_T and the average within-population SD (SD_S) for our two population model, as a function of the number of SNPs considered. While SD_S converges to zero, SD_T asymptotes to a value greater than zero, revealing the underlying structure. We note here that the rate of convergence to zero for SD_S is highly dependent on the diversity of each population—for lower diversity the within-population SD converges faster (and thus tends to be

lower for any finite number of loci). This factor also influences the rate of convergence of E_{ST} to its asymptotic value (see main text, Fig 3).

To further substantiate SD_T as a measure of structure, we would like to characterize the relation of SD_T to F_{ST} , both formulated as expressions of allele frequencies from two populations. We will proceed numerically, as our goal here is merely to get a qualitative intuition into the association of the two statistics.

We have from Eqs (8), (9) and (10), that asymptotically as $n \rightarrow \infty$, or practically under a high number of SNP loci,

$$SD_T = \sqrt{S} = \sqrt{\frac{1}{4}(\mu_W - \mu)^2 + \frac{1}{4}(\mu_Y - \mu)^2 + \frac{1}{2}(\mu_Z - \mu)^2} \tag{12}$$

where,

$$\mu = \frac{1}{4}\mu_W + \frac{1}{4}\mu_Y + \frac{1}{2}\mu_Z$$

$$\mu_W = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n 4(-p_i^4 + 2p_i^3 - 2p_i^2 + p_i)$$

$$\mu_Y = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n 4(-q_i^4 + 2q_i^3 - 2q_i^2 + q_i)$$

$$\mu_Z = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n 2(2p_i q_i^2 + 2q_i p_i^2 - 2p_i^2 q_i^2 - 4p_i q_i + p_i + q_i)$$

And from Tal (2013, Eq 13) we use the most common expression for F_{ST} across any number of n SNPs,

$$F_{ST} = \frac{\sum_{i=1}^n (p_i - q_i)^2}{\sum_{i=1}^n (p_i + q_i)(2 - p_i - q_i)} \tag{13}$$

Under the standard assumption that SNP frequencies are modeled on a *Beta* distribution with parameters deriving from some historical process (see [5, 22]) we sample a large number of sets of SNP frequencies for two subpopulations, each set generated from two *Beta* distributions with some randomized parameters. For each set we compute the pair SD_T (Eq 12) and F_{ST} (Eq 13) to generate a scatter plot of their association. Fig 9A–9D depicts several typical instances of such a simulation, demonstrating that the correlation of the two statistics is substantial when the total population is comprised of panmictic subpopulations,

$$0 \ll \rho[SD_T, F_{ST}] < 1$$

and is lower if the subpopulations have varying degrees of structure. Extending the numerical analysis for any number of *panmictic* multiple populations indicates that SD_T closely follows F_{ST} as a measure of structure.

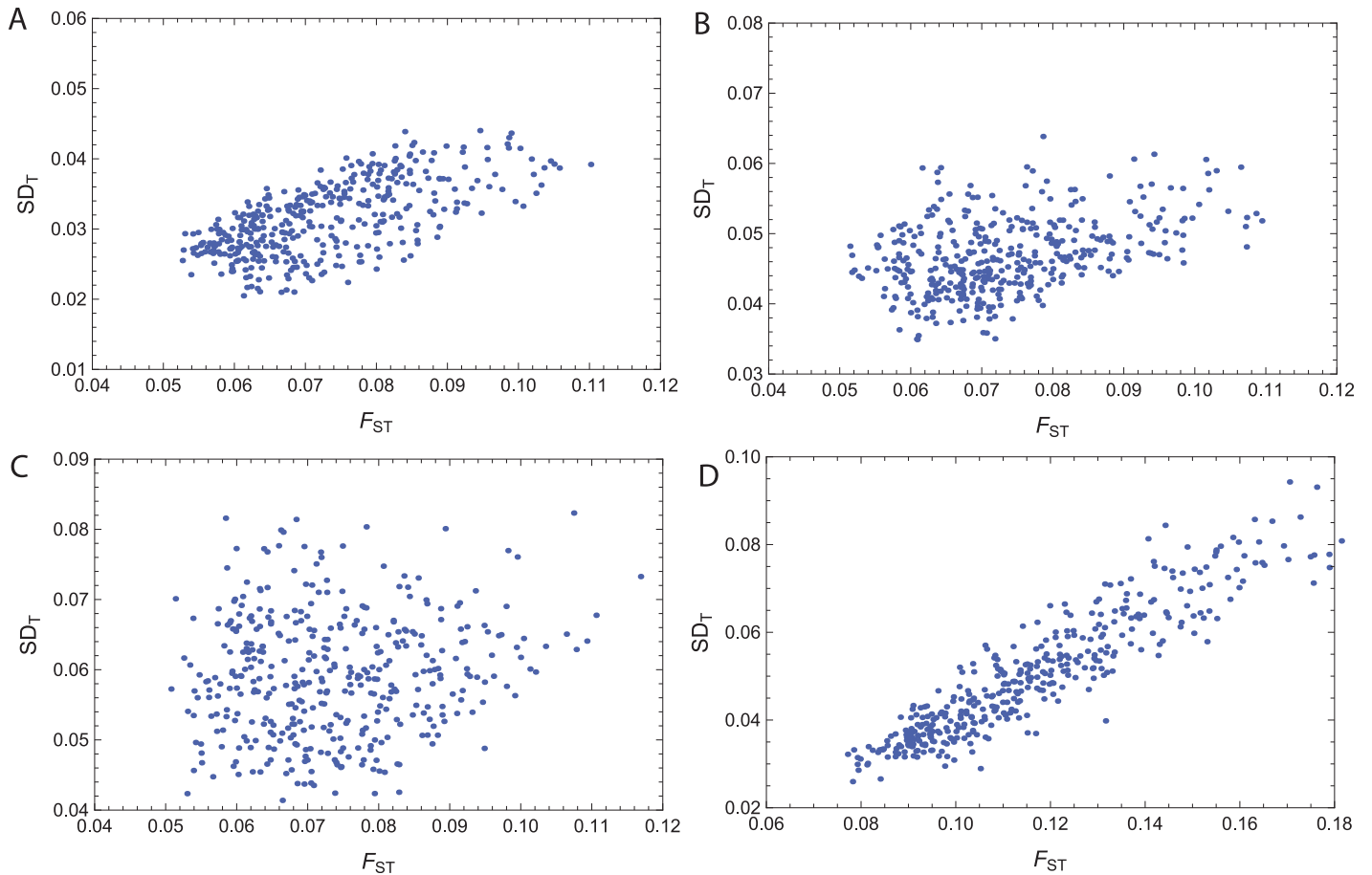


Fig 9. Scatter plots indicating a positive correlation for SD_T and F_{ST} . Each dot represents the two statistics computed for data sampled from our population model with 1000 SNPs and allele frequencies from Beta distributions. The Pearson product-moment correlation coefficient of F_{ST} and SD_T is 0.67 for plot (A) with panmictic populations, 0.38 for plot (B) with slightly varying structure in subpopulations, 0.14 for plot (C) with highly varying structure in subpopulations, and for 0.94 for plot (D) with three panmictic populations.

doi:10.1371/journal.pone.0160413.g009

We may generalize this result to more than two populations in a straightforward manner using induction. For instance, for three populations [Eq \(12\)](#) becomes,

$$SD_T = \sqrt{\frac{1}{9}(\mu_{W1} - \mu)^2 + \frac{1}{9}(\mu_{W2} - \mu)^2 + \frac{1}{9}(\mu_{W3} - \mu)^2 + \frac{2}{9}(\mu_{B1} - \mu)^2 + \frac{2}{9}(\mu_{B2} - \mu)^2 + \frac{2}{9}(\mu_{B3} - \mu)^2} \quad (14)$$

where,

$$\mu = \frac{1}{9}\mu_{W1} + \frac{1}{9}\mu_{W2} + \frac{1}{9}\mu_{W3} + \frac{2}{9}\mu_{B1} + \frac{2}{9}\mu_{B2} + \frac{2}{9}\mu_{B3}$$

$$\mu_{W1} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n 4(-p_i^4 + 2p_i^3 - 2p_i^2 + p_i)$$

$$\mu_{W2} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n 4(-q_i^4 + 2q_i^3 - 2q_i^2 + q_i)$$

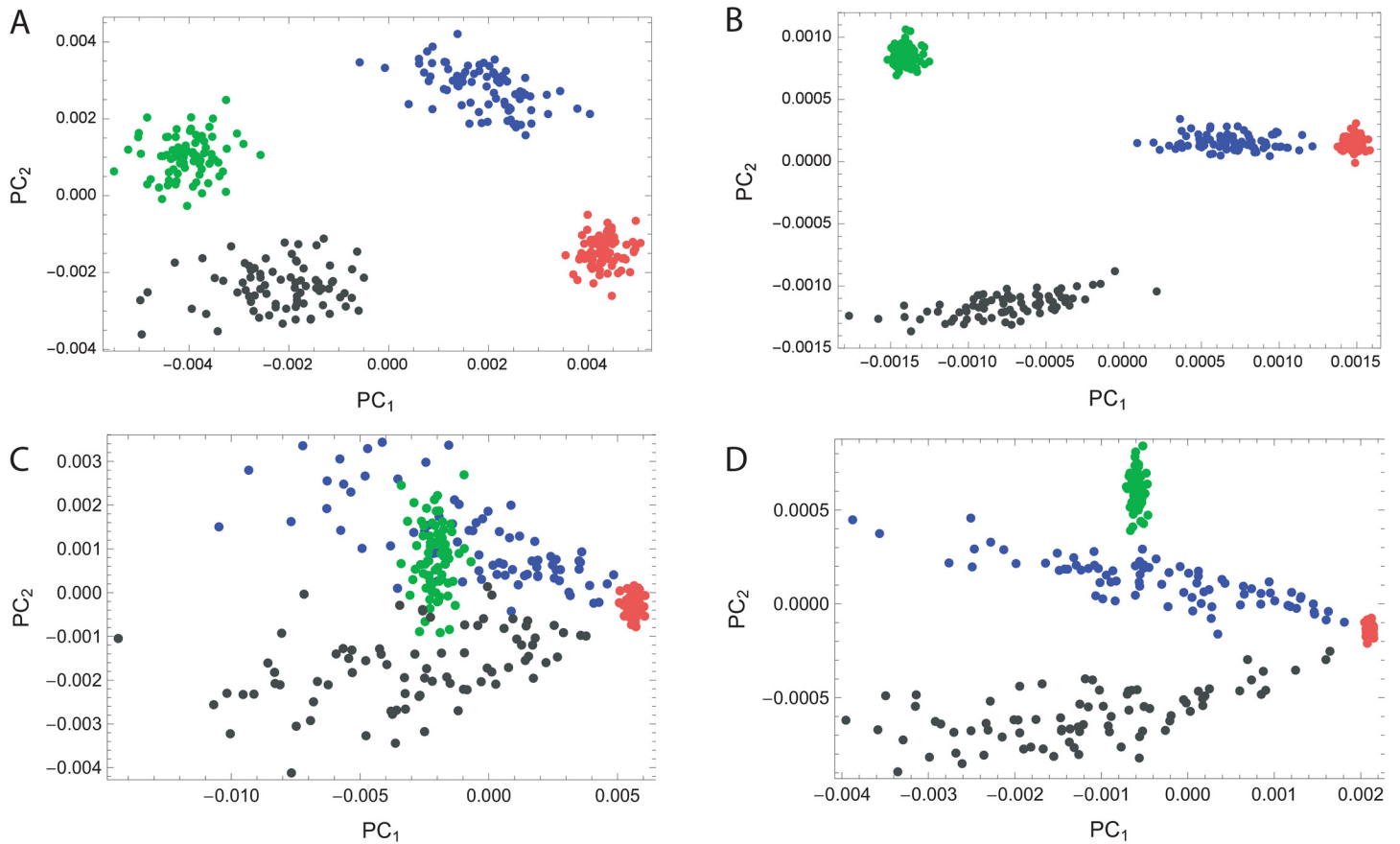


Fig 10. PCA plots from simulated SNP data of four populations (80 samples each) demonstrating the much pronounced decrease in SD_s for panmictic populations (red and green) relative to structured ones (black and blue), for two different patterns of internal structure, as the number of SNP loci processed by the PCA scheme is increased. (A-B) from 1K SNPs to 8K SNPs, where structure results from some random linkage disequilibrium (LD) pattern. (C-D) from 1K SNPs to 8K SNPs, where structure results from an LD pattern resembling admixture.

doi:10.1371/journal.pone.0160413.g010

$$\mu_{W3} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n 4(-r_i^4 + 2r_i^3 - 2r_i^2 + r_i)$$

$$\mu_{B1} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n 2(2p_i q_i^2 + 2q_i p_i^2 - 2p_i^2 q_i^2 - 4p_i q_i + p_i + q_i)$$

$$\mu_{B2} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n 2(2p_i r_i^2 + 2r_i p_i^2 - 2p_i^2 r_i^2 - 4p_i r_i + p_i + r_i)$$

$$\mu_{B3} = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n 2(2r_i q_i^2 + 2q_i r_i^2 - 2r_i^2 q_i^2 - 4r_i q_i + r_i + q_i)$$

and where F_{ST} for 3 panmictic populations with SNP loci is derived in the same manner as

Eq (13),

$$F_{ST} = \frac{\sum_{i=1}^n 2(p_i^2 + q_i^2 + r_i^2 - p_i q_i - p_i r_i - r_i q_i)}{\sum_{i=1}^n (p_i + q_i + r_i)(3 - p_i - q_i - r_i)} \quad (15)$$

A further perspective into SD_T as an unsupervised measure of internal structure is afforded by a qualitative comparison with principal component analysis (PCA) plots on data generated by the model. PCA is an *unsupervised* technique, essentially a dimensionality reduction procedure, used to emphasize the directions of greatest variation and bring out any strong patterns in a dataset. It can be used as a ‘preprocessing’ stage for clustering high-dimensional data, such as characteristic of population genetic samples. In such a setting, the first principal components tend to also extract existing substructure within the data in the form of clusters [14]. But more crucially to our goals, the relative dispersion of clusters on a PCA plot is highly associated with their internal structure, i.e., departures from panmixia, with increasing number of loci (and asymptotically, panmictic clusters would diminish to a single dot). This property is congruent with the convergence of SD_T to some value strictly greater than zero for non-panmictic populations. This is depicted in the four PCA plots of the same populations under increasing SNP count in Fig 10A–10D.

Finally, through numerical simulation of our model, we can see how varying degrees of internal structure (simulated by controlling the LD patterns) result in different asymptotic levels of E_{ST} (Fig 11). This serves to substantiate the empirical analysis depicted in Fig 3 of the main text.

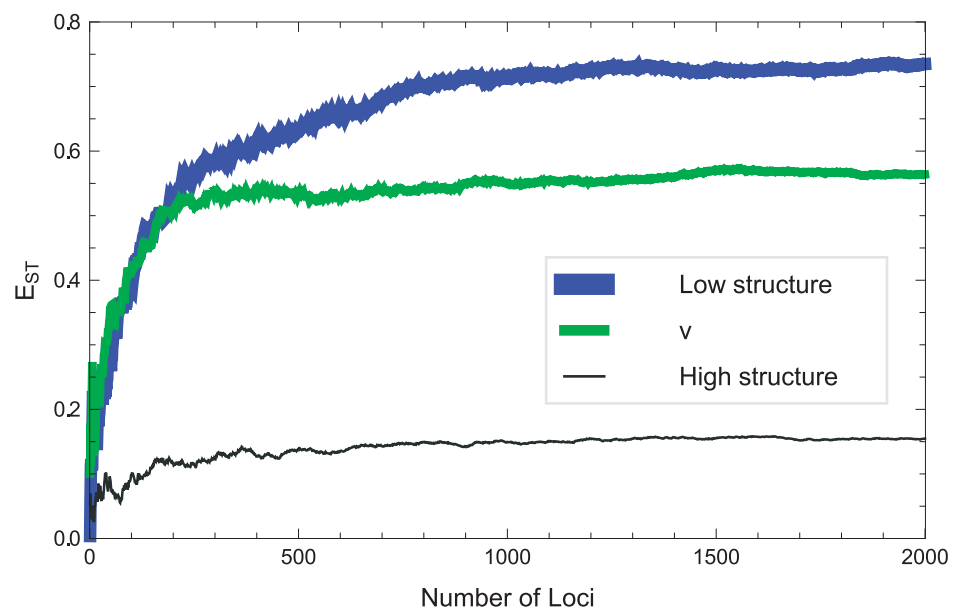


Fig 11. A numerical simulation of a model for E_{ST} for two structured populations (with $F_{ST} = 0.05$). E_{ST} was computed using the formulation in Eq (2) of Materials and Methods.

doi:10.1371/journal.pone.0160413.g011

Supporting Information

S1 Table. E_{ST} min-max, F_{ST} and E_{BT} in 60 HGDP population pairs ranked by E_{BT} .
(TIF)

S1 Fig. F_{ST} and E_{ST} within 12 local regions. Calculated from a single pair of populations per region: Israel (Bedouins vs. Druze), Italy (North Italians vs. Tuscans), Mexico (Maya vs. Pima), France (Basque vs. French), Brazil (Karitiana vs. Surui), Pakistan (Burusho vs. Kalash) East Asia (Cambodian vs. Mongola), Europe (Russians vs. Sardinians), Oceania (Melanesians vs. Papuans), Pygmies (Biaka vs. Mbuti), Russia (Russians vs. Yakut) and Southern Africans (South African Bantu vs. San). The most obvious discrepancy between F_{ST} and E_{ST} is in Brazil, with a high F_{ST} and moderate E_{ST} . The Druze and Bedouin of Israel live within a few hundred km of each other, speak the same language, and have the lowest E_{ST} among these 12 pairs, yet have a somewhat higher F_{ST} (several times higher than between the two Italian populations from Northern Italy and Tuscany and almost twice as high as between the French and Basques).
(TIF)

S2 Fig. Standard deviations (SD) of heterozygosity and pairwise genetic distances. From 660,755 SNPs in 53 HGDP populations. Excessive SD of genetic distance (blue) compared to SD of heterozygosity (red), as in the San and Naxi samples, implies the inclusion of relatives.
(TIF)

S3 Fig. Individual standard deviations in six HGDP populations. Each column represents the SD between a single individual and all other samples in the given population. Tuscans ($n = 7$), Italians ($n = 12$), Naxi ($n = 8$), Colombian ($n = 7$), Surui ($n = 8$), and Karitiana ($n = 13$). The “twin towers” in the Naxi batch are inferred to be a pair of close relatives in an otherwise panmictic population sample. These two individuals stick out like a sore thumb, while similarly related individuals are harder to identify among the Native American samples due to a higher base-level of structure in these population samples.
(TIF)

S4 Fig. The SD of pairwise distance plotted against the SD of heterozygosity. Generated from the entire HGDP dataset (938 individuals from 53 populations). The red diagonal line represents the linear trend line of the standard deviation of heterozygosity. Populations above this line are inferred to have more genetic structure than expected from heterozygosity, implying that relatives may have been included in the samples. Native American populations, highlighted in light blue, appear to have moderate or moderately high levels of relatives included among their samples.
(TIF)

S5 Fig. Pairwise F_{ST} and E_{ST} vs. geographic distance from the two Amazonian tribes to various global HGDP populations with increasing distance from the Amazon.
(TIF)

S6 Fig. Neighbor-Joining trees of individual similarities. Generated from 660,755 SNPs. Individual branches are black, inter-population branches are red, and intra-population branches are blue. A. Complete trees. B. Zoom into trees with individual branches (black) removed.
(TIF)

S7 Fig. Pairwise population distance charts. Each sample is represented by a red or blue string and each point on each string reflects distance between a pair of samples. Points that fall far below the rest are inferred to reflect close relatives.
(TIF)

S8 Fig. Superimposed distance plots of Uygur and Adygei (top) and Surui and Maya (bottom). This is the same kind of plot as in [S7 Fig](#), with each string representing a single individual. Despite a high F_{ST} of 0.09 ($E_{ST} = 0.52$), some Mayan individuals (red) are genetically closer to some Surui individuals (blue) than to some fellow Mayan individuals ($\omega > 0$), presumably due to outbreeding (some Mayan individuals have significant European admixture, which increases distances among Mayans). There is no such overlap between Uygur (yellow) and Adygei (black) samples ($\omega = 0$) despite a much lower pairwise F_{ST} of 0.02 ($E_{ST} = 0.79$). (TIF)

S9 Fig. Overview of F_{ST} , E_{ST} and E_{BT} among 60 HGDP population pairs (660,755 SNPs). Negative E_{STmin} (yellow) and E_{STmean} (orange) would imply that close relatives were included among these samples. Of the 60 population pairs in the analysis, 12 (20%) have negative E_{STmin} values and 6 have negative E_{STmean} values. $E_{STmedian}$, E_{STmax} , and E_{BT} cover virtually the entire 0–1 range with no negative values in these samples. The general trend is $F_{ST} < E_{STmin} < E_{STmean} < E_{STmedian} < E_{STmax}$. E_{BT} (gray) is usually somewhere between $E_{STmedian}$ (red) and E_{STmax} (black). (TIF)

S10 Fig. Mean heterozygosity as a function of sample size. Heterozygosity in various HGDP populations with sample size increasing from 1 to 15. All samples were included in populations with less than 15 samples (namely 7 in Colombians, 8 in Surui, and 13 in Karitiana). (TIF)

S11 Fig. Pairwise F_{ST} , E_{ST} and E_{BT} as a function of sample size. Differentiation was estimated in two population pairs: French-Japanese and Surui-Karitiana, with population sample sizes ranging from $n = 2$ to $n = 8$. French-Japanese estimates were also taken at $n = 15$ and $n = 28$ due to their larger samples. F_{ST} and E_{BT} start at $n = 2$; E_{ST} starts at $n = 3$, the minimal sample size for estimating the standard deviation of pairwise distances. (TIF)

Acknowledgments

We thank Alan Templeton for helpful advice, Tat Dat Tran for assisting with one of the proofs in the appendix, Lior Lesch for software support, Sagi Abelson for help with the MATLAB script and two anonymous reviewers. KS wishes to acknowledge the Israel Science Foundation (grant 189/05) and the Beutler Fund at Rambam Medical Center for research support.

Author Contributions

Conceived and designed the experiments: YG.

Analyzed the data: YG OT.

Wrote the paper: YG.

Provided assessment and correction to the mathematics and crucial corrections to the interpretation and manuscript: SR. Refinement of the concept: KS. Guidance in relevant population genetics applications: KS. Revision of manuscript: KS.

References

1. Wright S (1978) Evolution and the Genetics of Populations. Vol. 4, Variability Within and Among Natural Populations. University of Chicago Press, Chicago.

2. Nei M (1973) Analysis of gene diversity in subdivided populations. *Proc Natl Acad Sci USA*, 70, 3321–3323. PMID: [4519626](#)
3. Weir BS, Cockerham CC (1984) Estimating F-Statistics for the Analysis of Population Structure. *Evolution* 38 (6): 1358.
4. Hudson RR, Slatkin M, Maddison WP (1992) Estimation of levels of gene flow from DNA sequence data. *Genetics* 132 (2), 583–589. PMID: [1427045](#)
5. Tal O (2013) Two complementary perspectives on inter-individual genetic distance. *Biosystems*, 111: 18–36. doi: [10.1016/j.biosystems.2012.07.005](#) PMID: [22898797](#)
6. Hedrick PW (2005) A standardized genetic differentiation measure. *Evolution*, 59, 1633–1638. PMID: [16329237](#)
7. Jost L (2008) G_{ST} and its relatives do not measure differentiation. *Mol Ecol*, 17, 4015–4026. PMID: [19238703](#)
8. Whitlock M (2011) $G'st$ and D do not replace F_{st} . *Mol Ecol* 20:1083–109. doi: [10.1111/j.1365-294X.2010.04996.x](#) PMID: [21375616](#)
9. Lewontin RC (1972) The apportionment of human diversity. *Evol Biol*, 6, 381–98.
10. Witherspoon DJ, Wooding S, Rogers AR, Marchani EE, Watkins WS, Batzer MA, et al. (2007) Genetic Similarities Within and Between Human Populations. *Genetics* 176(1): 351–9. PMID: [17339205](#) doi: [10.1534/genetics.106.067355](#)
11. Edwards AWF (2003) Human genetic diversity: Lewontin's fallacy. *Bioessays*, 25, 798–801. PMID: [12879450](#)
12. Tal O (2012). The cumulative effect of genetic markers on classification performance: insights from simple models. *J. Theor. Biol.* 293 (January), 206–218.
13. Nelis M, Esko T, Mägi R, Zimprich F, Zimprich A, Toncheva D, et al. (2009) Genetic structure of Europeans: a view from the North-East. *PLoS One* 4: e5472. doi: [10.1371/journal.pone.0005472](#) PMID: [19424496](#)
14. Patterson N, Price AL, Reich D (2006) Population Structure and Eigenanalysis. *PLoS Genet* 2: 2074–2093. doi: [10.1371/journal.pgen.0020190](#)
15. Rosenberg NA, Pritchard JK, Weber JL, Cann HM, Kidd KK, e Zhivotovsky LA, et al. (2002) Genetic structure of human populations. *Science* 298: 2381–2385. PMID: [12493913](#)
16. Cann HM, de Toma C, Cazes L, Legrand MF, Morel V, Piouffre L (2002) A human genome diversity cell line panel. *Science* 296 (5566): 261–2. PMID: [11954565](#)
17. Rosenberg NA (2006) Standardized subsets of the HGDP-CEPH Human Genome Diversity Cell Line Panel, accounting for atypical and duplicated samples and pairs of close relatives. *Ann Hum Genet*, 70: 841–847. PMID: [17044859](#)
18. Jorde LB, Wooding SP (2004) Genetic variation, classification, and “race”. *Nat Genet*, 36, S28–32. PMID: [15508000](#)
19. Waples RS (1991) Definition of 'species' under the Endangered Species Act: Application to Pacific salmon. U.S. Department of Commerce NOAA Technical Memorandum, NMFS, F/NWC–194.
20. Darwin C (1871) *The Descent of Man and Selection in Relation to Sex*. London: John Murray.
21. Gao X, Martin ER (2009) Using allele sharing distance for detecting human population stratification. *Hum Hered*, 68, 182–191. doi: [10.1159/000224638](#) PMID: [19521100](#)
22. Nakamura T, Shoji A, Fujisawa H, Kamatani N (2005) Cluster analysis and association study of structured multilocus genotype data. *J Hum Genet* 50: 53–61. doi: [10.1007/s10038-004-0220-x](#) PMID: [15696377](#)