# Signatures of Archaic Adaptive Introgression in Present-Day Human Populations

Fernando Racimo,[‡,1] Davide Marnetto,[2] and Emilia Huerta-Sánchez[*,3]

[1]Department of Integrative Biology, University of California Berkeley, Berkeley, CA

[2]Department of Molecular Biotechnology and Health Sciences, University of Torino, Turin, Italy

[3]School of Natural Sciences, University of California Merced, Merced, CA

[‡]Present address: New York Genome Center, New York, NY

[*]**Corresponding author**: E-mail: ehuerta-sanchez@ucmerced.edu

**Associate editor**: John Novembre

## Abstract

**Comparisons of DNA from archaic and modern humans show that these groups interbred, and in some cases received an evolutionary advantage from doing so. This process—adaptive introgression—may lead to a faster rate of adaptation than is predicted from models with mutation and selection alone. Within the last couple of years, a series of studies have identified regions of the genome that are likely examples of adaptive introgression. In many cases, once a region was ascertained as being introgressed, commonly used statistics based on both haplotype as well as allele frequency information were employed to test for positive selection. Introgression by itself, however, changes both the haplotype structure and the distribution of allele frequencies, thus confounding traditional tests for detecting positive selection. Therefore, patterns generated by introgression alone may lead to false inferences of positive selection. Here we explore models involving both introgression and positive selection to investigate the behavior of various statistics under adaptive introgression. In particular, we find that the number and allelic frequencies of sites that are uniquely shared between archaic humans and specific present-day populations are particularly useful for detecting adaptive introgression. We then examine the 1000 Genomes dataset to characterize the landscape of uniquely shared archaic alleles in human populations. Finally, we identify regions that were likely subject to adaptive introgression and discuss some of the most promising candidate genes located in these regions.**

*Key words*: neanderthal, denisova, adaptive introgression, ancient DNA.

## Introduction

There is now a large body of evidence supporting the idea that certain modern human populations admixed with archaic groups of humans after expanding out of Africa. In particular, non-African populations have 1–2% Neanderthal ancestry (Green et al. 2010; Prüfer et al. 2014), and Melanesians and East Asians have 3% and 0.2% ancestry, respectively, from Denisovans (Reich et al. 2010; Meyer et al. 2012; Prüfer et al. 2014).

Recently, it has become possible to identify the fragments of the human genome that were introgressed and survive in present-day individuals (Prüfer et al. 2014; Sankararaman 2014; Vernot and Akey 2014; Sankararaman et al. 2016; Vernot et al. 2016). Researchers have also detected which of these introgressed regions are present at high frequencies in certain present-day non-African populations. Some of these regions are likely to have undergone positive selection in those populations after they were introgressed, a phenomenon known as adaptive introgression (AI). One particularly striking example of AI is the gene *EPAS1* (Hu et al. 2003) which confers a selective advantage in Tibetans by making them less prone to hypoxia at high altitudes (Beall et al. 2010; Bigham et al. 2010; Yi et al. 2010; Peng et al. 2011; Wang et al. 2011; Xu et al. 2011; Jeong et al. 2014; Hackinger et al. 2016). The selected Tibetan haplotype is likely to have been introduced in the human gene pool by Denisovans or a population closely related to them (Huerta-Sánchez et al. 2014; Huerta-Sanchez and Casey 2015).

In this study, we first use simulations to assess the power to detect AI using different exploratory summary statistics that do not require the introgressed fragments to be identified *a priori*. Some of these are inspired by the signatures observed in *EPAS1*, which contains an elevated number of sites with alleles uniquely shared between the Denisovan genome and Tibetans. We then apply these statistics to real human genomic data from phase 3 of the 1000 Genomes Project (Auton et al. 2015), to detect AI in human populations, and find candidate genes. While these statistics are sensitive to adaptive introgression, they may also be sensitive to other phenomena that generate genomic patterns similar to those generated by AI, like ancestral population structure and incomplete lineage sorting. These processes, however, should not generate long regions of the genome where haplotypes from the source and the recipient population are highly similar. As additional confirmation that the candidates we found with our statistics are generated by AI, we explored the haplotype structure of some of the most promising candidates,

**Open Access**

**Article**

**Table 1.** Summary statistics mentioned in the main text.

| Statistic | Explanation | References |
|---|---|---|
| $D$ | D-statistic: measures excess allele sharing between a test population and an outgroup using a sister population that is more closely related to the test than the outgroup | (Green et al. 2010; Durand et al. 2011) |
| $f_D$ | Similar to the D-statistic, but serves to better control for local variation in diversity patterns if one is interested in finding loci with excess ancestry from an admixing population. | (Martin et al. 2015) |
| $R_D$ | Average ratio of the sequence divergence between an individual from the source population and an individual from the admixed population, and the sequence divergence between an individual from the source population and an individual from the non-admixed population. This is computed by taking the average over all pairs of admixed and non-admixed individuals. | This study |
| $U_{A,B,C}(w,x,y)$ | Number of sites in which any allele is at a frequency lower than $w$ in panel A, higher than $x$ in panel B and equal to $y$ in panel C. | This study |
| $U_{A,B,C,D}(w,x,y,z)$ | Number of sites in which any allele is at a frequency lower than $w$ in panel A, higher than $x$ in panel B, equal to $y$ in panel C and equal to $z$ in panel D. | This study |
| $Q95_{A,B,C}(w,y)$ | 95% quantile of the distribution of derived allele frequencies in panel B, for sites where the derived allele is at a frequency lower than $w$ in panel A and equal to $y$ in panel C. | This study |
| $Q95_{A,B,C,D}(w,y,z)$ | 95% quantile of the distribution of derived allele frequencies in panel B, for sites where the derived allele is at a frequency lower than $w$ in panel A, equal to $y$ in panel C and equal to $z$ in panel D. | This study |
| $\pi$ | Expected heterozygosity, measured as the average of $2p(1-p)$ over all sites in a window, where $p$ is the frequency of an arbitrarily chosen allele. | (Crow et al. 1970) |
| $D'[intro]$ | A measure of linkage disequilibrium. Computed as $D/D_{max}$ where $D = p_{XY} - p_X p_Y$, $p_{XY}$ is the frequency of haplotype $XY$, $p_X$ is the frequency of allele $X$, $p_Y$ is the frequency of allele $Y$, and $D_{max}$ is the maximum theoretical value that D can take. $D'[intro]$ is computed only using frequencies from the introgressed panel. When we compute this in a window, the value is obtained by taking the average over all pairs of SNPs. | (Lewontin 1964) |
| $D'[comb]$ | $D'$ computed using haplotype and allele frequencies from the combination of the introgressed and non-introgressed panels. | (Lewontin 1964) |
| $r^2[intro]$ | A measure of linkage disequilibrium. Computed as $D^2/(p_X(1-p_X)p_Y(1-p_Y))$. $r^2[intro]$ is computed only using frequencies from the introgressed panel. When we compute this in a window, the value is obtained by taking the average over all pairs of SNPs. | (Hill and Robertson 1968) |
| $r^2[comb]$ | $r^2$ computed using haplotype and allele frequencies from the combination of the introgressed and non-introgressed panels. | (Hill and Robertson 1968) |

and used a probabilistic method (Seguin-Orlando et al. 2014) that infers introgressed segments along the genome by looking at the spatial arrangement of SNPs that are consistent with introgression. This allows us to verify that the candidate regions contain introgressed haplotypes at high frequencies: a hallmark of AI.
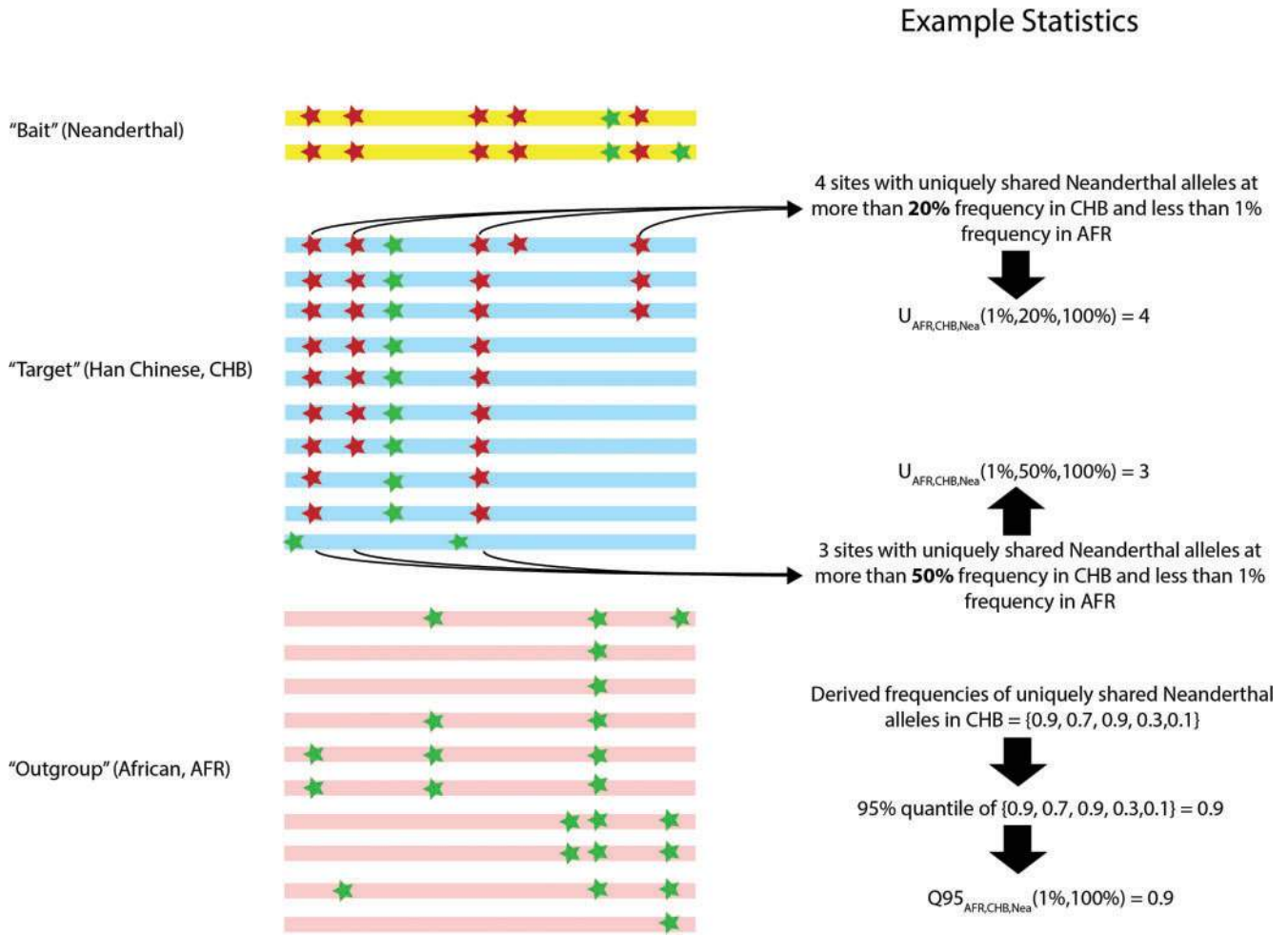
## Results

### Statistics for Detecting AI

We began by evaluating the performance of various statistics at detecting AI. In addition to testing statistics that have already been previously defined in the literature ($D$, $f_D$, $D'$, $r^2$, $\pi$), we define three new types of statistics that we find are particularly powerful at detecting AI (table 1). We briefly describe the new statistics here, but more extensive descriptions of all tested statistics can be found in the "Methods" section below.

First, in a region under AI, one would expect the sequence divergence between an individual from the source population and an admixed individual to be smaller than the sequence divergence between an individual from the source population and a nonadmixed individual. Thus, we define $R_D$ as the average ratio of the sequence divergence between a panel from the source population and a panel from the admixed population, over the sequence divergence between the source panel and a nonadmixed panel, computed in a window of the genome.

Second, under AI we would expect a large number of sites containing archaic alleles at high frequency in the admixed population, but absent or at low frequency in a nonadmixed population. Therefore, we define the statistic $U_{A,B,C}(w,x,y)$ to be equal to the number of sites within a genomic window where a sample C from an archaic source population has a particular allele at frequency $y$, and that allele is at a frequency smaller than $w$ in a panel A of a nonadmixed population but larger than $x$ in a panel B of an admixed population (fig. 1). Throughout the text, we denote panels A, B and C as the "outgroup", "target", and "bait" panels, respectively. If we have samples from two different archaic populations (e.g., a Neanderthal genome and a Denisova genome), we can define $U_{A,B,C,D}(w,x,y,z)$ as the number of sites where the archaic sample C has a particular allele at frequency $y$ and the archaic sample D has that allele at frequency $z$. Additionally, at those sites, the same allele should be at a frequency smaller than $w$ in an outgroup panel A and larger than $x$ in a target panel B (supplementary fig. S1, Supplementary Material online).

Finally, if we do not want to set a hard cutoff for what we consider "high-frequency" archaic alleles, we can just compute a summary statistic of the site-frequency spectrum in the target panel, conditional on the archaic allele being at low frequency in the outgroup. This statistic should be high when a region contains many alleles at especially high frequencies in the target. We therefore define $Q95_{A,B,C}(w,y)$ to be equal to the 95th percentile of derived frequencies in a target panel B

**FIG. 1.** Schematic illustration of the way the $U_{A,B,C}$ and $Q95_{A,B,C}$ statistics are calculated.

of all SNPs that have a derived allele frequency $y$ in the bait (archaic) panel C, but where the derived allele is at a frequency smaller than $w$ in an outgroup panel A from a nonadmixed population (fig. 1).
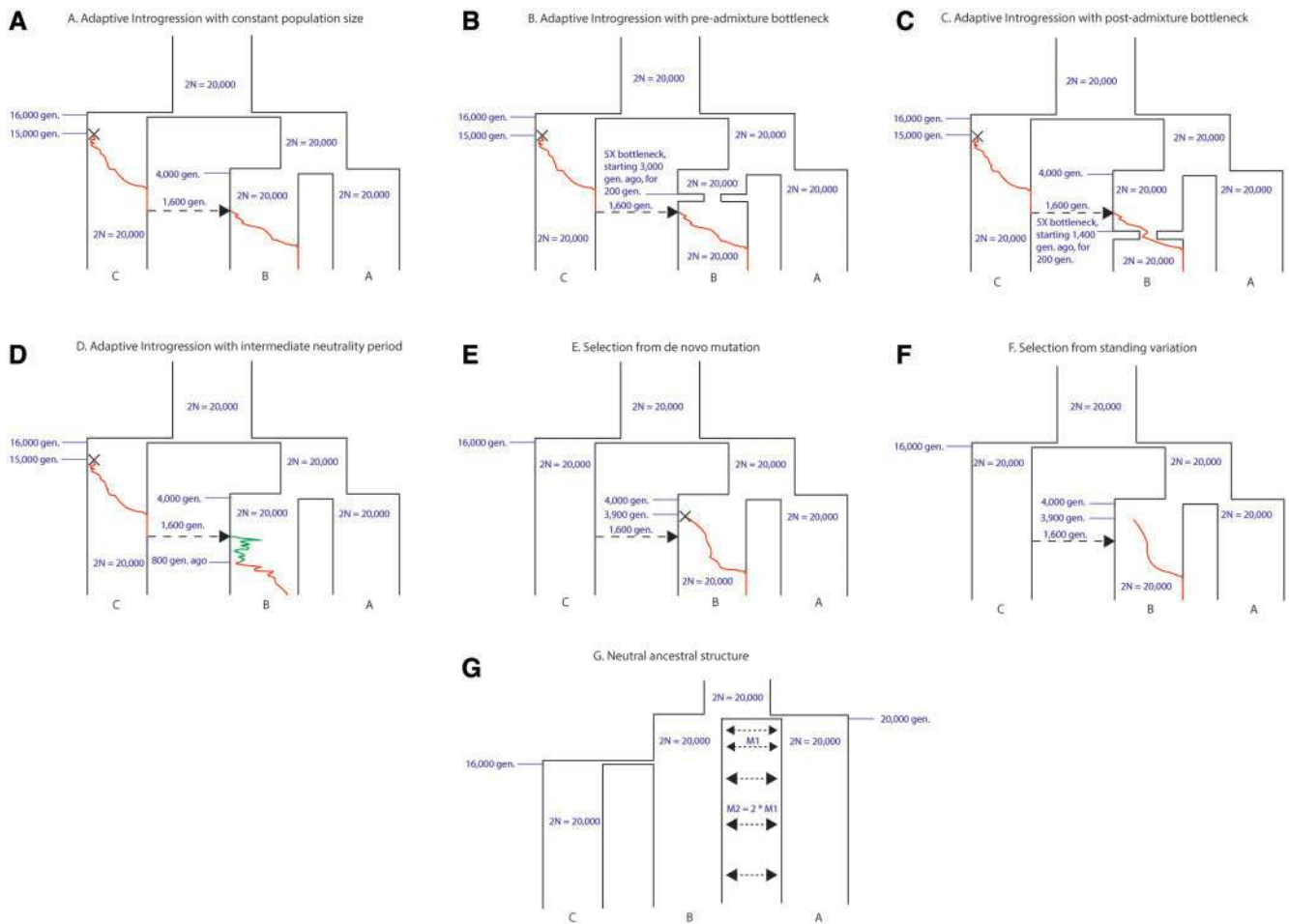
### Simulations under AI

We use simulations to assess the performance of the statistics mentioned above at detecting AI. Supplementary figures S3–S5, Supplementary Material online show the distribution of statistics that rely on patterns of shared allele configurations between source and introgressed populations ($\pi$, $D$, $f_D$, $U_{A,B,C}$, $Q95_{A,B,C}$, and $R_D$), for different choices of the selection coefficient $s$, and under 2%, 10%, and 25% admixture rates, respectively. For $Q95_{A,B,C}(w, 100\%)$ and $U_{A,B,C}(w, x, 100\%)$, we tested different choices of the outgroup cutoff $w$ (1%, 10%) and the target cutoff $x$ (0%, 20%, 50%, and 80%).

Some statistics, like $Q95_{A,B,C}(1\%, 100\%)$ and $f_D$ show strong separation between the selection regimes. For example, with an admixture rate of 2%, $Q95_{A,B,C}(1\%, 100\%)$ has 100% sensitivity at a specificity of 99%, for both $s = 0.1$ and $s = 0.01$. Some parameterizations of the $U$ statistic are not as effective, however. For example, $U_{A,B,C}(1\%, 0\%, 100\%)$ shows some power when the admixture rate is low (2%), but almost no power when the admixture rate is high (25%). This is because

setting the minimum frequency of the archaic allele in the test population at $x = 0\%$ means that any site with some archaic allele in the test panel will be counted, regardless of the allele frequency, so long as the archaic allele is at low frequency in the outgroup panel. At high admixture rates, low- and medium-frequency archaic alleles would naturally occur under neutrality, so they would not be informative about AI.

We also evaluated the effectiveness of LD-based statistics at detecting AI (supplementary fig. S6, Supplementary Material online). We tested the performance of $D'$ and $r^2$ by either computing each of these in the admixed panel only ($D'[intro]$, $r^2[intro]$) or in the combination of the admixed and nonadmixed panels ($D'[comb]$, $r^2[comb]$). Whereas $D'[intro]$, $D'[comb]$ and $r^2[comb]$ are modestly increased by AI, this is not the case with $r^2[intro]$ under strong selection and admixture regimes. This is because $r^2$ will tend to decrease if the minor allele frequency is very small, which will occur if this frequency is only measured in the population undergoing AI. In general, these statistics are not as powerful for detecting AI as allele configuration statistics like $U$ or $Q95$.

To jointly explore the power and specificity of all these statistics, we generated receiving operating characteristic (ROC) curves under various selection and admixture regimes (fig. 3 and supplementary fig. S7, Supplementary Material

**FIG. 2.** Demographic models described in the main text.

online). In general, $Q95_{A,B,C}(1\%, 100\%)$, $Q95_{A,B,C}(10\%, 100\%)$, and $f_D$ are very powerful statistics for detecting AI under strong ($s = 0.1$) and intermediate ($s = 0.01$) selection pressures. The number of uniquely shared sites $U_{A,B,C}(w, x, y)$ is also powerful, so long as the population in the target panel (B) is large. Additionally, for different choices of $x$, using $w = 1\%$ yields a more powerful statistic than using $w = 10\%$. We also tested AI scenarios with weak selection ($s = 0.001$), in which all statistics performed rather poorly, with $Q95$ and $f_D$ performing comparably better than the rest (supplementary fig. S8, Supplementary Material online). However, under these conditions, it is very unlikely that the selected allele will reach appreciable frequencies (supplementary fig. S2, Supplementary Material online), so the lack of sensitivity of all statistics here is largely a reflection of the fact that in most simulations the selected allele is not successful, especially when the probability of admixture is low. Conditioning on survival of the allele should therefore increase sensitivity.

We were also interested in the joint distribution of pairs of these statistics. Supplementary figure S9, Supplementary Material online shows the joint distribution of $Q95_{A,B,C}(1\%, 100\%)$ in the $y$-axis and four other statistics ($R_D$, $\pi$, $D$, and $f_D$) in the $x$-axis, under different admixture and selection regimes. One can observe, for example, that whereas $Q95_{A,B,C}(1\%, 100\%)$
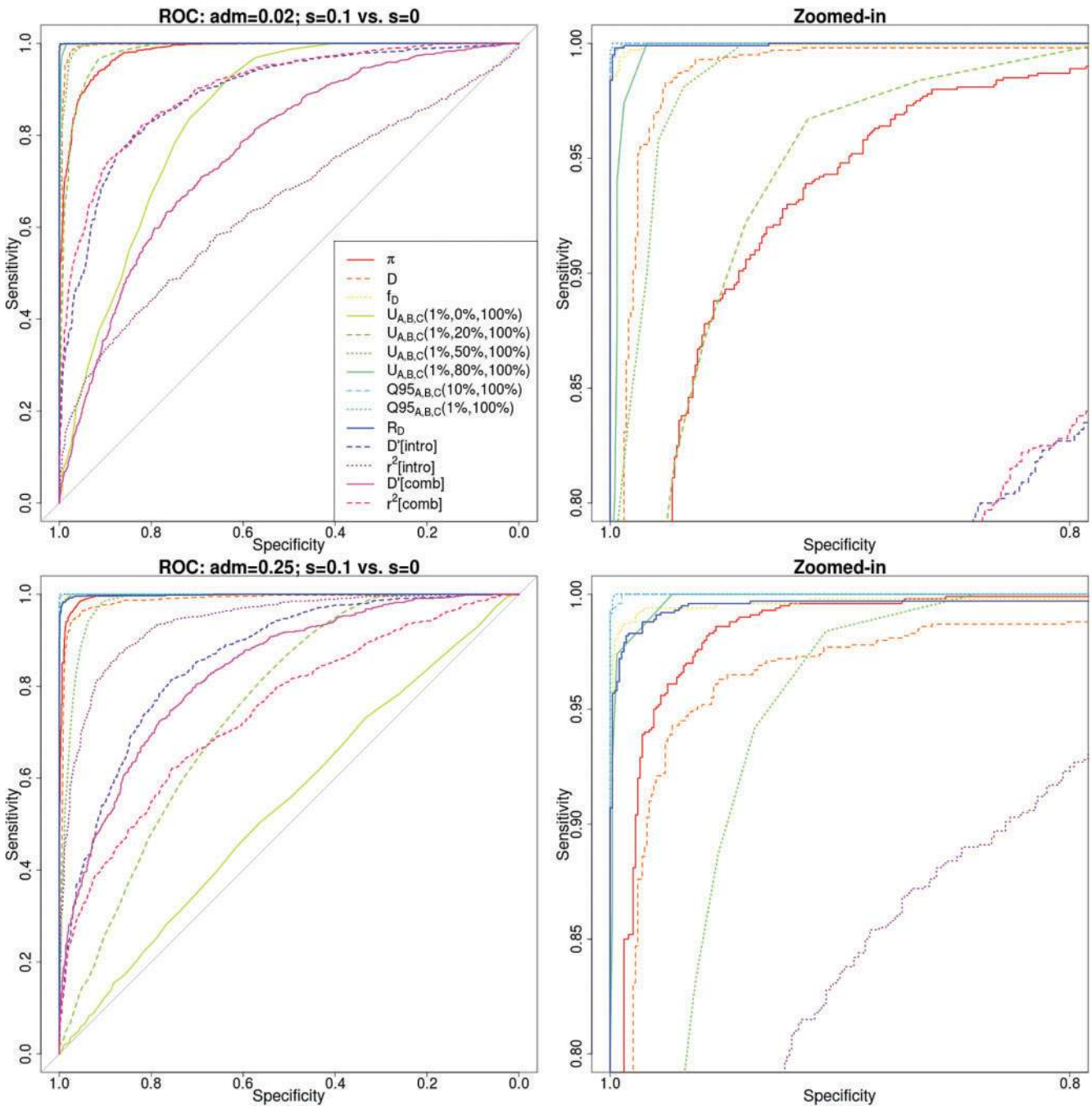
increases with increasing selection intensity and admixture rates, $\pi$ increases with increasing admixture rates, but decreases with increasing selection intensity. Thus, under AI the two forces cancel each other out, and we obtain a similar value of $\pi$ as under neutrality. Furthermore, the joint distributions of $Q95_{A,B,C}(1\%, 100\%)$ and $f_D$ or $R_D$ show particularly good separation among the different AI scenarios.

Another joint distribution that is especially good at separating different AI regimes is the combination of $Q95_{A,B,C}(w, 100\%)$ and $U_{A,B,C}(w, x, 100\%)$. In figure 4, we show this joint distribution, for different choices of $w$ (1% and 10%) and $x$ (20% and 50%). Here, with increasing intensity of selection and admixture, the number of uniquely shared sites and the quantile statistic increase, but the quantile statistic tends to only reach high values when selection is strong, even if admixture rates are low.

## Alternative Demographic Scenarios

We evaluated the performance of our statistics under various alternative demographic scenarios. First, we simulated a 5X bottleneck occurring in population B 1600 generations before the admixture event, and lasting 200 generations, to observe its effects on the power of the statistics for detecting AI (fig. 2B). Though we observe a reduction in power—most evident in the heterozygosity statistics—none of the statistics are very
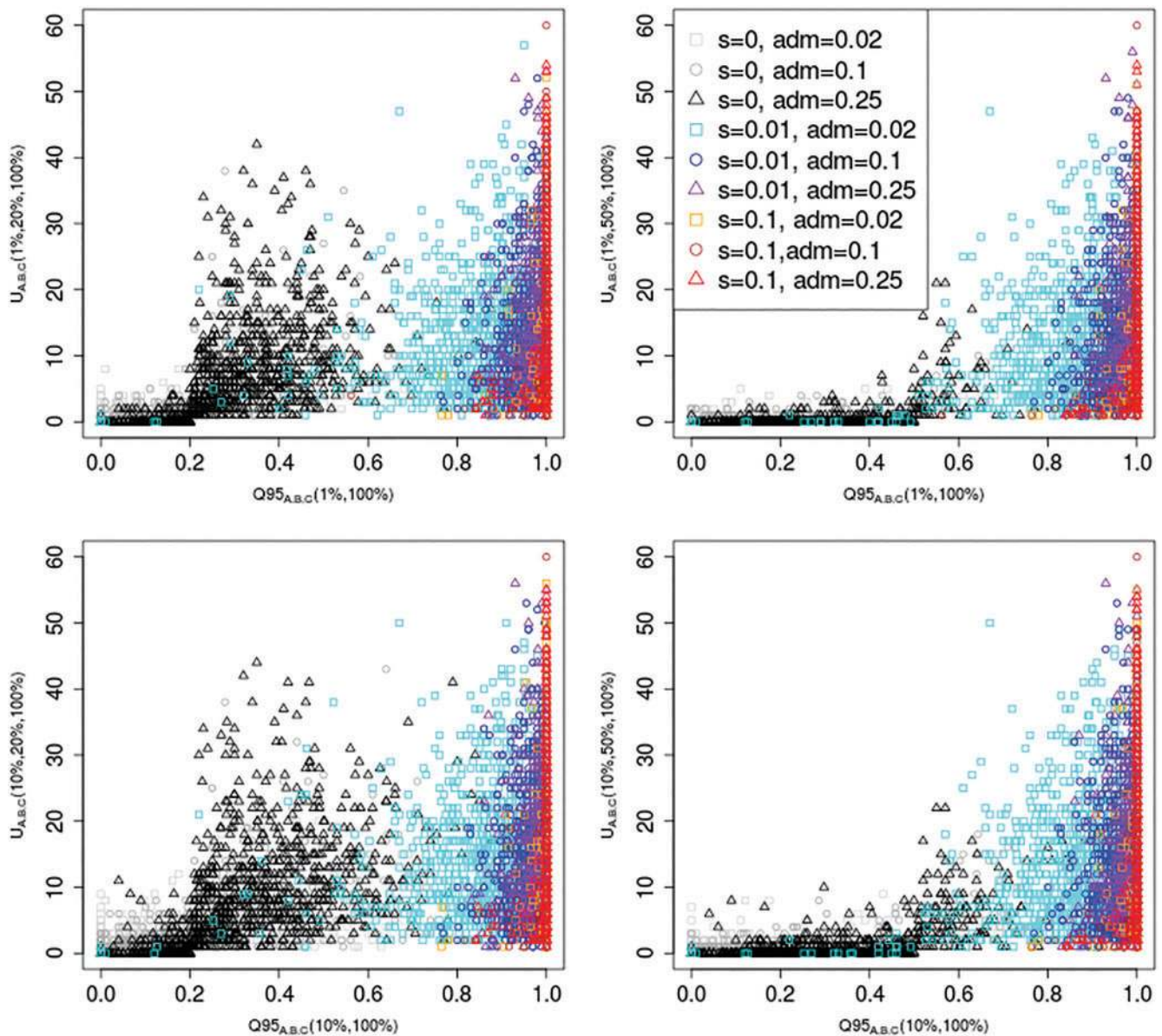
**Fig. 3.** Receiver operating characteristic curves for a scenario of adaptive introgression ($s = 0.1$) compared with a scenario of neutrality ($s = 0$), using 1,000 simulations for each case. Populations A and B split from each other 4,000 generations ago, and their ancestral population split from population C 16,000 generations ago. Population sizes were constant and set at $2N = 20,000$. The admixture event occurred 1,600 generations ago from population C to population B, at rate 2% (top panels) or 25% (bottom panels). The right panels are zoomed-in versions of the left panels.

strongly affected by this event (supplementary fig. S10, Supplementary Material online). We also simulated a bottleneck of equal size but occurring after the admixture event—starting 1,400 generations ago, and lasting 200 generations (fig. 2C). In this case, the sensitivity of all the statistics is strongly reduced when the admixture rate is low (supplementary fig. S11, Supplementary Material online). For example, when looking at the raw values of the $U_{A,B,C}$ and $Q95_{A,B,C}$ statistics, we observe that for low admixture rates the distribution under selection has a larger overlap with the

distribution under neutrality, which explains the low power (supplementary figs. S12 and S13, Supplementary Material online). Additionally, $U_{A,B,C}$ (but not $Q95$) seems to display more elevated values under neutrality in the bottleneck model than in the constant population size model. However, the relative performance of each statistic with respect to all the others does not appear to change substantially (supplementary fig. S11, Supplementary Material online).

We next explored a model where the introgressed haplotype was not immediately adaptive in the Eurasian

**Fig. 4.** Joint distribution of $Q95_{A,B,C}(w,y)$ and $U_{A,B,C}(w,x,y)$ for different choices of $w$ (1%, 10%) and $x$ (20%, 50%). We set $y$ to 100% in all cases. 100 individuals were sampled from panel A, 100 from panel B, and 2 from panel C. The demographic parameters were the same as in figure 3.

population, but instead underwent an intermediate period of neutral drift, before it became advantageous (fig. 2D). In such a situation, our power to detect AI is reduced, for all statistics (supplementary fig. S14, Supplementary Material online). This is particularly an issue when the admixture rate is low, as in those cases the starting frequency of the selected allele in the Eurasian population is low, so it is more likely to drift to extinction during the neutral period, before it can become advantageous.

We also evaluated the performance of our statistics under selective scenarios that did not involve adaptive introgression, to check which of them were sensitive to these models and which were not. Under a model of selection from de novo mutation (SDN, fig. 2E)—in which a single mutation appears in the receiving population after the split time between it and the nonadmixed population—the heterozygosity ($\pi$) and linkage disequilbrium statistics ($r^2[intro]$ and $D'[intro]$) are

the most sensitive ones (supplementary fig. S15, Supplementary Material online). This is expected, given that classical selective sweeps are known to strongly affect patterns of heterozygosity and linkage disequilibrium in the neighborhood of the selected site (Barton 1998; Kim and Stephan 2002; Kim and Nielsen 2004). Since all other statistics have very poor sensitivity to detect SDN, we expect to be able to distinguish signatures generated from SDN and AI. One caveat to this is the scenario in which a de novo selected mutation occurs on an introgressed haplotype immediately after an introgression event—before the haplotype has a chance to expand and recombine in the population—in which case our statistics will not be able to distinguish SDN from AI.

We also simulated a model of selection from standing variation (fig. 2F), by randomly selecting 20% of haplotypes within the introgressed population to be advantageous, after the split time between it and the nonintrogressed population. In this

case, all statistics perform poorly, especially when admixture is low. Interestingly, when admixture is high (supplementary fig. S16, Supplementary Material online), $Q95_{A,B,C}(1\%, 100\%)$ and $U_{A,B,C}(1\%, 0\%, 100\%)$ are the best performing statistics. This is likely because some of the haplotypes that are randomly chosen to be selected also happen to be ancestrally polymorphic and present in the archaic humans.

When we set ancestral structure to be our null model, we observe different behaviors depending on the strength of the migration rates. When the migration rates are strong (supplementary fig. S17, Supplementary Material online), we have excellent power to detect AI with several statistics, including $Q95_{A,B,C}(1\%, 100\%)$, $D$, $f_D$, $R_D$, and $U_{A,B,C}(1\%, 50\%, 100\%)$. When the rates are of medium strength (supplementary fig. S18, Supplementary Material online), the power is slightly reduced, but the same statistics are the ones that perform best. When the migration rates are weak—meaning ancestral structure is very strong—$Q95_{A,B,C}(1\%, 100\%)$ loses power, and the best-performing statistics are $R_D$, $D$, and $f_D$ (supplementary fig. S19, Supplementary Material online). We note, though, that the genome-wide $D$ observed under this last ancestral structure model ($D = 0.24$) is much more extreme than the genome-wide $D$ observed empirically between any Eurasian population and Neanderthals or Denisovans, suggesting that if there was ancestral structure between archaic and modern humans, it was likely not of this magnitude.

## Global Features of Uniquely Shared Archaic Alleles

Before identifying candidate genes for adaptive introgression, we investigated the frequency and number of uniquely shared alleles at the genome-wide level. Specifically, we wanted to know whether human populations varied in the number of sites with uniquely shared archaic alleles, and whether they also varied in the frequency distribution of these alleles. Therefore, we computed $U_{A,B,Nea,Den}(1\%,x,y,z)$ and $Q95_{A,B,Nea,Den}(1\%,y,z)$ for different choices of $x$, $y$, and $z$. We used different cutoffs for the frequency of the archaic allele ($x$) in the target population B: 0%, 20%, and 50%. To look for alleles uniquely shared with the Altai Neanderthal genome only, we set $y = 100\%$ and $z = 0\%$. To look for alleles uniquely shared with the Denisovan genome only, we set $y = 0\%$ and $z = 100\%$. Finally, to look for uniquely shared alleles matching both of the archaic genomes, we set $y = 100\%$ and $z = 100\%$.

We used each of the non-African panels in the 1000 Genomes Project phase 3 data (Auton et al. 2015) as the "target" panel (B), and chose the outgroup panel (A) to be the combination of all African populations (YRI, LWK, GWD, MSL, and ESN), excluding admixed African-Americans. We note that this is a conservative reference panel, as some of the African panels—like LWK—are from populations with a substantial amount of Eurasian ancestry (Auton et al. 2015), likely preventing the detection of introgressed segments at some loci.
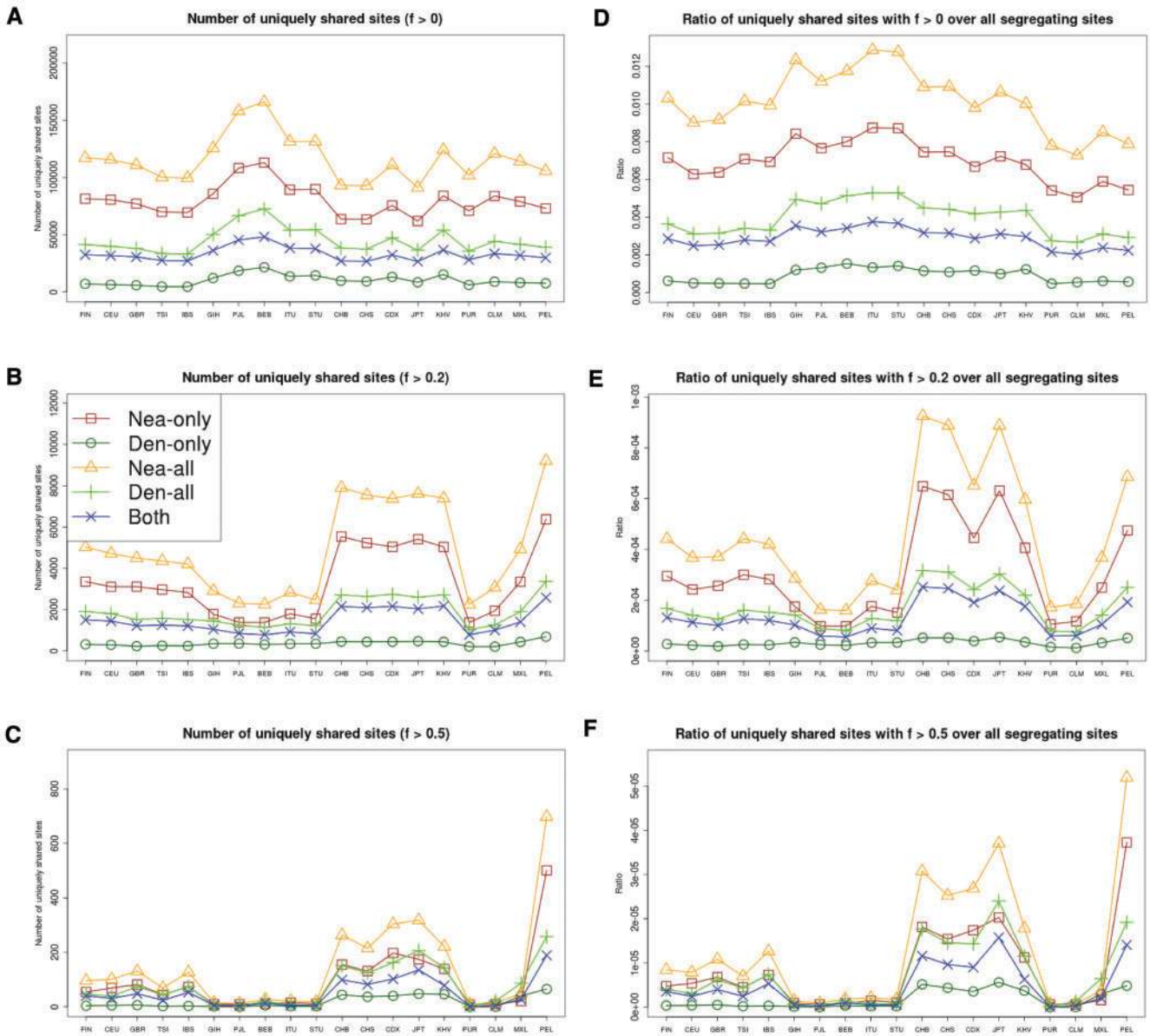
When setting $x = 0\%$ (i.e., not imposing a frequency cutoff in the target panel B), South Asians as a target population show the largest number of archaic alleles (fig. 5A). However, East Asians have a larger number of high-frequency uniquely

shared archaic alleles than Europeans and South Asians, for both $x = 20\%$ and $x = 50\%$ (fig. 5B–C). Population-specific D-statistics (using YRI as the nonadmixed population) also follow this trend (supplementary fig. S20, Supplementary Material online) and we observe this pattern when looking only at the X chromosome as well (supplementary fig. S21, Supplementary Material online). These results hold in comparisons with both archaic human genomes, but we observe a stronger signal when looking at Neanderthal-specific shared alleles. To correct for the fact that some panels have more segregating sites than others (and may therefore have more archaic-like segregating sites), we also scaled the number of uniquely shared sites by the total number of segregating sites per population panel (fig. 5D–F), and we see in general the same patterns, with the exception of a Peruvian panel, which we discuss further below. We also observe similar patterns when calculating $Q95_{A,B,Nea,Den}(1\%, y, z)$ genome-wide (supplementary fig. S22, Supplementary Material online). The elevation in $U_{A,B,Nea,Den}$ and $Q95_{A,B,Nea,Den}$ in East Asians may possibly result from higher levels of archaic ancestry in East Asians than in Europeans (Wall et al. 2013), which some studies argue could be due to additional admixture events occurring in East Asians (Kim and Lohmueller 2015; Vernot and Akey 2015;).

Surprisingly, the Peruvians (PEL) harbor the largest amount of high frequency mutations of archaic origin than any other single population, especially when using Neanderthals as bait (figs. 5B–C and supplementary fig. S21, Supplementary Material online). It is unclear whether this signal is due to increased drift or selection in this population. Skoglund and Jakobsson (2011) argue via simulations that if one analyzes a population with high amounts of recent genetic drift and excludes SNPs where the minor allele is at low frequency, some statistics that are meant to detect archaic ancestry—like $D$—may be artificially inflated. Our filtering procedure to select uniquely shared archaic alleles necessarily excludes sites where the archaic allele is at low frequency in the target panel, and the PEL panel comes from a population with a history of low effective population sizes (high drift) relative to other non-Africans (Auton et al. 2015), which could explain this pattern. This could also explain why the effect is not seen when $x = 0\%$ (fig. 5A), or when computing D-statistics (supplementary fig. S20, Supplementary Material online), both of which include sites with low-frequency alleles in their computation. Additionally, scaling the uniquely shared sites by the total number of segregating sites per population panel mitigates (but does not completely erase) this pattern. After scaling, PEL shows levels of archaic allele sharing within the range of the East Asian populations at $x = 20\%$ (fig. 5E), but is still the panel with the largest number of archaic sites at $x = 50\%$ (fig. 5F).

Furthermore, we plotted the values of $U_{AFR,X,Nea,Den}(1\%, x, y, z)$ and $Q95_{AFR,X,Nea,Den}(1\%, y, z)$ jointly for each population X, under different frequency cutoffs $x$. When $x = 0\%$, there is a generally inversely proportional relationship between the two scores (supplementary fig. S23, Supplementary Material online), but this becomes a directly proportional relationship when $x = 20\%$ (fig. 6) or $x = 50\%$ (supplementary fig. S27, Supplementary Material online). Here, we also clearly observe

**FIG. 5.** We computed the number of uniquely shared sites in the autosomes and the X chromosome between particular archaic humans and different choices of present-day non-African panels X (x-axis) from phase 3 of the 1000 Genomes Project. We used a shared frequency cutoff of 0% (A), 20% (B), and 50% (C). Nea-only = $U_{Afr,X,Nea,Den}(1\%, 20\%, 100\%, 0\%)$. Den-only = $U_{Afr,X,Nea,Den}(1\%, 20\%, 0\%, 100\%)$. Nea-all = $U_{Afr,X,Nea}(1\%, 20\%, 100\%)$. Den-all = $U_{Afr,X,Den}(1\%, 20\%, 100\%)$. Both = $U_{Afr,X,Nea,Den}(1\%, 20\%, 100\%, 100\%)$. Finally, we scaled each of the statistics from panels A to C by the number of segregating sites in each 1000 Genomes population panel, yielding panels D–F.

that PEL is an extreme panel with respect to both the number and frequency of archaic shared derived alleles, and that East Asian and American populations have more high-frequency archaic shared alleles than Europeans.
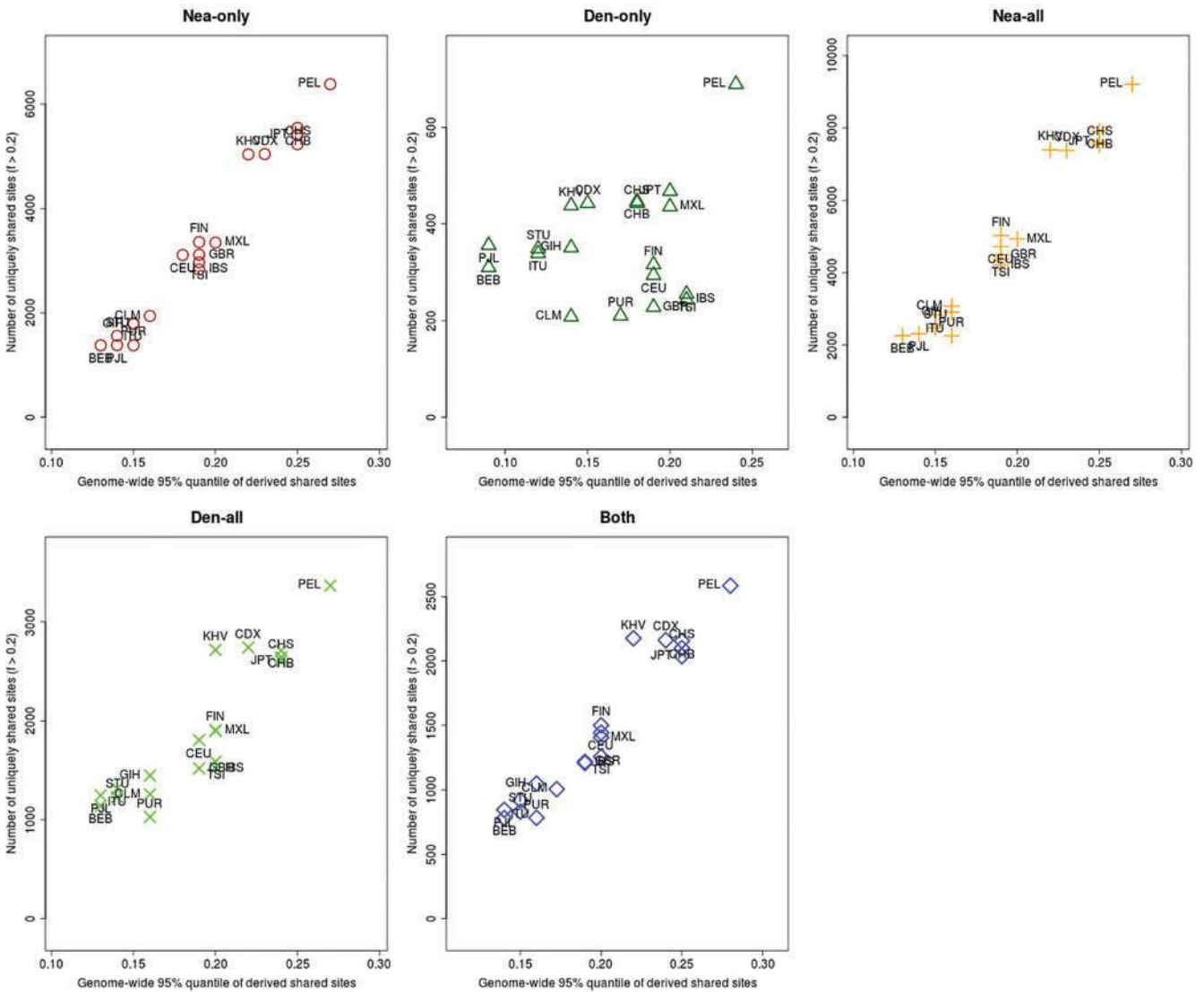
We checked via simulations if the observed excess of high frequency archaic derived mutations in Americans and especially Peruvians could be caused by genetic drift, as a consequence of the bottleneck that occurred in the ancestors of Native Americans as they crossed Beringia. We observe that if the introgressed population B undergoes a bottleneck, this can lead to a larger number of $U_{A,B,C}(w, x, y, z)$ for large values of $x$ (supplementary figs. S12, S13, and S24, Supplementary Material online). Indeed, population structure

analyses of the 1000 Genomes samples suggest that Peruvians have the largest amount of Native American ancestry (Auton et al. 2015) and show a bottleneck with a lack of recent population growth, which could explain this pattern. We also observe an increase in the variance of the distribution of U and Q95 in the presence of a bottleneck, especially when long and severe (supplementary figs. S25 and S26, Supplementary Material online).

### Candidate Regions for Adaptive Introgression

To identify adaptively introgressed regions of the genome, we computed $U_{A,B,C,D}(w, x, y, z)$ and $Q95_{A,B,C,D}(w, y, z)$ in 40 kb nonoverlapping windows along the genome, using the

**FIG. 6.** For each population panel from the 1000 Genomes Project, we jointly plotted the $U$ and $Q95$ statistics with an archaic frequency cutoff of $> 20\%$ within each population. Nea-only $= U_{Afr,X,Nea,Den}(1\%, 20\%, 100\%, 0\%)$ and $Q95_{Afr,X,Nea,Den}(1\%, 100\%, 0\%)$. Den-only $= U_{Afr,X,Nea,Den}(1\%, 20\%, 0\%, 100\%)$ and $Q95_{Afr,X,Nea,Den}(1\%, 0\%, 100\%)$. Nea-all $= U_{Afr,X,Nea}(1\%, 20\%, 100\%)$ and $Q95_{Afr,X,Nea}(1\%, 100\%)$. Den-all $= U_{Afr,X,Den}(1\%, 20\%, 100\%)$ and $Q95_{Afr,X,Den}(1\%, 100\%)$. Both $= U_{Afr,X,Nea,Den}(1\%, 20\%, 100\%, 100\%)$ and $Q95_{Afr,X,Nea,Den}(1\%, 100\%, 100\%)$.

low-coverage sequencing data from phase 3 of the 1000 Genomes Project (Auton et al. 2015). We used this window size because the mean length of introgressed haplotypes in Prüfer et al. (2014) was 44,078 bp (Supplementary Information 13). Our motivation was to find regions under AI in a particular panel B, using panel A as a nonintrogressed outgroup (generally Africans, unless otherwise stated). We used the high-coverage Altai Neanderthal genome (Prüfer et al. 2014) as bait panel C and the high-coverage Denisova genome (Meyer et al. 2012) as bait panel D. We deployed these statistics in three ways: (a) to look for Neanderthal-specific AI, we set $y = 100\%$ and $z = 0\%$; (b) to look for Denisova-specific AI, we set $y = 0\%$ and $z = 100\%$; (c) to look for AI matching both of the archaic genomes, we set $y = 100\%$ and $z = 100\%$ (supplementary fig. S1 and table S3, Supplementary Material online). To try to determine the adaptive pressure behind the putative AI event, we obtained

all the CCDS-verified genes located inside each window (Pruitt et al. 2009).

For guidance as to how high a value of $U$ and $Q95$, we would expect under neutrality, we used the simulations from figure 2 to obtain 95% empirical quantiles of the distribution of these scores under neutrality. Supplementary tables S1 and S2, Supplementary Material online show the 95% quantiles for these two statistics under various models of adaptive introgression and ancestral structure, for different choices of parameter values (see "Methods" section). When examining our candidates for AI below, we focused on windows whose values for $U_{A,B,Nea,Den}(w, x, y, z)$ and $Q95_{A,B,Nea,Den}(w, y, z)$ were both in the 99.9% quantile of their respective genome-wide distributions. We verified via simulations that, under a simple model of neutral admixture at a genome-wide rate of 2%, the estimated probability of obtaining values as high as these (or the false positive rate, FPR) was between

10.6% and 0%, depending on the target population chosen. The highest rates correspond to the African-American (ASW) admixed panel, as this panel contains high proportions of ancestry from the outgroup panel (unadmixed Africans) and are therefore not well-suited for our statistics. Excluding ASW, the highest estimated FPR was 5.5%.

We also calculated $D$ and $f_D$ along the same windows (using Africans as the nonadmixed population), and saw good agreement with the new statistics presented here (supplementary table S3, Supplementary Material online). Finally, we further validated the regions most likely to have been adaptively introgressed by searching for archaic tracts of introgression within them that were at high frequency, using a Hidden Markov Model (see below).

## Continental Populations

When focusing on adaptive introgression in continental populations, we first looked for uniquely shared archaic alleles specific to Europeans that were absent or almost absent ($< 1\%$ frequency) in Africans and East Asians. In addition, we also looked for uniquely shared archaic alleles in East Asians, which were absent or almost absent in Africans and Europeans. In this continental survey, we ignored Latin American populations as they have high amounts of European and African ancestry, which could confound our analyses. Figure 7 shows the number of sites with uniquely shared alleles for increasing frequency cutoffs in the introgressed population, and for different types of archaic alleles (Neanderthal-specific, Denisova-specific or common to both archaic humans). In other words, we calculated $U_{AFR,EUR,Nea,Den}(1\%, x, y, z)$ and $U_{AFR,EAS,Nea,Den}(1\%, x, y, z)$ for different values of $x$ (0%, 20%, 50%, and 80%) and different choices of $y$ and $z$, depending on which type of archaic alleles we were looking for. We observe that the regions in the extreme of the distributions for $x = 50\%$ corresponded very well to genes that had been previously found to be candidates for adaptive introgression from archaic humans in these populations, using more complex probabilistic methods (Sankararaman et al. 2014; Vernot and Akey 2014) or gene-centric approaches (Ding et al. 2013). These include *BNC2* (involved in skin pigmentation [Vanhoutteghem and Djian 2006; Jacobs et al. 2013]), *POU2F3* (involved in skin keratinocyte differentiation [Cabral et al. 2003; Takemoto et al. 2010]), *HYAL2* (involved in the response to UV radiation on human keratinocytes [Hašová et al. 2011]), *SIPA1L2* (involved in neuronal signaling [Spilker and Kreutz 2010]), and *CHMP1A* (a regulator of cerebellar development [Mochida et al. 2012]). To be more rigorous in our search for adaptive introgression, we looked at the joint distribution of the $U$ statistic and the $Q95$ statistic for the same choices of $w$, $y$, and $z$, and then selected the regions that were in the 99.9% quantiles of the distributions of both statistics (fig. 8, supplementary figs. S28 and S29, Supplementary Material online). We find that the strongest candidates here are *BNC2*, *POU2F3*, *SIPA1L2*, and the *HYAL2* region.

We also scanned for regions of the genome where South Asians (SAS) had uniquely shared archaic alleles at high frequency, which were absent or almost absent in Europeans, East Asians and Africans. In this case, we focused on $x = 20\%$ because we found that $x = 50\%$ left us with no candidate regions. Among the candidate regions sharing a large number of high-frequency Neanderthal alleles in South Asians, we find genes *ASTN2*, *SFMBT1*, *MUSTN1* and *MAML2* (supplementary fig. S30, Supplementary Material online). *ASTN2* is involved in neuronal migration (Wilson et al. 2010) and is associated with schizophrenia (Vrijenhoek et al. 2008; Wang et al. 2010). *SFMBT1* is involved in myogenesis (Lin et al. 2013) and is associated with hydrocephalus (Kato et al. 2011). *MUSTN1* plays a role in the regeneration of the muscoskeletal system (Krause et al. 2013). Finally, *MAML2* codes for a signaling protein (Lin et al. 2002; Wu et al. 2005), and is associated with cutaneous carcinoma (Winnes et al. 2007) and lacrimal gland cancer (Von Holstein et al. 2012).
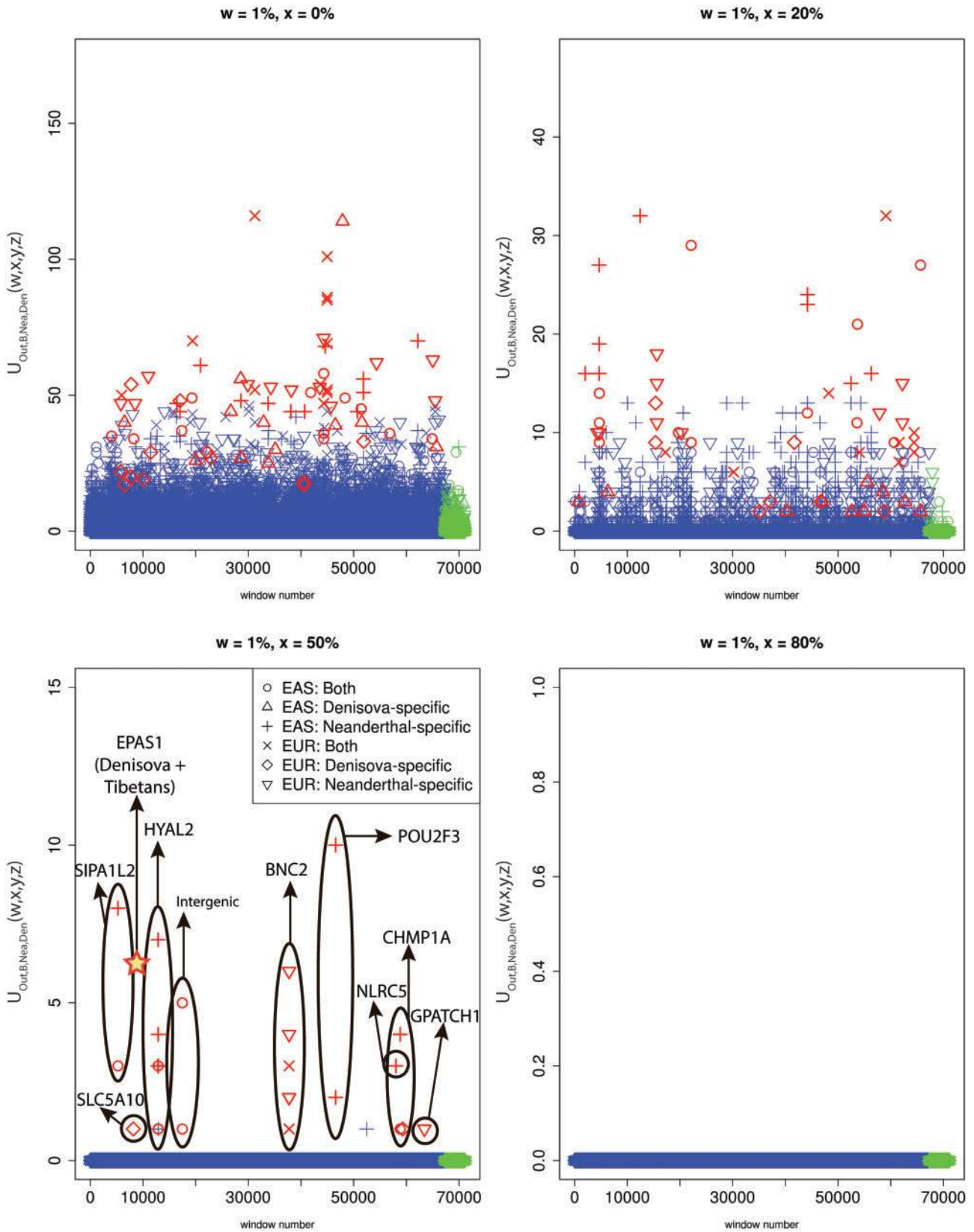
### Eurasia

We then looked for AI in all Eurasians (EUA = EUR + SAS + EAS, ignoring American populations) using Africans as the nonadmixed population (AFR, ignoring admixed African-Americans). Figure 8 shows the extreme outlier regions that are in the 99.9% quantiles for both $U_{EUA,AFR,Nea,Den}(1\%, 20\%, y, z)$ and $Q95_{EUA,AFR,Nea,Den}(1\%, y, z)$, whereas supplementary figure S31, Supplementary Material online shows the entire distribution. We focused on $x = 20\%$ because we found that $x = 50\%$ left us with almost no candidate regions. In this case, the region with by far the largest number of uniquely shared archaic alleles is the one containing genes *OAS1* and *OAS3*, involved in innate immunity (Knapp et al. 2003; Hamano et al. 2005; Fedetz et al. 2006; Lim et al. 2009). This region was previously identified as a candidate for AI from Neanderthals in non-Africans (Mendez et al. 2013). Another region that we recover and was previously identified as a candidate for AI is the one containing genes *TLR1* and *TLR6* (Dannemann et al. 2016; Deschamps et al. 2016). These genes are also involved in innate immunity and have been shown to be under positive selection in some non-African populations (Akira et al. 2006; Barreiro et al. 2009).

Interestingly, we find that a very strong candidate region in Eurasia contains genes *TBX15* and *WARS2*. This region has been associated with a variety of traits, including adipose tissue differentiation (Gburcik et al. 2012), body fat distribution (Heid et al. 2010; Liu et al., 2013, 2014; Shungin et al. 2015), hair pigmentation (Candille et al. 2004) facial morphology (Lausch et al. 2008; Pallares et al. 2015), ear morphology (Curry 1959), stature (Lausch et al. 2008), and skeletal development (Singh et al. 2005; Lausch et al. 2008). It was previously identified as being under positive selection in Greenlanders (Fumagalli et al. 2015), and it shows particularly striking signatures of adaptive introgression, so we devote a separate study to its analysis (Racimo et al. 2016).
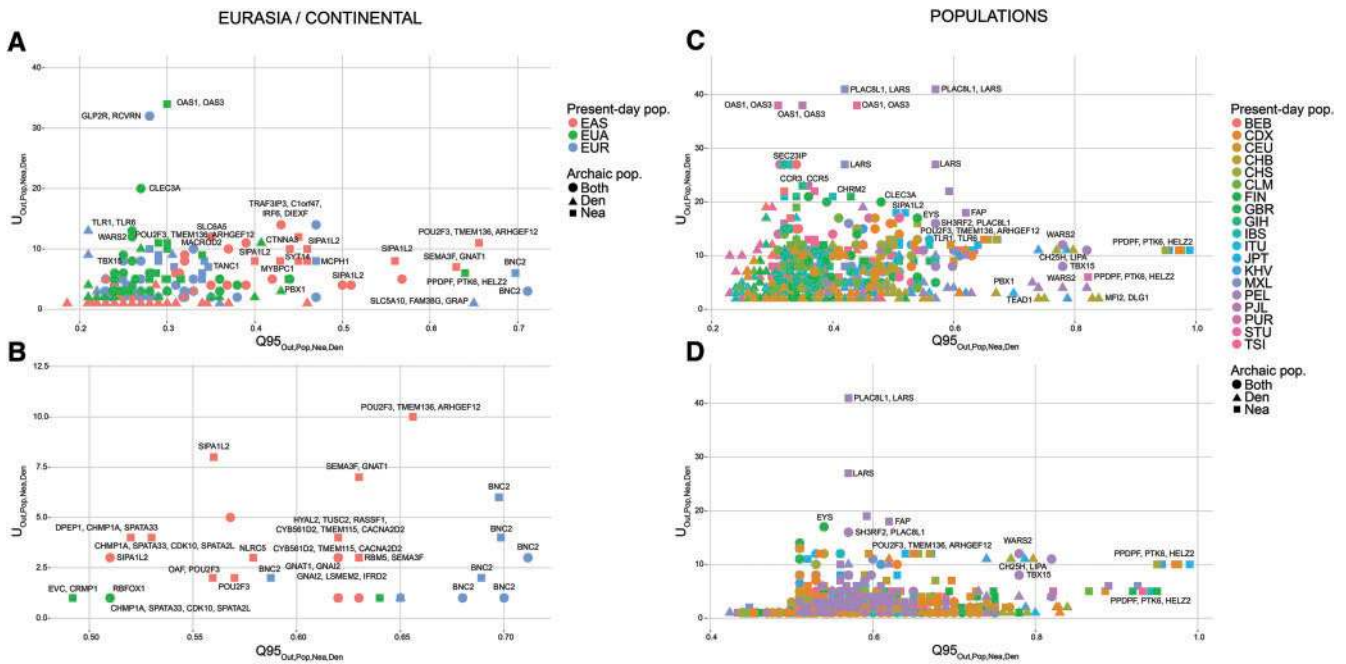
### Population-Specific Signals of Adaptive Introgression

To identify population-specific signals of AI, we looked for archaic alleles at high frequency in a particular non-African panel X, which were also at less than 1% frequency in all other

**Fig. 7.** We partitioned the genome into non-overlapping windows of 40 kb. Within each window, we computed $U_{Out,EUR,Nea,Den}(1\%, x, y, z)$ and $U_{Out,EAS,Nea,Den}(1\%, x, y, z)$, where Out = EAS + AFR for EUR as the target introgressed population, and Out = EUR + AFR for EAS as the target introgressed population. We searched for Neanderthal-specific alleles ($y = 100\%, z = 0\%$), Denisovan-specific alleles ($y = 0\%, z = 100\%$) and alleles present in both archaic genomes ($y = 100\%, z = 100\%$) that were uniquely shared with either EUR or EAS at frequencies above different cutoffs ($x = 0\%$, $x = 20\%$, $x = 50\%$, and $x = 80\%$). Windows that fall within the upper tail of the distribution for each modern-archaic population pair are colored in red ($P < 0.001$/number of pairs tested) and those that do not are colored in blue, except for those in the X chromosome, which

**FIG. 8.** We plotted the 40kb regions in the 99.9% highest quantiles of both the $Q95_{Out,Pop,Nea,Den}(1\%, y, z)$ and $U_{Out,Pop,Nea,Den}(1\%, x, y, z)$ statistics for different choices of target introgressed population (Pop) and outgroup non-introgressed population (Out), and different archaic allele frequency cutoffs within the target population ($x$). (A) We plotted the extreme regions for continental populations EUR (Out = EAS + AFR), EAS (Out = EUR + AFR), and Eurasians (EUA, Out = AFR), using a target population archaic allele frequency cutoff $x$ of 20%. (B) We plotted the extreme regions from the same statistics as in panel A, but with a more stringent target population archaic allele frequency cutoff $x$ of 50%. (C) We plotted the extreme regions for individual non-African populations within the 1000 Genomes data, using all African populations (excluding African-Americans) as the outgroup, and a cutoff $x$ of 20%. (D) We plotted the extreme regions from the same statistics as in panel C, but with a more stringent target population archaic allele frequency cutoff $x$ of 50%. Nea-only = $U_{Out,Pop,Nea,Den}(1\%, x, 100\%, 0\%)$ and $Q95_{Out,Pop,Nea,Den}(1\%, 100\%, 0\%)$. Den-only = $U_{Out,Pop,Nea,Den}(1\%, x, 0\%, 100\%)$ and $Q95_{Out,Pop,Nea,Den}(1\%, 0\%, 100\%)$. Both = $U_{Out,Pop,Nea,Den}(1\%, x, 100\%, 100\%)$ and $Q95_{Out,Pop,Nea,Den}(1\%, 100\%, 100\%)$.

non-African and African panels, excluding panel X (supplementary table S3, Supplementary Material online). This is a very restrictive requirement, and indeed, we only find a few windows in a single panel (PEL) with archaic alleles at more than 20% frequency. One of the regions with the largest number of uniquely shared Neanderthal sites in PEL contains gene *CHD2*, which codes for a DNA helicase (Woodage et al. 1997) involved in myogenesis (UniProtKB by similarity), and that is associated with epilepsy (Rauch et al. 2012; Carvill et al. 2013). We note, however, that the presence of these extreme regions could be due to the global elevation in the *U* statistic caused by higher levels of drift in Peruvians, as explained above.

## Shared Signals among Populations
In the previous section, we focused on regions where archaic alleles were uniquely at high frequencies in particular populations, but at low frequencies in all other populations. This precludes us from detecting AI regions that are shared across more than one non-African population. To address this, we conditioned on observing the archaic allele at less than 1%

frequency in a nonadmixed outgroup panel composed of all the African panels (YRI, LWK, GWD, MSL, and ESN), excluding African-Americans, and then looked for archaic alleles at high frequency in particular non-African populations. Unlike the previous section, we did not condition on the archaic allele being at low frequency in other non-African populations as well. The whole joint distributions of *U* and *Q95* for this choice of parameters for each non-African panel are shown in supplementary figs. S32–S50, Supplementary Material online whereas regions in the 99.9% quantile for both statistics are shown in figure 8.

Here, we recapitulate many of the findings from our Eurasian and continental-specific analyses above, like *TLR1/TLR6*, *BNC2*, *OAS1/OAS3*, *POU2F3*, *LIPA*, and *TBX15/WARS2* (fig. 8). For example, just as we found that *POU2F3* was an extreme region in the East Asian (EAS) continental panel, we separately find that almost all populations composing that panel (CHB, KHV, CHS, CDX, and JPT) have archaic alleles in that region at disproportionately high frequency, relative to their frequency in Africans. Additionally, we can learn things

**FIG. 7**. Continued
are in green. Ovals drawn around multiple points contain multiple windows with uniquely shared alleles that are contiguous. For comparison, the number of high frequency uniquely shared sites between Denisova and Tibetans is also shown (Huerta-Sánchez et al. 2014), although Tibetans are not included in the 1000 Genomes data and the region is 32 kb long, so this may be an underestimate.

we would not have detected at the continental level. For example, the Bengali from Bangladesh (BEB)—a South Asian population—also have archaic alleles at very high frequencies in the same genomic region.

We detected several genes that appear to show signatures of AI across various populations (fig. 8). One of the most extreme examples is a 120 kb region containing the LARS gene, with 76 uniquely shared Neanderthal alleles at < 1% frequency in Africans and > 50% frequency in Peruvians, which are also at > 20% frequency in Mexicans. LARS codes for a leucin-tRNA synthetase (Giles et al. 1980), and is associated with liver failure syndrome (Casey et al. 2012). Additionally, a region containing the gene ZFHX3 displays an elevated number of uniquely shared Neanderthal sites in PEL, and we also observe this when looking more broadly at East Asians (EAS) and—based on the patterns of inferred introgressed tracts (see below)—in various American (AMR) populations as well. ZFHX3 is involved in the inhibition of estrogen receptor-mediated transcription (Dong et al. 2010) and has been associated with prostate cancer (Sun et al. 2005).

We also find several Neanderthal-specific uniquely shared sites in American panels (PEL, CLM, MXL) in a region previously identified as harboring a risk haplotype for type 2 diabetes (chr17:6880001–6960000) (Sigma Type 2 Diabetes Consortium 2014). This is consistent with previous findings suggesting the risk haplotype was introgressed from Neanderthals and is specifically present at high frequencies in Latin Americans (Sigma Type 2 Diabetes Consortium 2014). The region contains gene SLC16A11, whose expression is known to alter lipid metabolism (Sigma Type 2 Diabetes Consortium 2014). We also find that the genes FAP/IFIH1 have signals consistent with AI, particularly in PEL. This region has been previously associated with type 1 diabetes (Qu et al. 2008; Liu et al. 2009). A previous analysis of this region has suggested that the divergent haplotypes in it resulted from ancestral structure or balancing selection in Africa, followed by local episodes of positive selection in Europe, Asia, and the Americas (Fumagalli et al. 2010). A more recent analysis has found this as a region of archaic AI in Melanesians as well (Vernot et al. 2016).

Another interesting candidate region contains two genes involved in lipid metabolism: LIPA and CH25H. We find a 40 kb region with 11 uniquely shared Denisovan alleles that are at low (< 1%) frequency in Africans and at very high (> 50%) frequency in various South and East Asian populations (JPT, KHV, CHB, CHS, CDX, and BEB). The Q95 and D statistics in this region are also high across all of these populations, and we also find this region to have extreme values of these statistics in our broader Eurasian scan. The LIPA gene codes for a lipase (Warner et al. 1980) and is associated with cholesterol ester storage disease (Klima et al. 1993) and Wolman disease (Aslanidis et al. 1996). In turn, the CH25H gene codes for a membrane hydroxylase involved in the metabolism of cholesterol (Lund et al. 1998) and associated with Alzheimer's disease (Shibata et al. 2006) and antiviral activity (Liu et al. 2013).

Finally, we find a region harboring between 3 and 10 uniquely shared Neanderthal alleles (depending on the panel

used) in various non-African populations. This region was identified earlier by Sankararaman et al. (2014) and contains genes PPDPF, PTK6 and HELZ2. PPDPF codes for a probable regulator of pancreas development (UniProtKB by similarity). PTK6 codes for an epithelial signal transducer (Kamalati et al. 1996) and HELZ2 codes for a helicase that works as a transcriptional coactivator for nuclear receptors (Surapureddi et al. 2002; Tomaru et al. 2006).
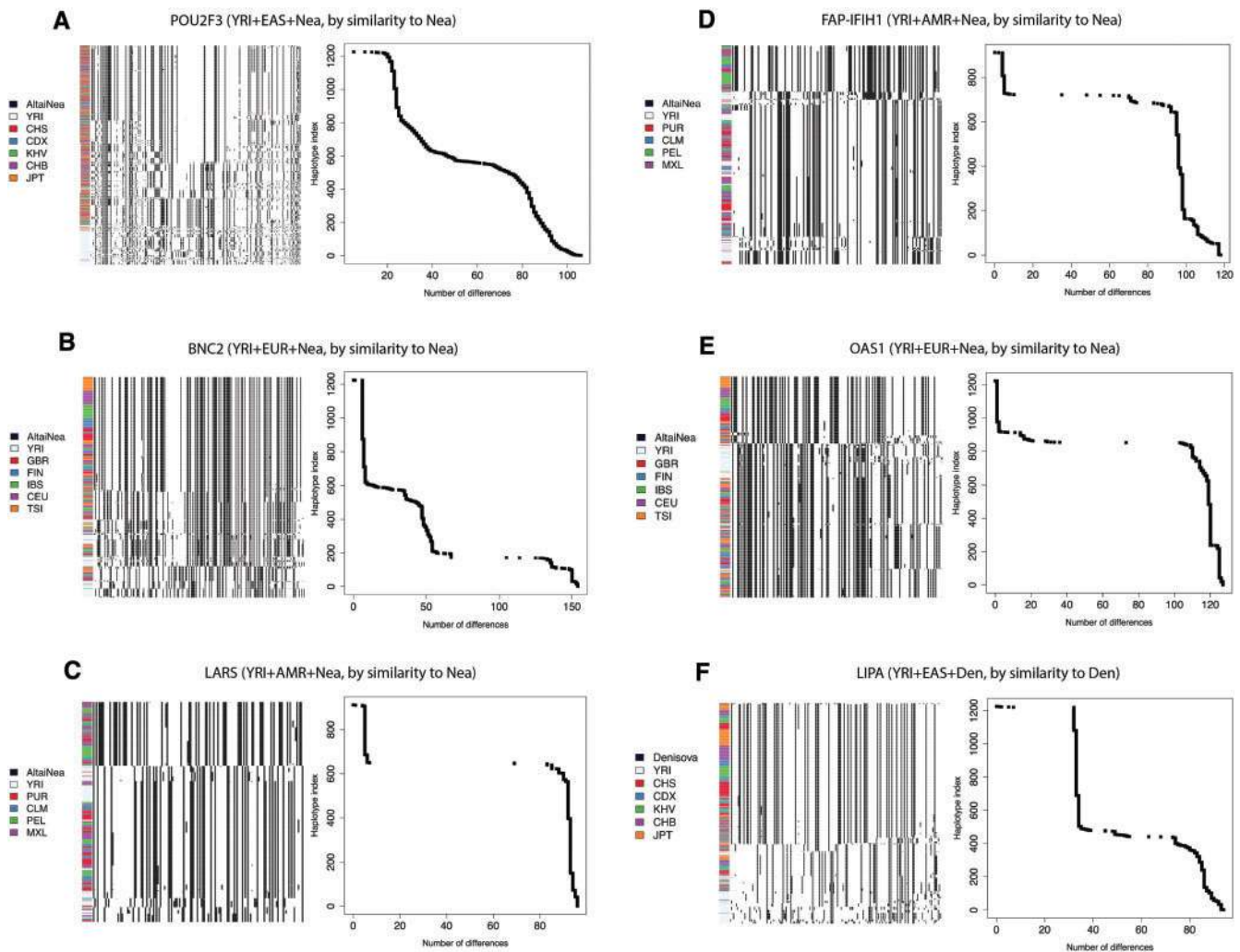
## The X Chromosome

Previous studies have observed lower levels of archaic introgression in the X chromosome relative to the autosomes (Sankararaman et al. 2014; Vernot and Akey 2014). Here, we observe a similar trend: compared with the autosomes, the X chromosome contains a smaller number of windows with sites that are uniquely shared with archaic humans (fig. 7). For example, for $w = 1\%$ and $x = 20\%$, we observe that, in Europeans, 0.4% of all windows in the autosomes have at least one uniquely shared site with Neanderthals or Denisovans, whereas only 0.05% of all windows in the X chromosome have at least one uniquely shared site ($P = 4.985 \times 10^{-4}$, chi-squared test assuming independence between windows). The same pattern is observed in East Asians ($P = 1.852 \times 10^{-8}$).

Nevertheless, we do identify some regions in the X chromosome exhibiting high values for both $U_{A,B,C,D}(w,x,y,z)$ and $Q95_{A,B,C,D}(w,y,z)$. For example, a region containing gene DHRSX contains a uniquely shared site where a Neanderthal allele is at < 1% frequency in Africans, but at > 50% frequency in a British panel (GBR). The site is also at high frequency (29%–47%) in the other European panels, but never as high as in GBR (55%). It is also surrounded by five neighboring SNPs that have intermediate Neanderthal allele frequencies (24%–41%) in GBR. Another region contains the gene DMD and harbors two uniquely shared sites where two archaic (Denisovan/Neanderthal) alleles are also at low (< 1%) frequency in Africans but at > 50% frequency in Peruvians. DHRSX codes for an oxidoreductase enzyme (Persson et al. 2009) whereas DMD is a well-known gene because mutations in it cause muscular dystrophy (Wood et al. 1987), and was also previously identified as having signatures of archaic introgression in non-Africans (Yotova et al. 2011). We note, however, that our simulations do not account for the particular inheritance and recombination patterns of the X chromosome, so caution should be taken when calling these regions as under AI.

## Introgressed Haplotypes in Candidate Loci

We inspected the haplotype patterns of candidate loci with support in favor of AI. We display the haplotypes for selected populations at seven regions: POU2F3 (fig. 9A), BNC2 (fig. 9B), LARS (fig. 9C), FAP/IFIH1 (fig. 9D), OAS1 (fig. 9E), LIPA (fig. 9F), and SLC16A11 (supplementary fig. S51C, Supplementary Material online). We included continental populations that show a large number of uniquely shared archaic alleles, and included YRI as a representative African population. We then clustered and ordered the haplotypes by similarity to the closest archaic genome (Altai Neanderthal or Denisova) (fig.

**Fig. 9.** We explored the haplotype structure of six candidate regions with strong evidence for AI. For each region, we applied a clustering algorithm to the haplotypes of particular human populations and then ordered the clusters by decreasing similarity to the archaic human genome with the larger number of uniquely shared sites (see "Methods" Section). We also plotted the number of differences to the archaic genome for each human haplotype and sorted them simply by decreasing similarity. In the latter case, no clustering was performed, so the rows in the cumulative difference plots do not necessarily correspond to the rows in the adjacent haplotype structure plots. *POU2F3*: chr11:120120001–120200000. *BNC2*: chr9:16720001–16760000. *LARS*: chr5:145480001–145520000. *FAP/IFIH1*: chr2:163040001–163120000. *OAS1*: chr12:113360001–113400000. *LIPA*: chr10:90920001–90980000.

9). As can be observed, all these regions tend to show sharp distinctions between the putatively introgressed haplotypes and the nonintrogressed ones. This is also evident when looking at the cumulative number of differences of each haplotype to the closest archaic haplotype, where we see a sharp rise in the number of differences, indicating strong differentiation between the two sets of haplotypes. Additionally, the YRI haplotypes tend to predominantly belong to the nonintrogressed group, as expected.

## Consequences of Relaxing the Outgroup Frequency Cutoff

When using a more lenient cutoff for the outgroup panel (10% maximum frequency, rather than 1%), we find a few genes that display values of the *U* statistic that are suggestive of AI, and that have been previously found to be under strong positive selection in particular human populations (Voight et al. 2006; Pickrell et al. 2009). The most striking examples are

*TYRP1* in EUR (using EAS + AFR as outgroup) and *OCA2* in EAS (using EUR + AFR as outgroup) (supplementary table S3, Supplementary Material online). Both of these genes are involved in pigmentation. We caution, however, that the reason why they carry archaic alleles at high frequency may simply be because their respective selective sweeps pushed an allele that was segregating in both archaic and modern humans to high frequency in modern humans, but not necessarily via introgression.

In fact, *TYRP1* only stands out as an extreme region for the number of archaic shared alleles in EUR when using the lenient 10% cutoff, but not when using the more stringent 1% cutoff. When looking at these SNPs in more detail, we find that their allele frequency in Africans (~20%) is even higher than in East Asians (~1%), largely reflecting population differentiation across Eurasia due to positive selection (Pickrell et al. 2009), rather than adaptive introgression. When

exploring the haplotype structure of this gene (supplementary fig. S51B, Supplementary Material online), we find one haplotype that shows similarities to archaic humans but is at low frequency. In the combined YRI + EUR panel, just 6.37% of all haplotypes have 36 or less differences to the Neanderthal genome, and this number is roughly the point of transition between the archaic-like and the nonarchaic-like haplotypes (supplementary fig. S51B, Supplementary Material online). There is a second—more frequent—haplotype that is more distinct from archaic humans but present at high frequency in Europeans. The uniquely shared sites obtained using the lenient ($<$ 10%) allele frequency outgroup cutoff are tagging both haplotypes together, rather than just the highly differentiated archaic-like haplotype.

*OCA2* has several sites with uniquely shared alleles in EAS (AFR + EUR as outgroup) when using the lenient 10% cutoff, but only a few (2) shared archaic sites when using the $<$ 1% outgroup frequency cutoff. When exploring the haplotype structure of this gene, we fail to find a clear-cut differentiation between putatively introgressed and nonintrogressed haplotypes, so the evidence for adaptive introgression in this region is also weak. A close inspection of its haplotype structure shows that *OCA2* does not show a large number of differences between the haplotype classes that are closer and those that are distant from the archaic humans (supplementary fig. S51A, Supplementary Material online).

Finally, using the lenient outgroup cutoff of $<$ 10% and a target cutoff of $>$20%, we find the gene with the highest number of uniquely shared sites among all the populations and cutoffs we tested: *MUC19*. This region is rather impressive in containing 115 sites where the archaic alleles are shared between the Mexican panel (MXL) and the Denisovan genome at more than 20% frequency, when using all populations that are not MXL as the outgroup. However, the actual proportion of individuals that contain a Denisova-like haplotype (though highly differentiated from the rest of present-day human haplotypes) is very small. Only 11.86% of haplotypes in the combined YRI + AMR panel show 69 differences or less to the closest archaic genome (Denisova), and the next closest haplotype has 134 differences (supplementary fig. S51D, Supplementary Material online).

Overall, a finer investigation of these three cases suggests that using a lenient outgroup frequency cutoff may lead to misleading inferences. Nevertheless, the haplotype structure of these genes and their relationship to their archaic human counterparts are quite unusual. It remains to be determined whether these patterns could be caused by either positive selection or introgression alone, or whether a combination of these or other demographic forces is required to explain them.

## Inferred Introgressed Tracts

We used an HMM (Seguin-Orlando et al. 2014) to verify that the strongest candidate regions effectively contained archaic segments of a length that would be consistent with introgression after the population divergence between archaic and modern humans. For each region, we used the closest archaic genome (Altai Neanderthal or Denisova) as the putative

source of introgression. We then plotted the inferred segments in non-African continental populations for genes with strong evidence for AI. Among these, genes with Neanderthal as the closest source (supplementary figs. S52–S59, Supplementary Material online) include: *POU2F3* (EAS,SAS), *BNC2* (EUR), *OAS1* (Eurasians), *LARS* (AMR), *FAP/IFIH1* (PEL), *CHD2* (PEL), *TLR1-6* (EAS), and *ZFHX3* (PEL). Genes with Denisova as the closest source (supplementary figs. S60 and S61, Supplementary Material online) include: *LIPA* (EAS, SAS, and AMR) and *MUSTN1* (SAS).

## Testing for Enrichment in Genic Regions

We aimed to test whether uniquely shared archaic alleles at high frequencies were enriched in genic regions of the genome. We looked at archaic alleles at high frequency in any of the non-African panels that were also at low frequency ($<$ 1%) in Africans. As background, we used all archaic alleles that were at any frequency larger than 0 in the same non-African populations, and that were also at low frequency in Africans. We then tested whether the high-frequency archaic alleles tended to occur in genic regions more often than expected.

SNPs in introgressed blocks will tend to cluster together and have similar allele frequencies, which could cause a spurious enrichment signal. To correct for the fact that SNPs at similar allele frequencies will cluster together (as they will tend to co-occur in the same haplotypes), we performed linkage disequilibrium (LD) pruning using two methods. In one (called "LD-1"), we downloaded the approximately independent European LD blocks published in Berisa and Pickrell (2016). For each set of high frequency derived sites, we randomly sampled one SNP from each block. In a different approach (called "LD-2"), for each set of high frequency derived sites, we subsampled SNPs such that each SNP was at least 200 kb apart from each other. We then tested these two types of LD-pruned SNP sets against 1,000 SNP sets of equal length that were also LD-pruned and that were obtained randomizing frequencies and collecting SNPs in the same ways as described above.

Regardless of which LD method we used, we find no significant enrichment in genic regions for high-frequency ($>$50%) Neanderthal alleles (LD-1 $P =$ 0.352, LD-2 $P = 0.161$) or Denisovan alleles (LD-1 $P = 0.348$, LD-2 $P = 0.192$). Similarly, we find no enrichment for medium-to-high-frequency ($>$20%) Neanderthal alleles (LD-1 $P = 0.553$, LD-2 $P = 0.874$) or Denisovan alleles (LD-1 $P = 0.838$, LD-2 $P = 0.44$).

## Discussion

Here, we carried out one of the first investigations into the joint dynamics of archaic introgression and positive selection, to develop statistics that are informative of AI. We find that one of the most powerful ways to detect AI is to look at both the number and allele frequency of mutations that are uniquely shared between the introgressed and the archaic populations. Such mutations should be abundant and at high frequencies in the introgressed population if AI occurred. In particular, we identified two novel summaries of the data

that capture this pattern quite well: the statistics Q95 and U. These statistics can recover loci under AI and are easy to compute from genomic data, as they do not require phasing.

We have also studied the general landscape of archaic alleles and their frequencies in present-day human populations. By scanning the present-day human genomes from phase 3 of the 1000 Genomes Project (Auton et al. 2015) using these and other summary statistics, we were able to recapitulate previous AI findings (like the TLR [Dannemann et al. 2016; Deschamps et al. 2016] and OAS regions [Mendez et al. 2013]) as well as identify new candidate regions for AI in Eurasia (like the LIPA gene and the FAP/IFIH1 region). These mostly include genes involved in lipid metabolism, pigmentation and innate immunity, as observed in previous studies (Khrameeva et al. 2014; Sankararaman et al. 2014; Vernot and Akey 2014). Phenotypic changes in these systems may have allowed archaic humans to survive in Eurasia during the Pleistocene, and may have been passed on to present-day human populations during their expansion out of Africa.

When using more lenient definitions of what we consider to be "uniquely shared archaic alleles" we find sites containing these alleles in genes that have been previously found to be under positive selection (like OCA2 and TYRP1) but not necessarily under adaptive introgression. Whereas these do not show as strong signatures of adaptive introgression as genes like BNC2 and POU2F3, their curious haplotype patterns and their relationship to archaic genomes warrants further exploration.

We tested whether uniquely shared archaic alleles at high frequencies in non-Africans were significantly more likely to be found in genic regions, relative to all shared archaic alleles, but did not find a significant enrichment. Though this suggests archaic haplotypes subject to AI may not be preferentially found near or inside genes, it may also be a product of a lack of power, or of the fact that not all uniquely shared archaic alleles may be truly introgressed. As mentioned before, some of these alleles may be present due to incomplete lineage sorting, which could add noise to the test signal. A more rigorous—and possibly more powerful—test could involve testing whether HMM-inferred introgressed archaic segments at high frequency tend to be found in genic regions, relative to all inferred introgressed archaic segments, controlling for features like the length of introgressed segments and the sensitivity of the HMM to different regions of the genome. However, we did not pursue this line of research further.

In this study, we have mostly focused on positive selection for archaic alleles. One should remember, though, that a larger proportion of introgressed genetic material was likely maladaptive to modern humans, and therefore selected against. Indeed, two recent studies have shown that negative selection on archaic haplotypes may have reduced the initial proportion of archaic material present in modern humans immediately after the hybridization event(s) (Harris and Nielsen 2015; Juric et al. 2015).

Another caveat is that some regions of the genome display patterns that could be consistent with multiple introgression events, followed by positive selection on one or more distinct archaic haplotypes (Dannemann et al. 2016). In this study, we have simply focused on models with a single pulse of admixture—followed immediately by selection or with an intermediate neutrality period in the introgressed population. We have not considered complex scenarios with multiple sources of introgression. Additionally, the currently limited availability of high-coverage archaic human genomes may prevent us from detecting AI events for which the source may not have been closely related to the sequenced Denisovan or Altai Neanderthal genomes. This may include other Neanderthal or Denisovan subpopulations, or other (as yet unsampled) archaic groups that may have lived in Africa and Eurasia.

It is also worth noting that positive selection for archaic haplotypes may be due to heterosis, rather than adaptation to particular environments (Harris and Nielsen 2015). That is, archaic alleles may not have been intrinsically beneficial, but simply protective against deleterious recessive modern human alleles, and therefore selected after their introduction into the modern human gene pool. The degree of dominance of deleterious alleles in humans remains elusive, so it is unclear how applicable this model would be to archaic admixture in humans.

Many of the statistics we introduced in this study have their drawbacks: notably, they depend on simulations to assess significance and some—like U—may be sensitive to local variation in mutation rates across the genome. Nevertheless, they serve as useful exploratory tools, as they highlight a characteristic signature left by AI in present-day human genomes. Future avenues of research could involve developing ways to incorporate uniquely shared sites into a robust test of selection that specifically targets regions under AI. For example, one could think about modifying statistics based on local between-population population differentiation, like PBS (Yi et al. 2010), so that they are only sensitive to allele frequency differences at sites that show signatures of archaic introgression.

Finally, whereas this study has largely focused on human AI, several other species also show suggestive signatures of AI (Hedrick 2013). Assessing the extent and prevalence of AI and uniquely shared sites in other biological systems could provide new insights into their biology and evolutionary history. This may also serve to better understand how populations of organisms respond to introgression events, and to derive general principles about the interplay between admixture and natural selection.

## Methods

### Summary Statistics Sensitive to Adaptive Introgression

Several statistics have been previously deployed to detect AI events (reviewed in Racimo et al. [Racimo et al. 2015]). We briefly describe these below, as well as three new statistics tailored specifically to find this signal (table 1). One of the simplest approaches consists of applying the D statistic (Green et al. 2010; Durand et al. 2011) locally over windows of the genome. The D statistic was originally applied to compare a single human genome against another human

genome, so as to detect excess shared ancestry between one of the genomes and a genome from an outgroup population. Application of this statistic comparing non-Africans and Africans served as one of the pieces of evidence in support of Neanderthal admixture into non-Africans. However, it can also be computed from large panels of multiple individuals instead of single genomes. This form of the $D$ statistic has been applied locally over windows of the genome to detect regions of excess shared ancestry between an admixed population and a source population (Kronforst et al. 2013; Smith and Kronforst 2013).

The $D$ statistic, however, can be confounded by local patterns of diversity, as regions of low diversity may artificially inflate the statistic even when a region was not adaptively introgressed. To correct for this, Martin et al. (2015) developed a similar statistic called $f_D$ which is less sensitive to differences in diversity along the genome. Both of these patterns exploit the excess relatedness between the admixed and the source population.

AI is also expected to increase linkage disequilibrium (LD), as an introgressed fragment that rises in frequency in the population will have several closely linked loci that together will be segregating at different frequencies than they were in the recipient population before admixture. Thus, two well-known statistics that are informative about the amount of LD in a region—$D'$ and $r^2$—could also be informative about adaptive introgression. To apply them over regions of the genome, we can take the average of each of the two statistics over all SNP pairs in a window. In the section below, we calculate these statistics in two ways: (a) using the introgressed panel only ($D'[intro]$ and $r^2[intro]$), and (b) using the combination of the introgressed and the nonintrogressed panels ($D'[comb]$ and $r^2[comb]$). The first way (intro) should capture patterns of within-population LD in the introgressed population under AI, whereas the second way (comb) should capture patterns of global LD across both populations. If the introgressed population has a particular set of archaic haplotypes at high frequency that are highly differentiated from the nonarchaic haplotypes in the nonintrogressed panel, we expect the second way to be more powerful at distinguishing AI from neutrality.

We also introduce three new statistics that one would expect, *a priori*, to be particularly effective at identifying windows of the genome that are likely to have undergone adaptive introgression: $R_D$, $U$ and $Q95$. $R_D$ is computed by calculating—in a window of the genome—the ratio of the sequence divergence between an individual from the source population and an admixed individual, and the sequence divergence between an individual from the source population and a nonadmixed individual. One can then take the average of this ratio over all individuals in the admixed and nonadmixed panels. This average should be larger if the introgressed haplotype is present in a large number of individuals of the admixed population. We call this statistic $R_D$.

Second, for a window of arbitrary size, let $U_{A,B,C}(w, x, y)$ be defined as the number of sites where a sample $C$ (the "bait") from an archaic source population (which could be as small as a single diploid individual) has a particular allele at frequency $y$, and that allele is at a frequency smaller than $w$ in a sample $A$

(the "outgroup") of a population but larger than $x$ in a sample $B$ (the "target") of another population (fig. 1). In other words, we are looking for sites that contain alleles shared between an archaic human genome and a test population, but absent or at very low frequencies in an outgroup (usually nonadmixed) population. For example, suppose we are looking for Neanderthal adaptive introgression in the Han Chinese (CHB). In that case, we can consider CHB as our target panel, and use Africans as the outgroup panel and a single Neanderthal genome as the bait. If $U_{AFR,CHB,Nea}(1\%, 20\%, 100\%) = 4$ in a window of the genome, that means there are four sites in that window where the Neanderthal genome is homozygous for a particular allele and that allele is present at a frequency smaller than 1% in Africans but larger than 20% in Han Chinese. In other words, there are four sites that are uniquely shared at more than 20% frequency between Han Chinese and Neanderthal, but not with Africans.

This statistic can be further parametrized if we have samples from two different archaic populations (e.g., a Neanderthal genome and a Denisova genome). In that case, we can define $U_{A,B,C,D}(w, x, y, z)$ as the number of sites where the archaic sample $C$ has a particular allele at frequency $y$ and the archaic sample $D$ has that allele at frequency $z$. In addition, the same allele should be at a frequency smaller than $w$ in an outgroup panel $A$ and larger than $x$ in a target panel $B$ (supplementary fig. S1, Supplementary Material online). For example, if we were interested in looking for Neanderthal-specific AI, we could set $y = 100\%$ and $z = 0\%$, to find alleles uniquely shared with Neanderthal, but not Denisova. If we were interested in archaic alleles shared with both Neanderthal and Denisova, we could set $y = 100\%$ and $z = 100\%$.

Another statistic that we found to be useful for finding AI events is $Q95_{A,B,C}(w, y)$, and is here defined as the 95th percentile of derived frequencies in an admixed sample $B$ of all SNPs that have a derived allele frequency $y$ in the archaic sample $C$, but where the derived allele is at a frequency smaller than $w$ in a sample $A$ of a nonadmixed population (fig. 1). For example, $Q95_{AFR,CHB,Nea}(1\%, 100\%) = 0.65$ means that if one computes the 95% quantile of all the Han Chinese derived allele frequencies of SNPs where the Neanderthal genome is homozygous derived and the derived allele has frequency smaller than 1% in Africans, that quantile will be equal to 0.65. The motivation for this statistic is that AI will produce archaic SNPs at high frequencies in the introgressed population. The 95th percentile should be an effective way of summarizing the frequencies of these SNPs while downweighting other SNPs that may also share the same allelic state as the archaic genomes, but that are segregating at low frequencies in the target panel and are therefore not informative about AI. In other words, it is a summary of the allele frequency spectrum in the introgressed population, conditional on only looking at alleles uniquely shared with the source population and at low frequency in the nonadmixed population. As before, we can generalize this statistic if we have a sample $D$ from a second archaic population. Then, $Q95_{A,B,C,D}(w, y, z)$ is the 95th percentile of derived frequencies in the sample $B$ of all SNPs that have a derived allele frequency

$y$ in the archaic sample C and derived allele frequency $z$ in the archaic sample D, but where the derived allele is at a frequency smaller than $w$ in the sample A (supplementary fig. S1, Supplementary Material online).

A common statistic that is indicative of population variation—expected heterozygosity ($\pi$)— was previously found to be affected by archaic introgression in a serial founder model of human history (DeGiorgio et al. 2009). We measured $\pi$ as the average of $2 * p * (1 - p)$ over all sites in a window, where $p$ is the sample derived allele frequency in the introgressed population.

## Simulations

None of these statistics have been explicitly vetted under scenarios of AI so far, though the performance of $D$ and $f_D$ has been previously evaluated for detecting local introgression (Martin et al. 2015). Therefore, we aimed to test how each of the statistics described above performed in detecting AI in a 40 kb window. We chose this window size because the mean length of introgressed haplotypes in Prüfer et al. (2014) was 44,078 bp (supplementary information S13, Supplementary Material online) and because 40kb is well above the length needed to reject incomplete lineage sorting for regions with moderate recombination rates (Huerta-Sánchez et al. 2014). We began by simulating a three population tree in Slim (Messer 2013) with constant $N_e = 10,000$, mutation rate equal to $1.5 \times 10^{-8}$ per bp per generation, recombination rate equal to $10^{-8}$ per bp per generation, and split times emulating the African-Eurasian and Neanderthal-modern human split times (4,000 and 16,000 generations ago, respectively). We allowed for admixture between the most distantly diverged population and one of the closely related sister populations, at different rates: 2%, 10%, and 25% (fig. 2A). We use the lower (2%) rate to represent the Neanderthal genome-wide admixture into Eurasians, with Africans as the nonadmixed population. The higher (10% and 25%) rates are meant to represent cases when a researcher is focusing on a particular region of the genome that has some *a priori* evidence for having been introgressed, thus pushing the local probability of introgression to high values, even though the genome-wide rate may be lower. Under each of the three admixture rate scenarios, we simulated regions that were evolving neutrally, regions where the central SNP was under weak positive additive selection ($s = 0.01$) and regions with a central SNP under strong selection ($s = 0.1$). We required the selected allele to be fixed in the archaic population prior to introgression, but allowed the allele to rise or decrease in frequency in the introgressed population, as determined by the strength of selection, its probability of entering the introgressed population and its starting frequency after introgression. Supplementary figure S2, Supplementary Material online shows the distributions of frequencies of the selected alleles in the introgressed population in the present.

We also tested how the statistics perform at detecting adaptive introgression when the alternative model is not a neutral introgression model, but a neutral model with ancestral structure (fig. 2G). We followed a model described in Huerta-Sánchez et al. (2014) and simulated a population in which an African population splits from archaic humans

before Eurasians, but is allowed to exchange migrants with them. Afterwards, we split Eurasians and archaic humans. At that point, we stop the previous migration and only allow for migration between the Eurasian and African populations until the present, at double the previous rate. This is meant to generate loci where Eurasians and archaic humans share a more recent common ancestor with each other than with Africans, but because of ancient shared ancestry, not recent introgression. We simulated three scenarios, in which we set the per-generation ancient migration rate to be 0.01, 0.001, and 0.0001, respectively, and the recent migration rate to be 0.02, 0.002, and 0.0002, respectively. We call these the strong-, medium-, and weak-migration scenarios, respectively. The stronger the migration, the weaker the ancestral structure, as archaic-shared segments in Eurasians will tend to be removed by migration with Africans.

## Plotting Haplotype Structure

The *Haplostrips* software (Marnetto et al. in prep.) was used to produce plots of haplotypes at candidate regions for AI. This software displays each SNP within a predefined region as a column, and each row represents a phased haplotype. Each haplotype is labeled with a color that corresponds to the 1000 Genomes panel of its carrier individual. The haplotypes were first hierarchically clustered via the single agglomerative method based on Manhattan distances, using the *stats* library in R. The resulting dendrogram of haplotypes was then reordered by decreasing similarity to a reference sequence constructed so that it contains all the derived alleles found in the archaic genome (Altai Neanderthal or Denisova). The reordering is performed using the minimum distance method, so that haplotypes with more derived alleles shared with the archaic population are at the top of the plot. Derived alleles are represented as black spots and ancestral alleles are represented as white spots. Variant positions were filtered out when the site in the archaic genome had mapping quality less than 30 or genotype quality less than 40, or if the minor allele had a population frequency smaller than 5% in each of the present-day human populations included in the plot.

## Hidden Markov Model

As haplotypes could look archaic simply because of ancestral structure or incomplete lineage sorting, we used a Hidden Markov Model (HMM) described in Seguin-Orlando et al. (2014) (which assumes an exponential distribution of admixture tract lengths [Pool and Nielsen 2009; Gravel 2012]), in order to verify that our candidate regions truly had archaic introgressed segments. This procedure also allowed us to confirm which of the archaic genomes was closest to the original source of introgression, as using a distant archaic source as input (e.g., the Denisova genome when the true source is closest to the Neanderthal genome) produced shorter or less frequent inferred segments in the HMM output than when using the closer source genome.

The HMM we used requires us to specify a prior for the admixture rate. We tried two priors: 2% and 50%. The first was chosen because it is consistent with the genome-wide admixture rate for Neanderthals into Eurasians. The second,

larger, value was chosen because each candidate region should *a priori* have a larger probability of being admixed, as they were found using statistics that are indicative of admixture in the first place. We observe almost no differences in the number of haplotypes inferred using either value. For example, for BNC2—a well-known candidate for AI (Sankararaman et al. 2014; Vernot and Akey 2014)—the frequency of sequences in EUR with inferred introgressed haplotypes under the 2% prior is 92.5%, whereas it is 93.1% under the 50% prior. However, the larger prior leads to longer and less fragmented introgressed chunks, as the HMM is less likely to transition into a nonintrogressed state between two introgressed states. Therefore, all figures we show below were obtained using a 50% admixture prior. The admixture time was set to 1,900 generations ago and the recombination rate parameter was set to the local recombination rate in each region, following the recombination rate map in Myers et al. (2005). A tract was called as introgressed if the posterior probability for introgression was higher than 90%. Under these parameters, the HMM has a specificity of 99.56%, a sensitivity of 36.07% and a false discovery rate of 1.15%.

## Supplementary Material

Supplementary figures S1–S61 and tables S1–S3 are available at *Molecular Biology and Evolution* online.

## Acknowledgments

## References

Akira S, Uematsu S, Takeuchi O. 2006. Pathogen recognition and innate immunity. *Cell* 124:783–801.

Aslanidis C, Ries S, Fehringer P, Büchler C, Klima H, Schmitz G. 1996. Genetic and biochemical evidence that CESD and wolman disease are distinguished by residual lysosomal acid lipase activity. *Genomics* 33:85–93.

Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO, Marchini JL, McCarthy S, McVean GA, Abecasis GR. 2015. A global reference for human genetic variation. *Nature* 526:68–74.

Barreiro LB, Ben-Ali M, Quach H, Laval G, Patin E, Pickrell JK, Bouchier C, Tichit M, Neyrolles O, Gicquel B, et al. 2009. Evolutionary dynamics of human Toll-like receptors and their different contributions to host defense. *PLoS Genet.* 5:e1000562.

Barton NH. 1998. The effect of hitch-hiking on neutral genealogies. *Genet Res.* 72:123–133.

Beall CM, Cavalleri GL, Deng L, Elston RC, Gao Y, Knight J, Li C, Li JC, Liang Y, McCormack M, et al. 2010. Natural selection on EPAS1 (hif2$\alpha$) associated with low hemoglobin concentration in Tibetan highlanders. *Proc Natl Acad Sci U S A.* 107:11459–11464.

Berisa T, Pickrell JK. 2016. Approximately independent linkage disequilibrium blocks in human populations. *Bioinformatics* 32:283–285.

Bigham A, Bauchet M, Pinto D, Mao X, Akey JM, Mei R, Scherer SW, Julian CG, Wilson MJ, Herráez DL, et al. 2010. Identifying signatures of natural selection in Tibetan and Andean populations using dense genome scan data. *PLoS Genet.* 6:e1001116.,

Cabral A, Fischer DF, Vermeij WP, Backendorf C. 2003. Distinct functional interactions of human Skn-1 isoforms with Ese-1 during keratinocyte terminal differentiation. *J Biol Chem.* 278:17792–17799.

Candille S, Van Raamsdonk CD, Chen C, Kuijper S, Chen-Tsai Y, Russ A, Meijlink F, Barsh GS. 2004. Dorsoventral patterning of the mouse coat by Tbx15. *PLoS Biol.* 2:E3.

Carvill GL, Heavin SB, Yendle SC, McMahon JM, O'Roak BJ, Cook J, Khan A, Dorschner MO, Weaver M, Calvert S, et al. 2013. Targeted resequencing in epileptic encephalopathies identifies de novo mutations in CHD2 and SYNGAP1. *Nat Genet.* 45:825–830.

Casey JP, McGettigan P, Lynam-Lennon N, McDermott M, Regan R, Conroy J, Bourke B, O'Sullivan J, Crushell E, Lynch S, et al. 2012. Identification of a mutation in LARS as a novel cause of infantile hepatopathy. *Mol Genet Metab.* 106:351–358.

Crow JF, Kimura M, et al. 1970. An introduction to population genetics theory. New York: Harper & Row.

Curry G. 1959. Genetical and developmental studies on droopy-eared mice. *J Embryol Exp Morphol.* 7:39–65.

Dannemann M, Andrés AM, Kelso J. 2016. Introgression of Neandertal- and Denisovan-like haplotypes contributes to adaptive variation in human toll-like receptors. *Am J Hum Genet.* 98:22–33.

DeGiorgio M, Jakobsson M, Rosenberg NA. 2009. Explaining worldwide patterns of human genetic variation using a coalescent-based serial founder model of migration outward from Africa. *Proc Natl Acad Sci U S A.* 106:16057–16062.

Deschamps M, Laval G, Fagny M, Itan Y, Abel L, Casanova JL, Patin E, Quintana-Murci L. 2016. Genomic signatures of selective pressures and introgression from archaic hominins at human innate immunity genes. *Am J Hum Genet.* 98:5–21.

Ding Q, Hu Y, Xu S, Wang J, Jin L. 2013. Neanderthal introgression at chromosome 3p21. 31 was under positive natural selection in East Asians. *Mol Biol Evol.* 31: 683–695.

Dong XY, Sun X, Guo P, Li Q, Sasahara M, Ishii Y, Dong JT. 2010. ATBF1 inhibits estrogen receptor (ER) function by selectively competing with AIB1 for binding to the ER in ER-positive breast cancer cells. *J Biol Chem.* 285:32801–32809.

Durand EY, Patterson N, Reich D, Slatkin M. 2011. Testing for ancient admixture between closely related populations. *Mol Biol Evol.* 28:2239–2252.

Fedetz M, Matesanz F, Caro-Maldonado A, Fernandez O, Tamayo J, Guerrero M, Delgado C, López-Guerrero J, Alcina A. 2006. OAS1 gene haplotype confers susceptibility to multiple sclerosis. *Tissue Antigens* 68:446–449.

Fumagalli M, Cagliani R, Riva S, Pozzoli U, Biasin M, Piacentini L, Comi GP, Bresolin N, Clerici M, Sironi M. 2010. Population genetics of IFIH1: ancient population structure, local selection, and implications for susceptibility to type 1 diabetes. *Mol Biol Evol.* 27:2555–2566.

Fumagalli M, Moltke I, Grarup N, Racimo F, Bjerregaard P, Jørgensen ME, Korneliussen TS, Gerbault P, Skotte L, Linneberg A, et al. 2015. Greenlandic Inuit show genetic signatures of diet and climate adaptation. *Science* 349:1343–1347.

Gburcik V, Cawthorn WP, Nedergaard J, Timmons JA, Cannon B. 2012. An essential role for Tbx15 in the differentiation of brown and "brite" but not white adipocytes. *Am J Physiol Endocrinol Metab.* 303:E1053–E1060.

Giles RE, Shimizu N, Ruddle FH. 1980. Assignment of a human genetic locus to chromosome 5 which corrects the heat sensitive lesion associated with reduced leucyl-tRNA synthetase activity in ts025Cl Chinese hamster cells. *Somatic Cell Genet.* 6:667–687.

Gravel S. 2012. Population genetics models of local ancestry. *Genetics* 191:607–619.

Green RE, Krause J, Briggs AW, Maricic T, Stenzel U, Kircher M, Patterson N, Li H, Zhai W, Fritz MHY, et al. 2010. A draft sequence of the Neandertal genome. *Science* 328:710–722.

Hackinger S, Kraaijenbrink T, Xue Y, Mezzavilla M, van Driem G, Jobling MA, de Knijff P, Tyler-Smith C, Ayub Q, et al. 2016. Wide distribution

and altitude correlation of an archaic high-altitude-adaptive EPAS1 haplotype in the Himalayas. *Hum Genet.* 135: 393–402.

Hamano E, Hijikata M, Itoyama S, Quy T, Phi NC, Long HT, Van Ban V, Matsushita I, Yanai H, Kirikae F, et al. 2005. Polymorphisms of interferon-inducible genes oas-1 and MxA associated with SARS in the Vietnamese population. *Biochem Biophys Res Commun.* 329:1234–1239.

Harris K, Nielsen R. 2015. The genetic cost of Neanderthal introgression. *bioRxiv* 030387.

Hašová M, Crhák T, Šafránková B, Dvořáková J, Muthnỳ T, Velebnỳ V, Kubala L. 2011. Hyaluronan minimizes effects of UV irradiation on human keratinocytes. *Arch Dermatol Res.* 303:277–284.

Hedrick PW. 2013. Adaptive introgression in animals: examples and comparison to new mutation and standing variation as sources of adaptive variation. *Mol Ecol.* 22:4606–4618.

Heid IM, Jackson AU, Randall JC, Winkler TW, Qi L, Steinthorsdottir V, Thorleifsson G, Zillikens MC, Speliotes EK, Mägi R, et al. 2010. Meta-analysis identifies 13 new loci associated with waist-hip ratio and reveals sexual dimorphism in the genetic basis of fat distribution. *Nat Genet.* 42:949–960.

Hill W, Robertson A. 1968. Linkage disequilibrium in finite populations. *Theor Appl Genet.* 38:226–231.

Hu CJ, Wang LY, Chodosh LA, Keith B, Simon MC. 2003. Differential roles of hypoxia-inducible factor 1α (HIF-1α) and HIF-2α in hypoxic gene regulation. *Mol Cell Biol.* 23:9361–9374.

Huerta-Sanchez E, Casey FP. 2015. Archaic inheritance: supporting high altitude life in Tibet. *J Appl Physiol.* 119:1129–1134.

Huerta-Sánchez E, Jin X, Bianba Z, Peter BM, Vinckenbosch N, Liang Y, Yi X, He M, Somel M, Ni P, et al. 2014. Altitude adaptation in Tibetans caused by introgression of Denisovan-like DNA. *Nature* 512:194–197.

Jacobs LC, Wollstein A, Lao O, Hofman A, Klaver CC, Uitterlinden AG, Nijsten T, Kayser M, Liu F. 2013. Comprehensive candidate gene study highlights UGT1A and BNC2 as new genes determining continuous skin color variation in Europeans. *Hum Genet.* 132:147–158.

Jeong C, Alkorta-Aranburu G, Basnyat B, Neupane M, Witonsky DB, Pritchard JK, Beall CM, Di Rienzo A. 2014. Admixture facilitates genetic adaptations to high altitude in Tibet. *Nat Commun.* 5:3281.

Juric I, Aeschbacher S, Coop G. 2015. The strength of selection against Neanderthal introgression. *bioRxiv* 030148.

Kamalati T, Jolin HE, Mitchell PJ, Barker KT, Jackson LE, Dean CJ, Page MJ, Gusterson BA, Crompton MR. 1996. Brk, a breast tumor-derived non-receptor protein-tyrosine kinase, sensitizes mammary epithelial cells to epidermal growth factor. *J Biol Chem.* 271:30956–30963.

Kato T, Sato H, Emi M, Seino T, Arawaka S, Iseki C, Takahashi Y, Wada M, Kawanami T. 2011. Segmental copy number loss of SFMBT1 gene in elderly individuals with ventriculomegaly: a community-based study. *Intern Med.* 50:297–303.

Khrameeva EE, Bozek K, He L, Yan Z, Jiang X, Wei Y, Tang K, Gelfand MS, Prufer K, Kelso J, et al. 2014. Neanderthal ancestry drives evolution of lipid catabolism in contemporary Europeans. *Nat Commun.* 5:3584.

Kim BY, Lohmueller KE. 2015. Selection and reduced population size cannot explain higher amounts of Neandertal ancestry in East Asian than in European human populations. *Am J Hum Genet.* 96:454–461.

Kim Y, Nielsen R. 2004. Linkage disequilibrium as a signature of selective sweeps. *Genetics* 167:1513–1524.

Kim Y, Stephan W. 2002. Detecting a local signature of genetic hitchhiking along a recombining chromosome. *Genetics* 160:765–777.

Klima H, Ullrich K, Aslanidis C, Fehringer P, Lackner K, Schmitz G. 1993. A splice junction mutation causes deletion of a 72-base exon from the mRNA for lysosomal acid lipase in a patient with cholesteryl ester storage disease. *J Clin Invest.* 92:2713–2718.

Knapp S, Yee L, Frodsham A, Hennig B, Hellier S, Zhang L, Wright M, Chiaramonte M, Graves M, Thomas H, et al. 2003. Polymorphisms in interferon-induced genes and the outcome of hepatitis C virus infection: roles of MxA, OAS-1 and PKR. *Genes Immun.* 4:411–419.

Krause M, Moradi J, Coleman S, D'Souza D, Liu C, Kronenberg M, Rowe D, Hawke T, Hadjiargyrou M. 2013. A novel GFP reporter mouse reveals Mustn1 expression in adult regenerating skeletal muscle,

activated satellite cells and differentiating myoblasts. *Acta Physiol.* 208:180–190.

Kronforst MR, Hansen ME, Crawford NG, Gallant JR, Zhang W, Kulathinal RJ, Kapan DD, Mullen SP. 2013. Hybridization reveals the evolving genomic architecture of speciation. *Cell Rep.* 5:666–677.

Lausch E, Hermanns P, Farin HF, Alanay Y, Unger S, Nikkel S, Steinwender C, Scherer G, Spranger J, Zabel B, et al. 2008. TBX15 mutations cause craniofacial dysmorphism, hypoplasia of scapula and pelvis, and short stature in cousin syndrome. *Am J Hum Genet.* 83:649–655.

Lewontin R. 1964. The interaction of selection and linkage. I. general considerations; heterotic models. *Genetics* 49:49.

Lim JK, Lisco A, McDermott DH, Huynh L, Ward JM, Johnson B, Johnson H, Pape J, Foster GA, Krysztof D, et al. 2009. Genetic variation in OAS1 is a risk factor for initial infection with West Nile virus in man. *PLoS Pathog.* 5:e1000321.

Lin S, Shen H, Li JL, Tang S, Gu Y, Chen Z, Hu C, Rice JC, Lu J, Wu L. 2013. Proteomic and functional analyses reveal the role of chromatin reader SFMBT1 in regulating epigenetic silencing and the myogenic gene program. *J Biol Chem.* 288:6238–6247.

Lin SE, Oyama T, Nagase T, Harigaya K, Kitagawa M. 2002. Identification of new human mastermind proteins defines a family that consists of positive regulators for notch signaling. *J Biol Chem.* 277:50612–50620.

Liu CT, Buchkovich ML, Winkler TW, Heid IM, Borecki I, Fox CS, Mohlke KL, North KE, Cupples LA, A. A. A. G. Consortium, et al. 2014. Multi-ethnic fine-mapping of 14 central adiposity loci. *Hum Mol Genet.* 23:4738–4744.

Liu CT, Monda KL, Taylor KC, Lange L, Demerath EW, Palmas W, Wojczynski MK, Ellis JC, Vitolins MZ, Liu S, et al. 2013. Genome-wide association of body fat distribution in African ancestry populations suggests new loci. *PLoS Genet.* 9:e1003681.

Liu S, Wang H, Jin Y, Podolsky R, Reddy MPL, Pedersen J, Bode B, Reed J, Steed D, Anderson S, et al. 2009. IFIH1 polymorphisms are significantly associated with type 1 diabetes and IFIH1 gene expression in peripheral blood mononuclear cells. *Hum Mol Genet.* 18:358–365.

Liu SY, Aliyari R, Chikere K, Li G, Marsden MD, Smith JK, Pernet O, Guo H, Nusbaum R, Zack JA, et al. 2013. Interferon-inducible cholesterol-25-hydroxylase broadly inhibits viral entry by production of 25-hydroxycholesterol. *Immunity* 38:92–105.

Lund EG, Kerr TA, Sakai J, Li WP, Russell DW. 1998. cDNA cloning of mouse and human cholesterol 25-hydroxylases, polytopic membrane proteins that synthesize a potent oxysterol regulator of lipid metabolism. *J Biol Chem.* 273:34316–34327.

Martin SH, Davey JW, Jiggins CD. 2015. Evaluating the use of ABBA–BABA statistics to locate introgressed loci. *Mol Biol Evol.* 32:244–257.

Mendez FL, Watkins JC, Hammer MF. 2013. Neandertal origin of genetic variation at the cluster of OAS immunity genes. *Mol Biol Evol.* 30:798–801.

Messer PW. 2013. SLiM: simulating evolution with selection and linkage. *Genetics* 194:1037–1039.

Meyer M, Kircher M, Gansauge MT, Li H, Racimo F, Mallick S, Schraiber JG, Jay F, Prüfer K, de Filippo C, et al. 2012. A high-coverage genome sequence from an archaic Denisovan individual. *Science* 338:222–226.

Mochida GH, Ganesh VS, de Michelena MI, Dias H, Atabay KD, Kathrein KL, Huang HT, Hill RS, Felie JM, Rakiec D, et al. 2012. CHMP1A encodes an essential regulator of BMI1-INK4A in cerebellar development. *Nat Genet.* 44:1260–1264.,

Myers S, Bottolo L, Freeman C, McVean G, Donnelly P. 2005. A fine-scale map of recombination rates and hotspots across the human genome. *Science* 310:321–324.

Pallares LF, Carbonetto P, Gopalakrishnan S, Parker CC, Ackert-Bicknell CL, Palmer AA, Tautz D. 2015. Mapping of craniofacial traits in outbred mice identifies major developmental genes involved in shape determination. *PLOS Genet.* 11:e1005607.

Peng Y, Yang Z, Zhang H, Cui C, Qi X, Luo X, Tao X, Wu T, Chen H, Shi H, et al. 2011. Genetic variations in Tibetan populations and high-altitude adaptation at the Himalayas. *Mol Biol Evol.* 28:1075–1081.

Persson B, Kallberg Y, Bray JE, Bruford E, Dellaporta SL, Favia AD, Duarte RG, Jörnvall H, Kavanagh KL, Kedishvili N, et al. 2009. The SDR (short-chain dehydrogenase/reductase and related enzymes) nomenclature initiative. *Chem Biol Interact.* 178:94–98.

Pickrell JK, Coop G, Novembre J, Kudaravalli S, Li JZ, Absher D, Srinivasan BS, Barsh GS, Myers RM, Feldman MW, et al. 2009. Signals of recent positive selection in a worldwide sample of human populations. *Genome Res.* 19:826–837.

Pool JE, Nielsen R. 2009. Inference of historical changes in migration rate from the lengths of migrant tracts. *Genetics* 181:711–719.

Prüfer K, Racimo F, Patterson N, Jay F, Sankararaman S, Sawyer S, Heinze A, Renaud G, Sudmant PH, de Filippo C, et al. 2014. The complete genome sequence of a Neanderthal from the Altai mountains. *Nature* 505:43–49.

Pruitt KD, Harrow J, Harte RA, Wallin C, Diekhans M, Maglott DR, Searle S, Farrell CM, Loveland JE, Ruef BJ, et al. 2009. The consensus coding sequence (CCDS) project: identifying a common protein-coding gene set for the human and mouse genomes. *Genome Res.* 19:1316–1323.

Qu HQ, Marchand L, Grabs R, Polychronakos C. 2008. The association between the IFIH1 locus and type 1 diabetes. *Diabetologia* 51:473–475.

Racimo F, Gokhman D, Fumagalli M, Hansen T, Moltke I, Albrechtsen A, Carmel L, Huerta-Sanchez E, Nielsen R. 2016. Archaic adaptive introgression in TBX15/WARS2. *bioRxiv* 033928.

Racimo F, Sankararaman S, Nielsen R, Huerta-Sánchez E. 2015. Evidence for archaic adaptive introgression in humans. *Nat Rev Genet.* 16:359–371.

Rauch A, Wieczorek D, Graf E, Wieland T, Endele S, Schwarzmayr T, Albrecht B, Bartholdi D, Beygo J, Di Donato N, et al. 2012. Range of genetic mutations associated with severe non-syndromic sporadic intellectual disability: an exome sequencing study. *Lancet* 380:1674–1682.

Reich D, Green RE, Kircher M, Krause J, Patterson N, Durand EY, Viola B, Briggs AW, Stenzel U, Johnson PL, et al. 2010. Genetic history of an archaic hominin group from Denisova cave in Siberia. *Nature* 468:1053–1060.

Sankararaman S, Mallick S, Dannemann M, Prüfer K, Kelso J, Pääbo S, Patterson N, Reich D. 2014. The genomic landscape of Neanderthal ancestry in present-day humans. *Nature* 507:354–357.

Sankararaman S, Mallick S, Patterson N, Reich D. 2016. The combined landscape of Denisovan and Neanderthal ancestry in present-day humans. *Curr Biol.* 26: 1241–1247.

Seguin-Orlando A, Korneliussen TS, Sikora M, Malaspinas AS, Manica A, Moltke I, Albrechtsen A, Ko A, Margaryan A, Moiseyev V, et al. 2014. Genomic structure in Europeans dating back at least 36,200 years. *Science* 346:1113–1118.

Shibata N, Kawarai T, Lee JH, Lee HS, Shibata E, Sato C, Liang Y, Duara R, Mayeux RP, St George-Hyslop PH, et al. 2006. Association studies of cholesterol metabolism genes (CH25H, ABCA1 and CH24H) in Alzheimer's disease. *Neurosci Lett.* 391:142–146.

Shungin D, Winkler TW, Croteau-Chonka DC, Ferreira T, Locke AE, Mägi R, Strawbridge RJ, Pers TH, Fischer K, Justice AE, et al. 2015. New genetic loci link adipose and insulin biology to body fat distribution. *Nature* 518:187–196.

Sigma Type 2 Diabetes Consortium. 2014. Sequence variants in SLC16A11 are a common risk factor for type 2 diabetes in Mexico. *Nature* 506:97–101.

Singh MK, Petry M, Haenig B, Lescher B, Leitges M, Kispert A. 2005. The t-box transcription factor Tbx15 is required for skeletal development. *Mech Dev.* 122:131–144.

Skoglund P, Jakobsson M. 2011. Archaic human ancestry in East Asia. *Proc Natl Acad Sci U S A.* 108:18301–18306.

Smith J, Kronforst MR. 2013. Do Heliconius butterfly species exchange mimicry alleles? *Biol Lett.* 9:20130503.

Spilker C, Kreutz MR. 2010. Rapgaps in brain: multipurpose players in neuronal Rap signalling. *Eur J Neurosci.* 32:1–9.

Sun X, Frierson HF, Chen C, Li C, Ran Q, Otto KB, Cantarel BM, Vessella RL, Gao AC, Petros J, et al. 2005. Frequent somatic mutations of the transcription factor ATBF1 in human prostate cancer. *Nat Genet.* 37:407–412.

Surapureddi S, Yu S, Bu H, Hashimoto T, Yeldandi AV, Kashireddy P, Cherkaoui-Malki M, Qi C, Zhu YJ, Rao MS, et al. 2002. Identification of a transcriptionally active peroxisome proliferator-activated receptor α-interacting cofactor complex in rat liver and characterization of PRIC285 as a coactivator. *Proc Natl Acad Sci U S A.* 99:11836–11841.

Takemoto H, Tamai K, Akasaka E, Rokunohe D, Takiyoshi N, Umegaki N, Nakajima K, Aizu T, Kaneko T, Nakano H, et al. 2010. Relation between the expression levels of the POU transcription factors Skn-1a and Skn-1n and keratinocyte differentiation. *Jo Dermatol Sci.* 60:203–205.

Tomaru T, Satoh T, Yoshino S, Ishizuka T, Hashimoto K, Monden T, Yamada M, Mori M. 2006. Isolation and characterization of a transcriptional cofactor and its novel isoform that bind the deoxyribonucleic acid-binding domain of peroxisome proliferator-activated receptor-γ. *Endocrinology* 147:377–388.

Vanhoutteghem A, Djian P. 2006. Basonuclins 1 and 2, whose genes share a common origin, are proteins with widely different properties and functions. *Proc Natl Acad Sci U S A.* 103:12423–12428.

Vernot B, Akey JM. 2014. Resurrecting surviving Neandertal lineages from modern human genomes. *Science* 343:1017–1021.

Vernot B, Akey JM. 2015. Complex history of admixture between modern humans and Neandertals. *Am J Hum Genet.* 96:448–453.

Vernot B, Tucci S, Kelso J, Schraiber JG, Wolf AB, Gittelman RM, Dannemann M, Grote S, McCoy RC, Norton H, et al. 2016. Excavating Neandertal and Denisovan DNA from the genomes of Melanesian individuals. *Science.* DOI: 10.1126/science.aad9416.

Voight BF, Kudaravalli S, Wen X, Pritchard JK. 2006. A map of recent positive selection in the human genome. *PLoS Biol.* 4:e72.

Von Holstein SL, Fehr A, Heegaard S, Therkildsen MH, Stenman G. 2012. CRTC1-MAML2 gene fusion in mucoepidermoid carcinoma of the lacrimal gland. *Oncol Rep.* 27:1413–1416.

Vrijenhoek T, Buizer-Voskamp JE, van der Stelt I, Strengman E, Risk G, Sabatti C, van Kessel AG, Brunner HG, Ophoff RA, Veltman JA, et al. 2008. Recurrent CNVs disrupt three candidate genes in schizophrenia patients. *Am J Hum Genet.* 83:504–510.

Wall JD, Yang MA, Jay F, Kim SK, Durand EY, Stevison LS, Gignoux C, Woerner A, Hammer MF, Slatkin M. 2013. Higher levels of Neanderthal ancestry in East Asians than in Europeans. *Genetics* 194:199–209.

Wang B, Zhang YB, Zhang F, Lin H, Wang X, Wan N, Ye Z, Weng H, Zhang L, Li X, et al. 2011. On the origin of Tibetans and their genetic basis in adapting high-altitude environments. *PLoS One* 6:e17002.

Wang KS, Liu XF, Aragam N. 2010. A genome-wide meta-analysis identifies novel loci associated with schizophrenia and bipolar disorder. *Schizophr Res.* 124:192–199.

Warner TG, Dambach LM, Shin JH, O'Brien JS. 1980. Separation and characterization of the acid lipase and neutral esterases from human liver. *Am J Hum Genet* 32:869.

Wilson PM, Fryer RH, Fang Y, Hatten ME. 2010. Astn2, a novel member of the astrotactin gene family, regulates the trafficking of ASTN1 during glial-guided neuronal migration. *J Neurosci.* 30:8529–8540.

Winnes M, Mölne L, Suurküla M, Andrén Y, Persson F, Enlund F, Stenman G. 2007. Frequent fusion of the CRTC1 and MAML2 genes in clear cell variants of cutaneous hidradenomas. *Genes Chromosomes Cancer* 46:559–563.

Wood DS, Zeviani M, Prelle A, Bonilla E, Salviati G, Miranda AF, DiMauro S, Rowland LP. 1987. Is nebulin the defective gene product in Duchenne muscular dystrophy? *N Engl J Med.* 1987:107–108.

Woodage T, Basrai MA, Baxevanis AD, Hieter P, Collins FS. 1997. Characterization of the CHD family of proteins. *Proc Natl Acad Sci U S A.* 94:11472–11477.

Wu L, Liu J, Gao P, Nakamura M, Cao Y, Shen H, Griffin JD. 2005. Transforming activity of MECT1-MAML2 fusion oncoprotein is mediated by constitutive CREB activation. *EMBO J.* 24:2391–2402.

Xu S, Li S, Yang Y, Tan J, Lou H, Jin W, Yang L, Pan X, Wang J, Shen Y, et al. 2011. A genome-wide search for signals of high-altitude adaptation in Tibetans. *Mol Biol Evol.* 28:1003–1011.

Yi X, Liang Y, Huerta-Sanchez E, Jin X, Cuo ZXP, Pool JE, Xu X, Jiang H, Vinckenbosch N, Korneliussen TS, et al. 2010. Sequencing of 50 human exomes reveals adaptation to high altitude. *Science* 329:75–78.

Yotova V, Lefebvre JF, Moreau C, Gbeha E, Hovhannesyan K, Bourgeois S, Bédarida S, Azevedo L, Amorim A, Sarkisian T, et al. 2011. An X-linked haplotype of neandertal origin is present among all non-African populations. *Mol Biol Evol.* 28:1957–1962.