

Balancing selection on genomic deletion polymorphisms in humans

Authors: Alber Aqil¹, Leo Speidel^{2,3}, Pavlos Pavlidis⁴, Omer Gokcumen¹

Affiliations:

1. Department of Biological Sciences, University at Buffalo, Buffalo, NY. USA.

2. University College London, Genetics Institute, London, UK.

3. The Francis Crick Institute, London, UK.

4. Institute of Computer Science (ICS), Foundation of Research and Technology-Hellas, Heraklion, Crete, Greece.

Correspondence:

Omer Gokcumen, gokcumen@gmail.com

Abstract:

A key question in biology is why genomic variation persists in a population for extended periods. Recent studies have identified examples of genomic deletions that have remained polymorphic in the human lineage for hundreds of millennia, ostensibly owing to balancing selection. Nevertheless, genome-wide investigation of ancient and possibly adaptive deletions remains imperative. Here, we demonstrate an excess of polymorphisms in present-day humans that predate the modern human-Neanderthal split (ancient polymorphisms), which cannot be explained solely by selectively neutral scenarios. We analyze the adaptive mechanisms that underlie this excess in deletion polymorphisms. Using a previously published measure of balancing selection, we show that this excess of ancient deletions is largely owing to balancing selection. Based on the absence of signatures of overdominance, we conclude that it is a rare mode of balancing selection among ancient deletions. Instead, more complex scenarios involving spatially and temporally variable selective pressures are likely more common mechanisms. Our results suggest that balancing selection resulted in ancient deletions harboring disproportionately more exonic variants with GWAS associations. We further found that ancient deletions are significantly enriched for traits related to metabolism and immunity. As a by-product of our analysis, we show that deletions are, on average, more deleterious than single-nucleotide variants. We can now argue that not only is a vast majority of common variants shared among human populations, but a considerable portion of biologically relevant variants has been segregating among our ancestors for hundreds of thousands, if not millions, of years.

36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86

INTRODUCTION

The evolutionary forces that shape the allele frequency distribution of functional genetic variants remain a hotly debated issue. In humans, tens of thousands of common variants are reported to be associated with human diseases (Loos, 2020). However, the mainstream view remains that the majority of these functional genetic variants have had a negligible effect on reproductive fitness and that the frequency of these variants has fluctuated neutrally by drift over time (Bromberg et al., 2013; Dudley et al., 2012). Functional variants that have measurable fitness effects are often observed at a low frequency (Eyre-Walker, 2010). These low-frequency functional variants are considered to be in the process of being eliminated from the population by negative selection (Gibson, 2018; Lettre, 2014; Zeng et al., 2018). Nevertheless, an increasing number of studies are showing that more complex evolutionary histories (Benton et al., 2021; Mathieson and Mathieson, 2018) involving introgression from archaic hominins (McArthur et al., 2020), geography-specific adaptation (Hamid et al., 2021; Lachance and Tishkoff, 2013; Mendoza-Revilla et al., 2021), negative selection (Zeng et al., 2018), and polygenic selection (Barghi et al., 2020; Berg and Coop, 2014; Pritchard et al., 2010; Sella and Barton, 2019) may explain the allele frequencies of variants associated with complex diseases. In this context, we aim to test the hypothesis that balancing selection is a considerable force in shaping the allele frequencies of extant functional deletions in the human genome.

Balancing selection is a mode of natural selection that maintains a genomic polymorphism by overcoming the stochastic loss or fixation of one of the alleles by genetic drift (Fijarczyk and Babik, 2015; Fisher, 1922; Noonan et al., 2006). H.J. Muller was the first to discover balancing selection from his study of balanced lethals in *Drosophila* (Muller, 1918). Adaptive variational maintenance by balancing selection may be achieved in a number of ways. In a mechanism known as over-dominance (also called heterozygote advantage), the individual who is heterozygous for a certain variant has a higher fitness (Fisher, 1922; Wallace, 1970). In negative frequency-dependent selection, rarer variants confer higher fitness. This leads to a fluctuation of a variant's frequency in the population until an equilibrium is established, such that neither variant confers an advantage relative to the other (Smith Maynard et al., 1998; Takahashi and Kawata, 2013). Temporally varying selection, wherein the selection coefficient associated with an allele changes over time, can lead to the oscillation of this allele's frequency over time (Abdul-Rahman et al., 2021; Wittmann et al., 2017). Spatially varying selection, wherein the selection coefficient associated with an allele varies across geography, may fix this allele locally in one niche and eliminate it in another, leading to the global persistence of variation at the locus (Hedrick, 2006; Levene, 1953; Saitou et al., 2021a).

Unlike positive and negative selection, there is only a modest number of well-established instances of balancing selection (Charlesworth and Charlesworth 2016). In humans, these include polymorphisms of the *ABO* gene, which determine the A, B, and O blood groups (Ségurel et al. 2012), and polymorphisms in the major histocompatibility complex, which encodes cell-surface glycoproteins that display samples of peptides from within the cell on the cell's surface (Takahata et al. 1992). Two variants of *ERAP2*, which too is involved in the antigen-presenting pathway, have also been maintained under balancing selection (Andrés et al., 2010; Klunk et al., 2022). The classic example of recent, shorter-term balancing selection in humans is the maintenance (by over-dominance) of sickle-cell alleles at the β -globin locus in the regions of Africa where malaria is endemic (Allison, 1954a a, 1954b b; Hedrick, 2011). Similar reasoning applies to certain α -thalassemia alleles in parts of Southeast Asia where malaria is widespread (Qiu et al. 2013). In fact, the higher fitness of heterozygotes for thalassemia alleles in malaria-struck regions was presciently predicted by Haldane in 1949 (Lederberg, 1999). In the realm of structural variants, complex copy number variation in the human salivary agglutinin

87 genes (Alharbi et al., 2022), a regulatory deletion upstream of APOBEC3 gene family
88 (Gokcumen et al., 2013), and a deletion spanning *LCE3B* and *LCE3C* (Pajic et al., 2016), which
89 is associated with psoriasis, have been explicitly argued to be evolving in the human lineage
90 under balancing selection.

91
92 So far, most systematic investigations into balancing selection in modern humans have focused
93 primarily on genes (DeGiorgio et al., 2014; Soni et al., 2022) and on single nucleotide variants
94 (SNVs) (Siewert and Voight, 2020); (Bitarello et al., 2018; Siewert and Voight, 2017).

95 Additionally, some studies have focused exclusively on “long-term” balancing selection wherein
96 variants have been maintained in the human lineage since before the split from the chimpanzee
97 clade (Leffler et al., 2013). Others have focused on short-term or population-specific balancing
98 selection (Hedrick, 2011; Qiu et al., 2013). We set out to identify potential targets of balancing
99 selection that are structural in nature and that may not have been captured by earlier studies.

100 Thus, we concentrate our efforts on autosomal deletion polymorphisms (> 50bp) that have been
101 maintained in the human lineage since before the split, approximately ~700,000 years ago, of
102 anatomically modern humans (AMHs) from the lineage that led to both Neanderthals and
103 Denisovans (henceforth, collectively referred to as archaic hominins). In this study, we will use
104 the term “ancient polymorphisms” to refer to such polymorphisms. Focusing on such “medium-
105 term” balancing selection will likely allow us to capture more potential targets than could an
106 exclusive study of either “long-term” or “short-term” balancing selection. Moreover, deletions are
107 interesting in the context of selection: since a given deletion affects more nucleotides than a
108 single nucleotide variant (SNV), if a defined genomic window is indeed of adaptive importance,
109 deletions may have more profound functional consequences (Conrad et al., 2010; Saitou and
110 Gokcumen, 2020). Such functional outcomes may translate into non-trivial selection coefficients
111 either for or against the deletion. Additionally, deletions are relatively easier to both genotype
112 and analyze than are other structural variants, making an evolutionary analysis involving
113 deletions tractable.

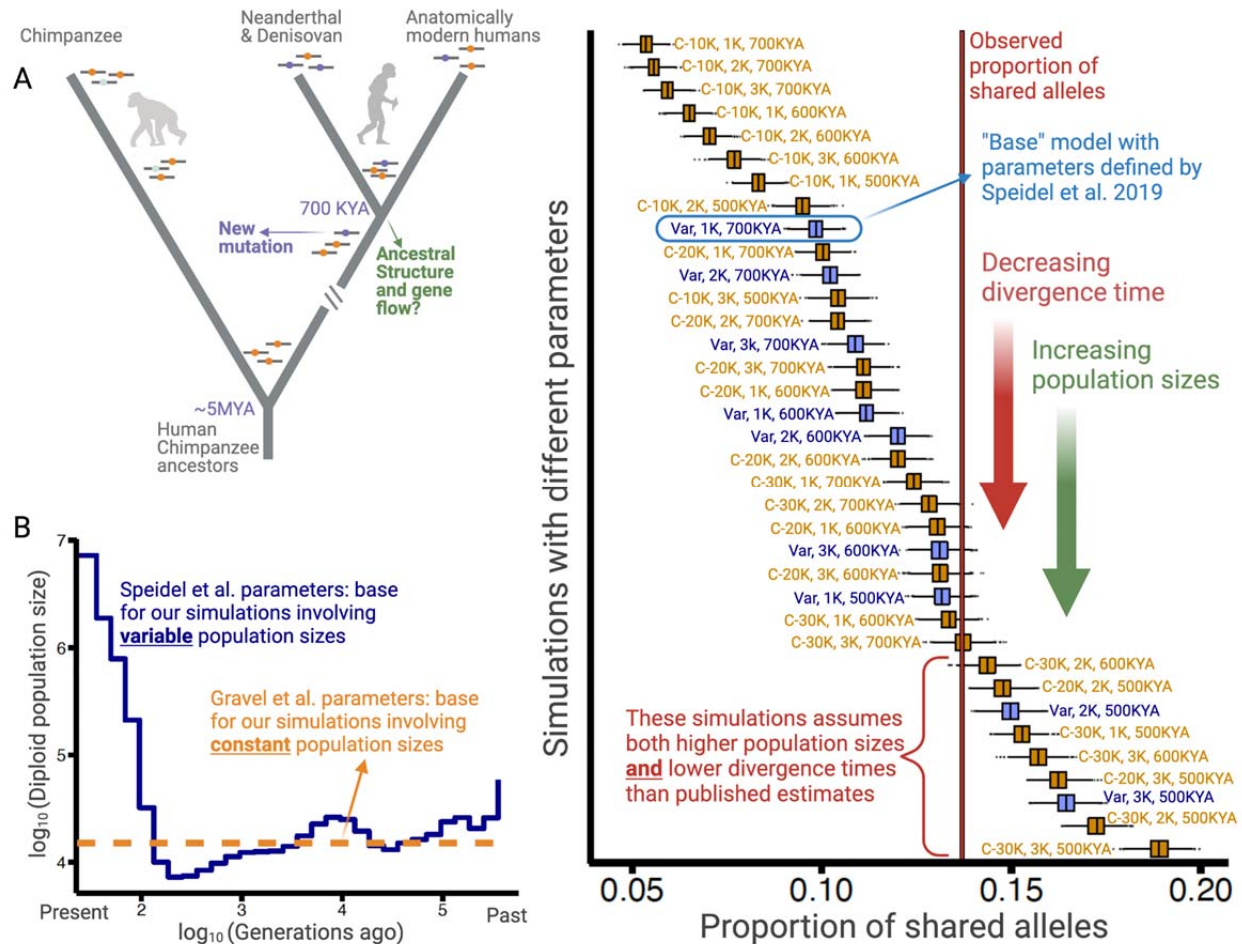
114

115 **RESULTS AND DISCUSSION**

116 **AMH exhibits a greater proportion of ancient polymorphisms than expected under** 117 **adaptive neutrality**

118 Older polymorphisms may be more likely than newer ones to exhibit signatures associated with
119 balancing selection because they have survived stochastic fixation or elimination for extended
120 periods. It is, therefore, possible that a certain proportion of human polymorphisms that are
121 older than the AMH-archaic split (~700,000 years), *i.e.*, ancient polymorphisms, (**Figure 1A**)
122 have been maintained by balancing selection. We tested this hypothesis by comparing the
123 proportion of ancient polymorphisms segregating in AMHs to a neutrally expected distribution of
124 this proportion. If this proportion is significantly higher in the observed data than under the
125 neutrally simulated data, and if we can reject other plausible explanations for difference, we can
126 conclude that some of the polymorphisms older than 700,000 years may have been maintained
127 by balancing selection.

128



129
130
131
132
133
134
135
136
137
138
139
140
141
142

Figure 1. Excess of ancient polymorphisms segregating in AMHs. **A.** A schematic representation of derived “ancient” variants (purple) that emerged before the AMH-archaic hominin divergence (and after hominin-chimp divergence), and have remained polymorphic in the AMH lineage. The ancestral variants are indicated as orange, and the derived chimpanzee-specific variants are indicated in light blue. **B.** The Speidel et al. and Gravel et al. simulation parameters. Speidel et al. provide parameters that involve varying population sizes for the YRI population. **C.** Expected distribution of the proportion of ancient polymorphisms in YRI under different models. Each distribution is labeled with three parameters in the form (*AMH- N_e* , *Archaic- N_e* , *time since archaic-AMH divergence*). The simulations where we used variable effective population size published by Speidel et al. are indicated by blue color and labeled “Var”. The simulations where *AMH- N_e* is constant are shown in orange, and provide the population size used. The vertical line represents the empirical proportion of ancient polymorphisms in YRI.

143
144
145
146
147
148
149
150
151
152
153

For this test, we focused on 28,291 randomly chosen SNVs (minor allele count > 1) in the Yoruba (YRI) population (1000 Genomes Project Consortium et al., 2015); see **Supplementary material** for a discussion of our rationale behind using SNVs rather than deletion polymorphisms, the variant class of our interest). We focused on random SNVs instead of all the SNVs in order to mitigate the biases that would be introduced due to linkage. A variant was classified as ancient if the derived allele was shared, by common descent, with at least one of the four high-coverage archaic hominin genomes (three Neanderthals and one Denisovan) (Mafessoni et al., 2020; Meyer et al., 2012; Prüfer et al., 2017, 2014). We found that the ancient SNVs make up 13.7% (3,894 SNVs) of the total. Note that we removed the recurrent SNVs from our analysis using a linkage-disequilibrium based approach (see **Methods**). To compare the proportion of ancient SNVs against neutral expectations, we used *ms* (Hudson, 2002) to

154 produce 2,000 runs of 20,000 neutrally simulated variants in the Yoruba population (see
155 **Methods** for details). Thereupon, in each run, we calculated the proportion of variants shared
156 with archaic hominins, producing a distribution of the expected proportion of ancient
157 polymorphisms under neutrality.

158
159 To ensure that our analysis is not biased by the idiosyncrasies of any particular model, we
160 performed these simulations using 36 distinct models. The models vary by three parameters: N_e
161 of Yoruba/AMH, N_e of archaic hominins, and the time of divergence between the AMH and the
162 archaic hominin lineage. The N_e for humans can be either constant (ranging from 10,000 to
163 30,000) or varying over time (Speidel et al., 2019) (**Figure 1B**). N_e of archaic hominins ranges
164 from 1,000 to 3,000; and the divergence time ranges from 500 to 700 kya (Bergström et al.,
165 2021; Mafessoni et al., 2020; Meyer et al., 2012; Prüfer et al., 2014). In the main text, we focus
166 on the model with variable AMH- N_e (Speidel et al., 2019), using the well-accepted archaic
167 hominin N_e of 1,000 and a divergence time of 700 kya. We refer to this as the “base model”.

168
169 Using this model, we find that the entire simulated distribution of the proportion of ancient
170 polymorphisms lies to the left of the empirical value of 13.7% (**Figure 1C**). These results hold
171 for all other models with realistic sets of parameters. Even when we change any single
172 parameter to unrealistic levels (e.g., either AMH- N_e = 30,000, or divergence time = 500 kya), we
173 still observe an excess of ancient polymorphisms. Neutral models can explain the empirical
174 proportion of ancient polymorphisms only when at least two parameters are tweaked in a less
175 realistic direction (e.g., both AMH- N_e = 30,000 and divergence time = 500 kya). Therefore, we
176 conclude that the high proportion (13.7%) of ancient polymorphisms cannot be explained by
177 realistic neutral scenarios.

178
179 There is a possibility that we are observing an empirical excess of ancient polymorphisms owing
180 to a high incidence of recurrent mutations in the AMH and archaic hominin lineages that
181 remained undetected by the linkage-disequilibrium-based approach we used to identify them.
182 Since CpG sites may be particularly prone to recurrent mutations, we calculated the proportion
183 of empirical ancient polymorphisms again using only A↔T SNVs. This analysis yielded a
184 proportion of 14.24%, not much different from the previously calculated 13.7%. Moreover, if
185 recurrence was the main cause of the observed excess of ancient polymorphisms, we would
186 expect this excess (relative to the simulated distribution) to be more pronounced among
187 polymorphisms with low derived allele frequencies. To test if this is the case, we repeated our
188 analysis using the base model for simulations, dividing the empirical and simulated SNVs into
189 derived allele frequency bins (**Figure 1—figure supplement 1**). We observed that the excess
190 of ancient polymorphisms is, in fact, most pronounced at high derived allele frequency. Both
191 these results combined suggest that our results are not biased due to undetected recurrent
192 mutations.

193
194 Next, we consider possible non-adaptive explanations for the observed excess of ancient
195 polymorphisms. First, we considered scenarios invoking structure in the population that was
196 ancestral to both AMHs and archaic hominins, while allowing gene flow between the latent
197 subgroups within the ancestral population (**Figure 1**; see **Methods** for details). We found that
198 the excess of ancient polymorphisms can be explained by structuring the ancestral population

199 into 3 distinct subpopulations, such that the fraction of each subgroup formed by the migrants of
200 each of the other subgroups, every generation, is below 0.0075% (**Figure 1—figure**
201 **supplement 2A-B**). However, the allele frequency spectrum for SNVs simulated with ancient
202 population structure significantly deviates from the observed allele frequency spectrum in that
203 the former overestimates the intermediate/common variants (**Figure 1—figure supplement**
204 **2C**). Therefore, invoking such structure to explain the excess of ancient polymorphisms may be
205 unrealistic. Another possible explanation comes from the evidence of introgression from early
206 modern human ancestors into Neanderthals to the exclusion of Denisovans (Posth et al., 2017).
207 Such admixture can increase the apparent proportion of ancient polymorphisms due to elevated
208 allele sharing with Neanderthals. However, we do not observe a higher proportion of derived
209 allele sharing (by common descent) with Denisovans than with Neanderthals. In fact, the
210 proportion of derived alleles shared with the Denisovan (9.75%) and the Altai Neanderthal
211 (10.28%) is higher than the proportion shared with Chagyrskaya and Vindija Neanderthals
212 (6.14% each), which is incompatible with such introgression as the prime cause of excessive
213 ancient polymorphisms in AMHs. We note that the differential allele sharing with the Denisovan
214 and Altai Neanderthal on one hand, and the Vindija and Chagyrskaya Neanderthal on the other
215 would be an interesting subject for future studies.

216
217 Overall, based on our current knowledge of ancient interactions and demographic history, our
218 analyses implicate balancing selection as a possible cause of the excess of observed ancient
219 polymorphisms. Next, we focus on deletion polymorphisms segregating in AMHs, categorize
220 them based on their evolutionary histories, and test whether ancient deletion polymorphisms are
221 enriched for targets of balancing selection.

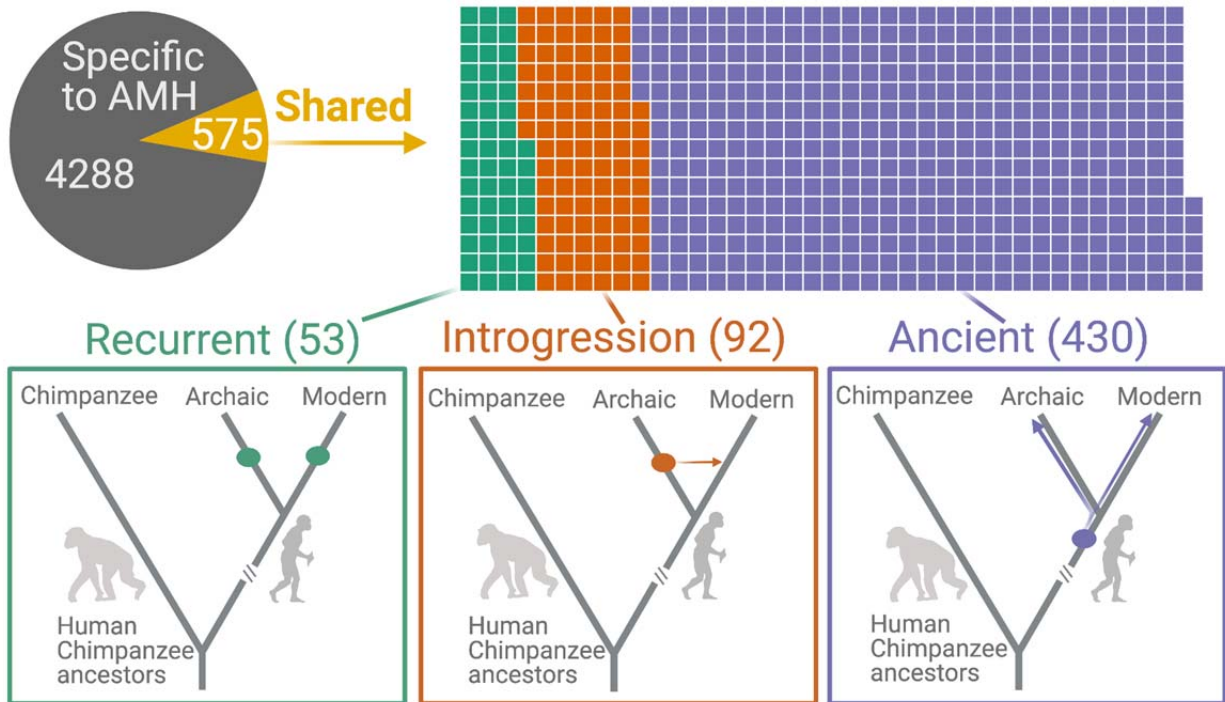
222 223 **Categorizing human deletions based on their evolutionary history**

224 Having established that AMHs exhibit an excess of ancient polymorphisms that cannot be
225 explained solely by non-adaptive causes, we identify ancient deletion polymorphisms among
226 AMHs. Since the vast majority of deletions in AMHs are derived relative to chimpanzees
227 (**Supplementary material**), this could be accomplished by identifying AMH deletions that are
228 shared with archaic hominins by common descent.

229
230 In this analysis, AMHs are represented by the YRI (Yoruba), CEU (Utah residents with Northern
231 and Western European ancestry), and CHB (Han Chinese in Beijing) from 1000 Genomes
232 Project Phase-3 dataset (1000 Genomes Project Consortium et al., 2015); and archaic hominins
233 are represented by the four available high-coverage (~30X) archaic hominin genomes
234 (Mafessoni et al., 2020; Meyer et al., 2012; Prüfer et al., 2017, 2014). Our choice of AMH
235 populations was guided by our wish to both sample from different regions, and use relatively
236 well-studied populations. We genotyped all AMH deletions in the archaic hominin genomes
237 using a read depth based pipeline (**Supplementary File 1**). We considered a deletion “shared”
238 if it was identified in at least one of the four archaic genomes. For our analysis, we used only the
239 deletions with an allele count greater than 1 in YRI, CEU, and CHB combined. Additionally, we
240 retained only 4,863 human deletion polymorphisms that are in linkage-disequilibrium (LD, $r^2 >$
241 0.9) with at least one SNV (**Supplementary File 2**). We imposed this LD requirement because
242 SNVs in LD with the deletion can enable us to distinguish the shared deletions that are
243 introgressed or recurrent from those that are shared by common descent.

244
245 We found that 575 (11.8%) AMH deletions were shared with archaic hominins, *i.e.*, identified in
246 at least one archaic hominin genome (**Figure 2**). We identified 53 instances of independent
247 emergence (*recurrent deletions*) in archaic hominin and AMH lineages, wherein no SNV that is
248 in LD with the deletion in AMHs accompanied the deletion in archaic hominins. In parallel, we
249 identified 92 deletions that were *introgressed* from archaic hominins into AMHs: the SNVs in LD

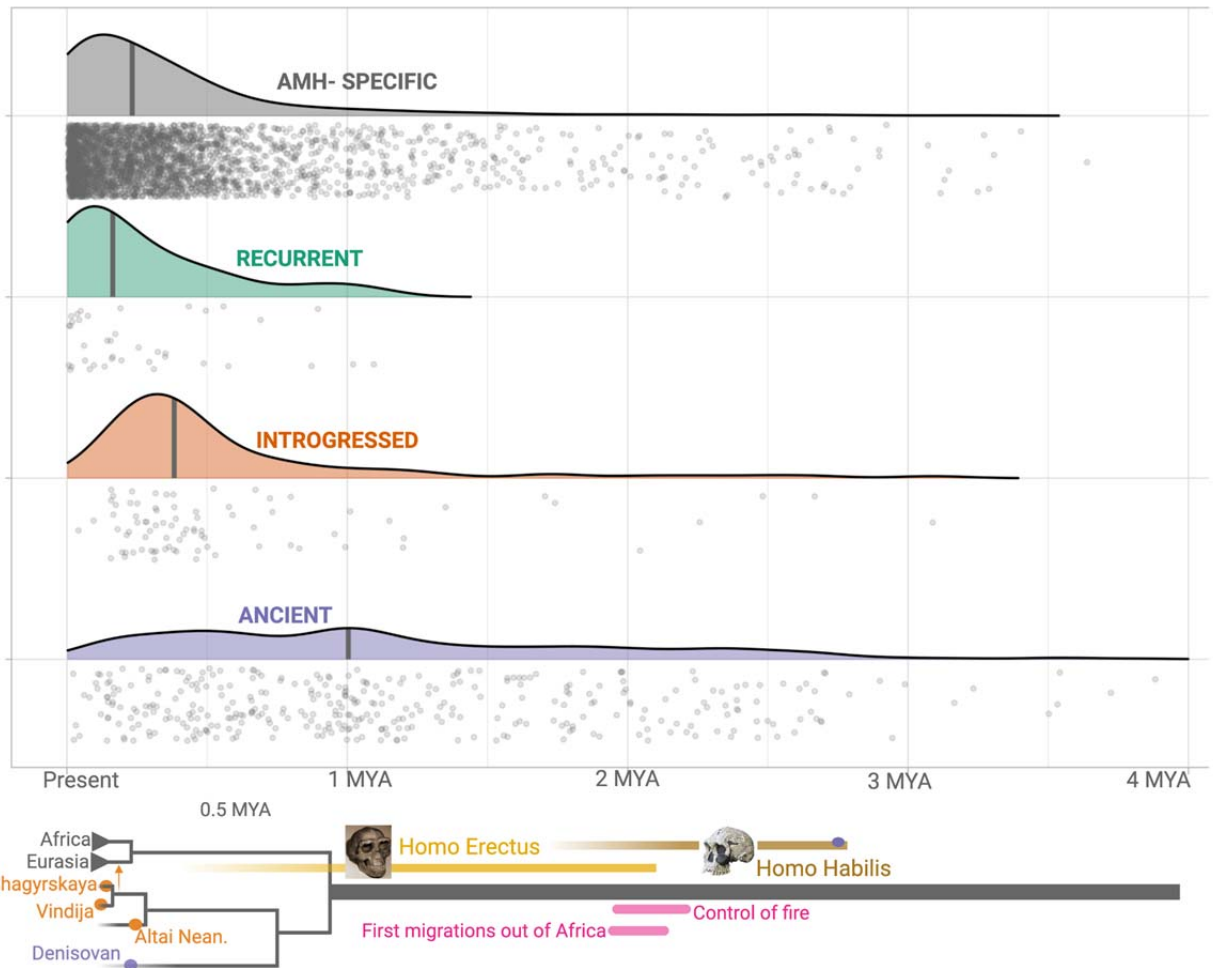
250 with these deletions were present in previously identified introgressed haplotypes (Taskent et al., 2017; Vernot and Akey, 2014). By this process of elimination, we found that 430 (8.8% of the total) shared deletions are ancient polymorphisms, *i.e.*, they are shared with archaic hominins by common descent and thus emerged at least ~700,000 years ago.



254
255 **Figure 2 - Deletions in AMHs that are shared with archaic hominins.** The top panel shows the categorization of
256 deletion polymorphisms as AMH-specific, recurrent (green), introgressed (orange), or ancient (purple). The
257 evolutionary histories of shared deletions are summarized schematically in the bottom panel.

258
259 To confirm that our pipeline for identifying ancient deletions (**Figure 2—figure supplement 1**)
260 has high accuracy, we estimated the ages of deletions, based on the ages of SNVs in LD. We
261 used two methods in parallel: 1) Human Genome Dating (Albers and McVean, 2020); and 2)
262 Relate (Speidel et al., 2019) (**see Methods**). If both our genotyping pipeline and categorization
263 of shared deletions (as recurrent, introgressed, or ancient) are sound, we should expect that
264 $Age(\text{human-specific}) \approx Age(\text{recurrent}) < Age(\text{introgressed}) < Age(\text{Ancient})$. Both methods
265 yielded the expected pattern of ages across the categories of deletions (**Figure 3, Figure 3—**
266 **figure supplement 1**). We found that the median age for ancient deletions, using Relate, is ~1
267 million years. About 15% (63) of these deletions are older than 2 million years. In contrast, the
268 median age of non-ancient deletions is ~239,000 years. As such, we infer that both our
269 genotyping pipeline and deletion categorization approach are sound. Counterintuitively, a small
270 number of “ancient” deletions have very recent dates. This may be due to instances of recent
271 soft sweeps involving some deletions leading to an increased length of the associated haplotype
272 and an artificial decrease in age. Secondly, some ancient deletions may have low frequencies,
273 which too, creates a downward bias in age. Lastly, this may be due to miscategorization of non-
274 ancient deletions as ancient. Next, we test if ancient deletions are enriched for targets of
275 balancing selection relative to non-ancient deletions.

276
277



278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299

Figure 3. Age estimates of the haplotypes harboring polymorphic deletions. The x-axis shows the age estimates, obtained using Relate, for the deletions. For orienting the reader regarding the age of these variants, we provide below a schematic phylogeny representing recent human evolution.

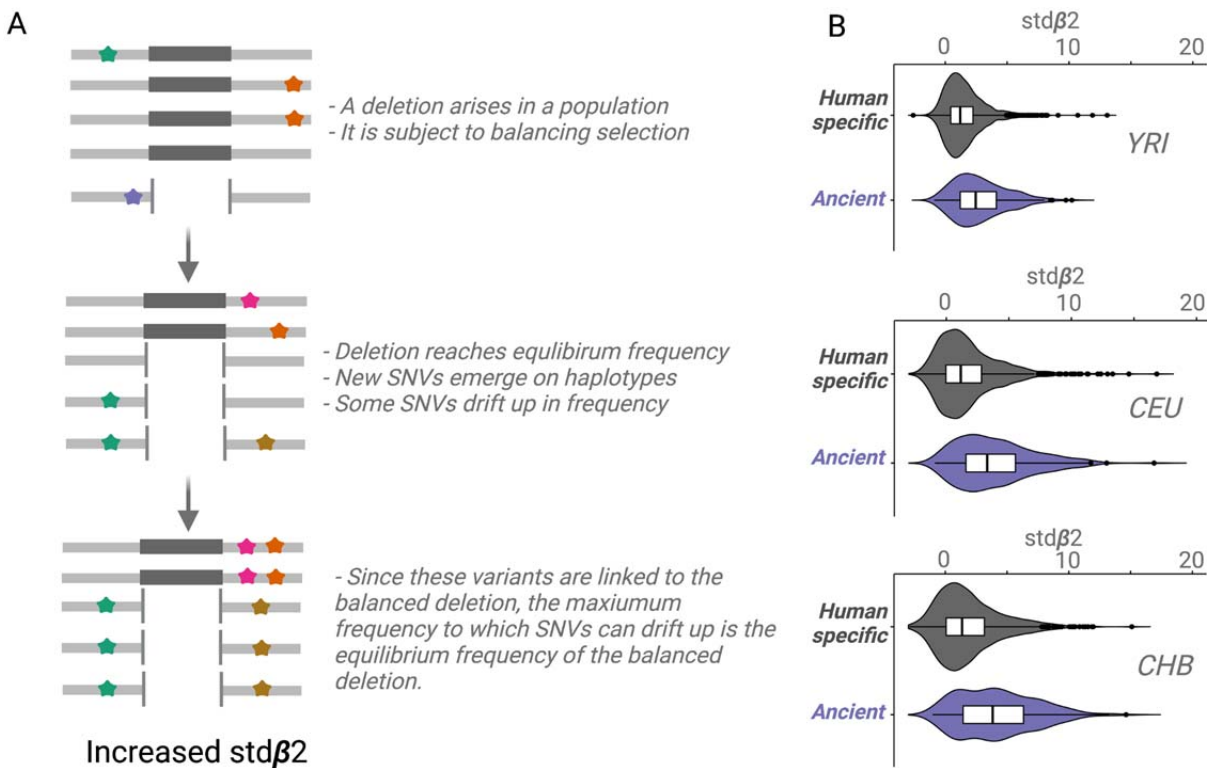
Ancient deletion polymorphisms are more likely to be targets of balancing selection than are non-ancient ones

We used the $\text{std}\sigma^2$ -statistic (Siewert and Voight, 2020) to test the hypothesis that ancient deletion polymorphisms are more likely than are non-ancient deletions to be targets of balancing selection. $\text{std}\sigma^2$ is a measure of balancing selection, that calculates the weighted average of the number of flanking derived variants, where weights are the similarity in frequency between the core allele and the flanking variants (Siewert and Voight, 2017).

For a conceptual understanding of $\text{std}\sigma^2$ (**Figure 4A**), suppose a deletion emerges and the resulting polymorphism is subject to balancing selection; the deletion will rise in frequency until it reaches a certain equilibrium frequency. New SNVs will emerge on the haplotypes carrying the deletion. Some of these SNVs will drift upward in frequency, but since these SNVs will be linked to the deletion, they too can only rise to the equilibrium frequency of the balanced deletion (Siewert and Voight, 2020, 2017). We refer to this type of drift as *Goldilocks drift*, since the linked SNVs drift upward to the “just-right” equilibrium frequency of the balanced deletion. Goldilocks drift thus leads to allelic class build-up (analogous to how hitchhiking leads to sweeps), which refers to a situation involving the fixation of many flanking variants within the set

300 of haplotypes carrying the deletion. $\text{Std}\beta_2$ value for a core variant may be thought of as the
 301 average intensity of *Goldilocks drift* experienced by SNVs around it. Therefore, a high $\text{std}\beta_2$
 302 value for a variant implies that it is either a target of balancing selection or close to a target of
 303 balancing selection.

304
 305 We observed that $\text{std}\beta_2$ estimates for ancient deletions are significantly larger than those for
 306 non-ancient deletions across YRI, CEU, and CHB populations ($p < 10^{-7}$, Wilcoxon) (**Figure 4B**).
 307 These results provide empirical evidence that ancient deletion polymorphisms are enriched for
 308 targets of balancing selection. Our results are consistent with other recent studies (Soni et al.,
 309 2022) that have argued that the role of balancing selection in explaining the maintenance of
 310 common variation in the human lineage is underappreciated.
 311



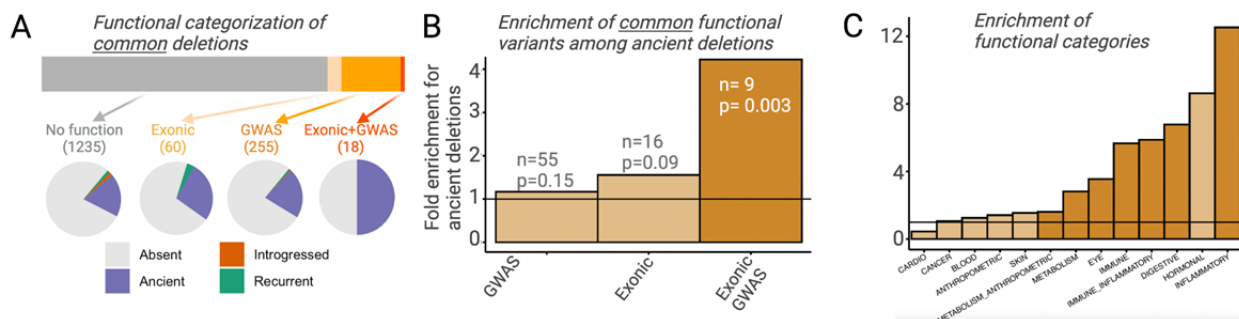
312
 313 **Figure 4. An empirical assessment of putative balancing selection among ancient deletions. A)** The
 314 conceptual framework in which $\text{std}\beta_2$ statistic works. The last step demonstrates “Goldilocks” drift (the process that
 315 results in allelic class build-up). **B)** A box plot for $\text{std}\beta_2$ for AMH-specific, versus ancient deletions (frequency > 5%
 316 in respective populations). Higher $\text{std}\beta_2$ values for older deletions represented in purple empirically show that older
 317 deletions are significantly enriched for targets of balancing selection. All comparisons are significant, $p < 10^{-7}$
 318 (Wilcoxon).

319
 320 Previous genome-wide balancing selection scans focused on either individual genes or single
 321 nucleotide variants (SNVs). Consequently, we do not expect to find many ancient deletions that
 322 have previously been reported as targets of balancing selection. Nevertheless, we investigated
 323 whether the exons in any of the genes that have been reported as targets of balancing selection
 324 in DeGiorgio et al. (2014) or Soni et al. (2022) overlap with ancient deletions. We found no
 325 overlaps. We also found that 77 common (> 5% in YRI, CEU, and CHB combined) ancient
 326 deletions were in LD ($r^2 > 0.9$) with SNVs that had high (in the 95th percentile) $\text{Std}\beta_2$ (Siewert
 327 and Voight, 2020) values in YRI, CEU, and CHB. This is unsurprising since this is the measure
 328 we used to show that ancient deletion polymorphisms are enriched for balancing selection

329 targets. Interestingly, one of the ancient deletions with a high associated $\text{Std}\sigma^2$ value overlaps
 330 a candidate region for balancing selection previously identified using the non-central deviation
 331 (*NCD*) method (Bitarello et al., 2018). This 433 bp deletion (esv3607090), which is 2 million
 332 years old and common across populations, deletes part of an intron of the *STK32A* gene. This
 333 could be an interesting subject for future studies. Regardless, a vast majority of common
 334 ancient deletions (73%) were not reported previously as balancing selection candidates and
 335 thus are novel targets for future studies.

336 337 338 **Phenotypic relevance of ancient deletion polymorphisms**

339 Selection can only act on a region of the genome by means of the phenotypic function it
 340 confers. It follows then that any adaptively maintained ancient polymorphisms must be
 341 functional. If an appreciable proportion of ancient deletion polymorphisms have evolved under
 342 balancing selection and more recent deletion polymorphisms have not, we should expect
 343 ancient deletions to be enriched for functional effects. To avoid biases introduced by different
 344 proportions of rare variants among ancient versus non-ancient deletions, we focus only on
 345 deletion with frequency > 5% in AMHs. For both ancient and non-ancient deletions, we
 346 investigated 1) whether a deletion intersects exons and 2) whether any of the SNVs in LD with it
 347 are associated with UK BioBank GWAS traits with $p < 10^{-8}$ (<http://www.nealelab.is/uk-biobank/>;
 348 **Figure 5A; see Methods**). We did not observe a significant increase in either the proportion of
 349 exonic (ancient=5.6%; non-ancient=3.6%) or GWAS-associated (ancient=19.1%; non-
 350 ancient=16.4%) ancient deletion, relative to non-ancient deletions (**Figure 5B**). However, when
 351 we classify a deletion as functional more conservatively, *i.e.*, it both intersects an exon and has
 352 a GWAS association (*i.e.*, one of the SNVs in LD with it has a GWAS association), we observe
 353 a 4.2-fold enrichment ($p=0.003$) of functional variants among ancient deletions (**Figure 5B**). In
 354 fact, out of the 18 common deletions (frequency > 5%) that both intersect genes and have
 355 GWAS associations, 9 (50%) are ancient. Further, we observed a phenotypic enrichment
 356 among ancient deletions for some GWAS trait categories: a 12.5-fold enrichment ($p=10^{-5}$) of
 357 traits related to inflammatory response and a 2.8-fold enrichment ($p=0.003$) of traits related to
 358 metabolism (**Figure 5C**).
 359



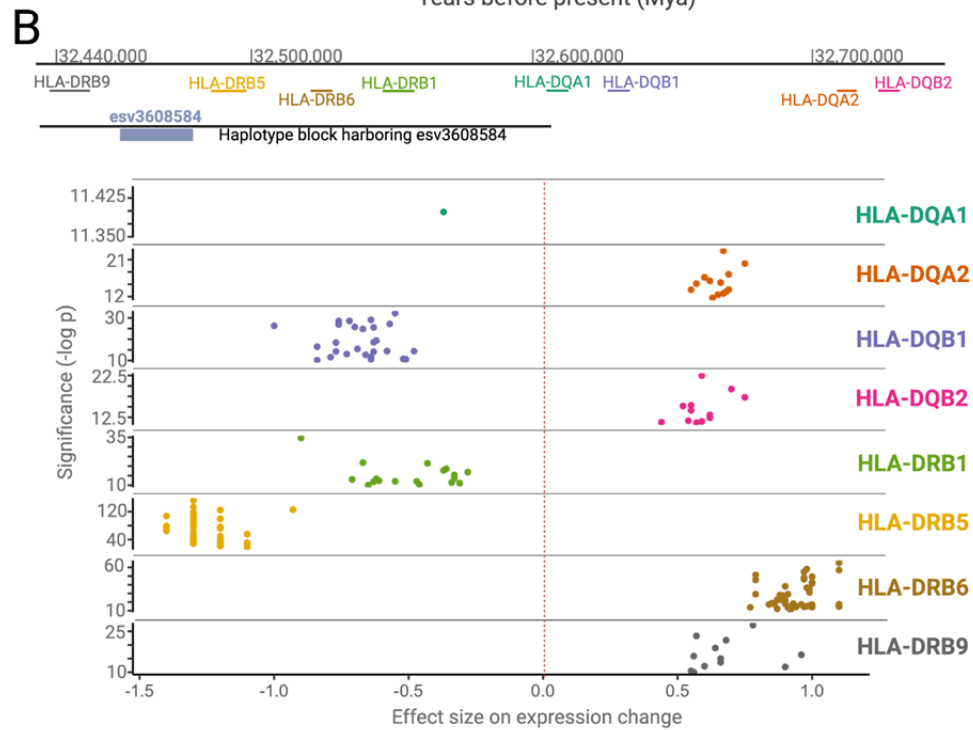
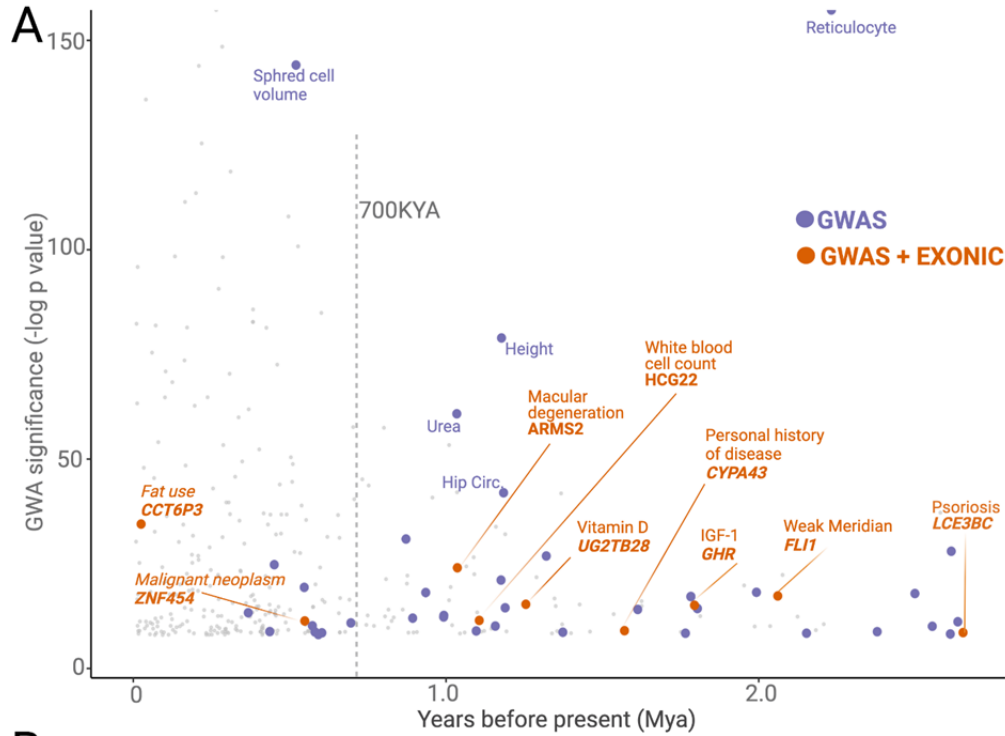
360
 361 **Figure 5. Functional enrichment among ancient deletions.** **A.** Functional categorization of common deletions.
 362 Within each category, the proportions of deletions falling under different evolutionary categories are shown in pie-
 363 charts. **B.** Permutation-based analysis of enrichment of functionality among ancient deletions, relative to non-ancient
 364 deletions. The black horizontal line indicates the expected ratio of 1.0. For each definition of functionality, the number
 365 of functional ancient deletions, and the p-value associated with the enrichment are provided. **C.** Permutation-based
 366 enrichment analysis for different phenotypic categories (based on GWAS) among ancient deletions, relative to non-
 367 ancient deletions. The black horizontal line indicates the expected ratio of 1.0. Dark orange indicates a statistically
 368 significant deviation from the expected ratio of 1.0. Light orange means no significant deviation from the expected
 369 ratio of 1.0.

370
 371 A focused literature review and analysis of functional effects associated with some of the
 372 ancient deletions revealed multiple mechanisms through which they affect function (**Figure 6A**).

373 Firstly, whole gene deletions may affect the function of entire environment-interacting gene
374 families. We found two ancient whole gene deletions: esv3587563 (deleting *LCE3B* and
375 *LCE3C*) and esv3600896 (deleting *UGT2B28*). The members of the *LCE3* and *UGT2B* gene
376 families mediate immune response and steroid metabolism, respectively; genes from both
377 families likely evolved under adaptive forces (de Guzman Strong et al., 2010; Pajic et al., 2016;
378 Starr et al., 2021; Xue et al., 2008). The functional consequence of whole gene deletions is, of
379 course, loss of function of the deleted genes. In addition, esv3587563 is associated with an
380 increase in the expression of *LCE3A* (de Guzman Strong et al., 2010; Pajic et al., 2016), while
381 esv3600896 is associated with an increase in the expression of *UGT2B11*. Thus, we propose
382 that whole-gene deletions of members of environment-interacting gene families may lead to the
383 functional “fine-tuning” of the entire gene family.

384
385 Secondly, we revealed dozens of potentially adaptive ancient deletions that both mediate gene
386 regulation and are associated with human traits. For example, we found multiple ancient
387 deletions that are proximal to the *HLA* locus, associated with immune-related phenotypes, and
388 affect the expression levels of nearby genes. One such example, the deletion esv3608584, is
389 noteworthy within the context of balancing selection because it affects the expression of
390 different *HLA* genes in opposite directions (**Figure 6B**). As such, this deletion may lead to
391 increased susceptibility to some pathogens while increasing the defenses against others.
392 Further, we observed that the ancient deletions in the *HLA* locus also lead to the expression of
393 different isoforms of *HLA* genes. Using the GTeX database, we found at least four other
394 instances where ancient deletions lead to the expression of different isoforms, including
395 deletions affecting the *HLA-DRB1-6*, *HLA-DOB*, *SIRPB1*, *GHR*, and *CYP3A43* genes. We
396 recently showed that the ancient deletion of the third exon of the growth hormone receptor gene
397 leads to the expression of a smaller version of growth hormone, which may be adaptive in times
398 of starvation (Saitou et al., 2021b). The *SIRPB1* gene encodes a glycosylated transmembrane
399 receptor protein (Kharitononkov et al., 1997), and its different isoforms may lead to the
400 recognition of different pathogens. Similarly, *CYP3A43*, a member of the cytochrome p450 gene
401 family, is involved in metabolizing external substances, and genetically determined isoforms
402 contribute to its functional variation in humans (Agarwal et al., 2008). Thus, ancient deletions
403 that lead to specific isoform expression may have been adaptively evolving to adjust the
404 function of environment-interacting genes across both geography and time. It is important to
405 acknowledge that these non-exonic deletions may not be the causal variant in the associated
406 haplotypes. Nevertheless, the full extent of deletion polymorphisms shaping the expression
407 levels and sculpting the isoform diversity at the genetic level remains a fascinating area of future
408 research.

409
410



411
 412 **Figure 6. A.** The significance levels ($-\log(p\text{-value})$) of phenotypic associations of deletions with GWAS traits as a
 413 function of their emergence time. Gray points indicate non-ancient deletions. Purple and orange points indicate non-
 414 exonic ancient deletions with GWAS hits and exonic ancient deletions with GWAS hits, respectively. The genes
 415 whose exons are covered by ancient deletions, and the traits associated with ancient deletions are mentioned in the
 416 plot. **B.** The significance levels ($-\log(p\text{-value})$) and sizes of expression level changes of nearby *HLA* genes associated
 417 with the presence of the deletion esv3608584. Each color refers to a different *HLA* gene. Each point in a given color
 418 represents a different tissue. Only those tissues whose expression level changes are statistically significant are
 419 shown here.

420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465

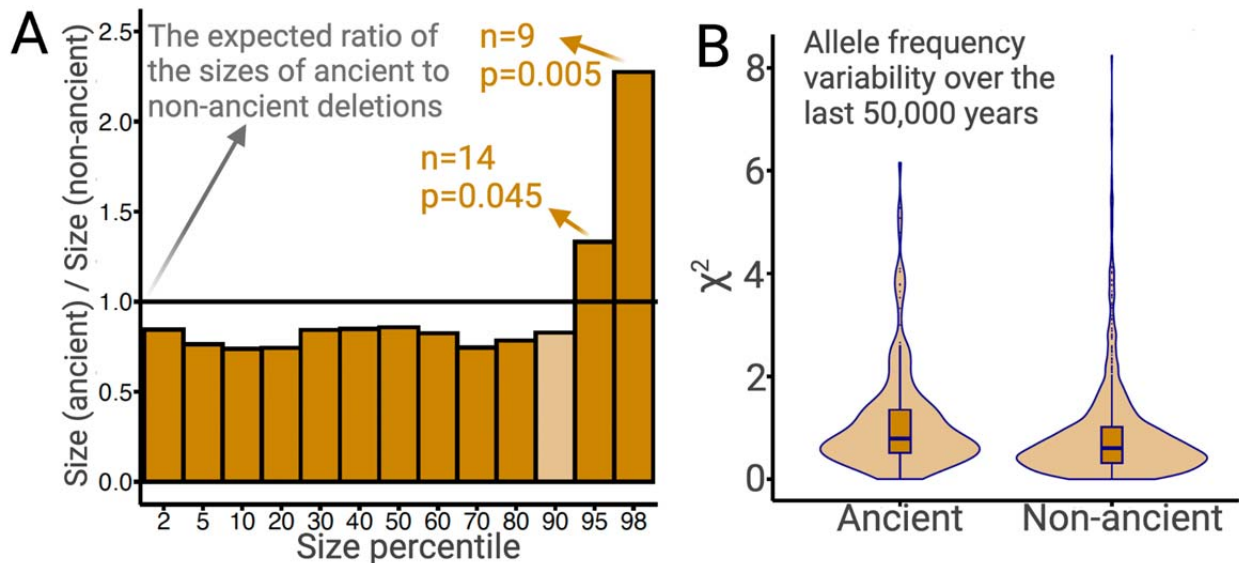
The effect of negative selection is stronger on deletions than SNVs

Based on previous work, we expect that deletions are more likely than SNVs to be under negative selection (Conrad et al., 2006; Kondrashov, 2017; Lin et al., 2015; Lin and Gokcumen, 2019). To investigate the magnitude of this effect, we compared the proportion of SNVs and deletion polymorphisms that are ancient. Applying the same bioinformatic pipeline to identify ancient polymorphisms in both cases, we found that 13.7% of SNVs and 9.6% of deletion polymorphisms are ancient in YRI. This result alone suggests that deletion polymorphisms are more likely than SNVs to be eliminated by negative selection, a trend that we expect to be more pronounced with increasing ages of polymorphisms. The greater intensity of negative selection acting on deletions implies that deletions are, in general, more deleterious than SNVs. It follows that younger (non-ancient) deletions currently segregating in human populations, which negative selection has not yet purged, are more likely to be deleterious (and perhaps disease-causing) than SNVs (Kondrashov, 2017).

The preceding argument makes intuitive sense since a given deletion spans more bases than does a SNV. If this intuition is correct, we expect that larger deletions should, on average, experience more intense negative selection. Since most large ancient deletions would have been purged by negative selection, we expect surviving ancient deletions to be, on average, smaller than non-ancient deletions. We test this using common deletions (frequency > 5% in YRI, CEU, and CHB combined). Ancient deletions are indeed 14% shorter than non-ancient deletions ($p=0.02$; permutation test). Nevertheless, there is an excess of long deletions among ancient deletions relative to non-ancient deletions (**Figure 7A**). In particular, the 95th and 98th size percentiles of ancient deletions are 33% ($p=0.04$; permutation test) and 128% ($p=0.005$; permutation test) larger than non-ancient deletions, respectively (see **Methods**). This excess of longer deletions is inconsistent with evolution under neutrality or negative selection. Therefore, the longest 5% of ancient deletions are excellent targets for future studies of balancing selection. In fact, 3 out of the 9 GWAS-associated common exonic deletions intersecting *SIRPB1*, *LCE3A*, *LCE3B*, and *UGT2B28* are in the 95th percentile of the size distribution of ancient deletions.

Strong overdominance is rare among ancient deletion polymorphisms

Having established that ancient deletion polymorphisms appear enriched for targets of balancing selection, we wanted to investigate whether classical overdominance is a common mechanism underlying this observation. To accomplish this, we first identified the genomic signatures that we expect to see in a region where a polymorphism has evolved under overdominance, and then we look for these signatures among ancient deletions. To identify the signatures of overdominance, we simulated sequence evolution under neutrality and overdominance (using a variety of selection coefficients), in turn, for variants that emerged one million years ago. We asked whether we can distinguish between neutrality and overdominance by calculating several population genetic statistics, including Tajima's D , π , θ , etc., on sequences generated from the neutral and overdominance simulations (see **Methods** for full list). We found that none of these statistics alone can distinguish overdominance from neutrality, even for strong selection coefficients (**Figure 7—figure supplement 1**)



466
467
468 **Figure 7. A.** The ratios of sizes of ancient deletions to those of non-ancient deletions at different size percentiles. The
469 black horizontal line refers to the expected ratio of 1.0. Dark orange bars refer to a statistically significant
470 (permutation test) deviation from the expected ratio. Light orange bars mean that the deviation from the extend ratio
471 of 1.0 is not statistically significant. **B.** The estimated measure of allele frequency change (χ^2) between 50,000 and
472 5,000 years before present in common ancient versus common non-ancient deletions. Ancient deletions have
473 significantly ($p=2 \times 10^{-7}$, Wilcoxon) higher frequency variability over the last 50,000 years.

474 Instead, we found that a distinct feature of overdominance is that the allele frequency rapidly
475 increases (similar to a selective sweep) until it reaches an equilibrium frequency, whereafter it
476 remains remarkably stable across time (**Figure 7—figure supplement 1**). In contrast, under
477 neutrality, a random change in allele frequency in every generation produced elevated noise
478 across time in allele frequency trajectories. To ascertain if overdominance is a common
479 mechanism of evolution for ancient deletions, we inferred the allele frequency trajectories of
480 ancient deletions using *Relate* (Speidel et al., 2019) and quantified the variation in allele
481 frequency between 5,000 and 50,000 years ago by squaring the standardized allele frequency
482 difference (χ^2) (**Methods**). We already know that ancient deletions are enriched for targets of
483 balancing selection. If large proportion of these balancing selection targets have evolved under
484 overdominance, we expect ancient deletions to have more stable allele frequencies across time,
485 relative to non-ancient deletions, leading to smaller χ^2 values on average. However, we do not
486 observe this trend among common (frequency > 5% in YRI, CEU, and CHB combined) deletions
487 (**Figure 7B**). Consequently, at least with our current resolution of allele frequency trajectory
488 estimation, we found no evidence for overdominance being the prime mode of balancing
489 selection operating on ancient polymorphisms in AMHs.

490
491 In fact, we observe that ancient deletion polymorphisms exhibit greater allele frequency
492 variation than do non-ancient deletions (**Figure 7B**; $p=2 \times 10^{-7}$, Wilcoxon). This suggests that a
493 large proportion of the instances of medium-term balancing selection likely involve temporally
494 and spatially variable selection, which lead to elevated levels of allele-frequency variation over
495 time. This is consistent with our locus-specific analyses of ancient deletion polymorphisms. For
496 example, we recently reported that the deletion (esv3604875) of the third exon of the human
497 growth hormone receptor gene (*GHR*) has evolved under temporally and geographically
498 variable adaptive constraints (Saitou et al., 2021b). In fact, this deletion is in the 93rd percentile
499 of χ^2 values. Ancient deletions like these are common and old, and also exhibit high population
500 differentiation. Collectively, we argue that such adaptive maintenance of ancient, functional
501 alleles may be due to varying selection trends across geographies and time.

502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527

CONCLUSION

This manuscript asked whether adaptive forces have maintained ancient deletion polymorphisms in humans. We provide evidence supporting the idea that AMHs exhibit an excess of ancient polymorphisms, relative to the neutral expectation. Using simulations and empirical data, we provide evidence for the notion that balancing selection is likely a considerable force shaping the extant functional deletion polymorphisms. We show that when functionality is defined conservatively, ancient deletions are disproportionately functional, compared to non-ancient deletions. In fact, 50% of such functionally relevant deletions are ancient. Additionally, ancient deletions are enriched for association related to metabolism and inflammatory response. Our results suggest that classical overdominance may not be the prime mode of balancing selection affecting the evolution of ancient deletions. Instead, geographically and temporally variable, as well as frequency-dependent selection may underlie the maintenance of ancient functional deletions. We also provide insights about the mechanisms by which a deletion could confer function: in addition to previously defined functional effects of deletions such as the loss of a gene's function and regulation of expression levels, we highlight multiple instances where the presence of ancient deletions lead to the expression of different isoforms. Overall, our study contributes to the growing body of evidence supporting the notion that balancing selection may be an important force in the evolution of genomic variation shared among human populations. These ancient variants are an important part of our legacy as a species: something we all share.

528 **DATA AVAILABILITY**

529 All data generated can be found in the supplementary files. The codes that we used to generate
530 our datasets and simulations can be found either in Methods or on our GitHub page
531 (<https://github.com/GokcumenLab>). The consolidated file including all significant UKBioBank
532 GWAS associations for SNVs is available at our FigShare
533 (https://figshare.com/articles/dataset/Table_S3_for_Aqil_et_al_2022/19606192).
534

535 **SUPPLEMENTARY FILES**

536 **Supplementary File 1:** Deletion genotypes in archaic hominin genomes
537 **Supplementary File 2:** Haplotype-level analysis of deletion polymorphisms
538 **Supplementary File 3:** Curated GWAS meta-categories
539

540
541 **FUNDING**

542 O.G. acknowledges support from the National Science Foundation (Grant No. 2123284). L.S. is
543 funded by a Sir Henry Wellcome fellowship (220457/Z/20/Z). This research was funded in part
544 by the Wellcome Trust. For the purpose of Open Access, the authors have applied a CC BY
545 public copyright license to any Author Accepted Manuscript version arising from this
546 submission.
547

548 **ACKNOWLEDGEMENTS**

549 We thank Dr. Victor Albert and Dr. Vincent Lynch for their careful reading of this manuscript. We
550 acknowledge Petar Pajic and Charikleia Karageorgiou for their insightful discussions throughout
551 the development of this project.

552 **ETHICS**

553 All data used in this study have been previously published in peer-reviewed journals.
554

555 **COMPETING INTERESTS**

556 We have no competing interests
557
558
559
560

561 **METHODS**

562

563 **Proportion of ancient polymorphisms in neutral simulations versus observation:**

564 We compared the proportion of ancient polymorphisms in randomly chosen YRI SNVs against
565 the expected distribution of this proportion under neutrality. We obtained SNVs from the 1000G
566 phase-3 vcf files (1000 Genomes Project Consortium et al., 2015) for the analysis described
567 above. Using a script written in AWK (Aho et al., 1978), we subsetted the vcf files to retain only
568 those biallelic SNVs that contained information about the ancestral/derived status of the two
569 alleles. We used the --keep option in vcftools (Danecek et al., 2011) to retain individuals only
570 from the Yoruba population. We then used the --mac option to retain only those SNVs for which
571 the minor allele count was greater than 1. The allele-count filter was used to exclude singletons
572 which could create spurious results for the linkage disequilibrium analysis described below. On
573 the resulting vcf files, we used the SelectVariants tool with the --select-random-fraction option in
574 GATK (Van der Auwera and O'Connor, 2020) to randomly retain 0.25% of the variants. This
575 resulted in a set of 38,231 SNVs. Next, we investigated whether the random SNVs are in
576 linkage disequilibrium with any other SNVs in their vicinity. In particular, we used the --hap-r2-
577 positions and --min-r2 option in vcftools, along with the ancestral/derived status of alleles, to
578 retain only those random SNVs wherein the derived allele was in linkage disequilibrium ($r^2 >$
579 0.9) with the derived allele of another SNV within a 50 kb radius of the random SNV. This
580 resulted in 28,491 SNVs. We only retained the random SNVs with variants in LD in their vicinity
581 to rule out cases of recurrence of the SNV between AMHs and archaic hominins, as described
582 below. There seems to be one bias that may be introduced by eliminating polymorphisms with
583 no SNVs in LD in AMHs: this would bias our analysis to regions with low recombination rates. In
584 these regions, we would expect higher background selection due to the Hill-Robertson effect,
585 leading to a deflation in the proportion of ancient polymorphisms. Since we observe a larger
586 than expected proportion of ancient polymorphisms despite this bias, we can conclude that this
587 bias only makes our analysis more conservative.

588

589 We then inspected the 28,491 random SNVs to see whether the derived alleles are shared with
590 any of the four high coverage archaic genomes (Altai Neanderthal, Vindija Neanderthal,
591 Chagyrskaya Neanderthal, and the Denisovan). We found that 4,616 SNVs (16.2%) had their
592 derived allele shared (either homo- or heterozygously) with at least one of the four archaic
593 hominins. Since we are only interested in polymorphisms older than 700,000 years, we want to
594 focus only on SNVs where the derived allele is shared with archaic hominins by *common*
595 *descent*. We thus excluded the SNVs where the derived allele emerged independently
596 (recurrence) in AMHs and archaic hominins in the following way. For each of the 4,616 core
597 SNVs with shared derived alleles in archaic hominins, we tested whether any of the derived
598 alleles in LD with the core derived allele is also present in any of the same archaic hominins
599 which carry the derived allele for the core SNV itself. If any of the archaic hominins contain both
600 the derived allele of the core SNV along with at least one derived allele in LD with the derived
601 allele of the core SNV, we classify that core SNV as “shared by common descent.” If any core
602 SNV has a derived allele shared with archaic hominins but not a derived allele in LD with the
603 core SNV, we classify the core SNV as “recurrent.” This approach yielded 3,894 SNVs (13.7%)
604 wherein the derived allele is shared with at least one of the archaic hominins by common
605 descent.

606

607 In order to investigate whether this percentage (13.7%) of ancient polymorphisms (700,000
608 years old polymorphisms) is significantly higher than the neutral expectation, we wanted to
609 calculate the same percentage for a set of neutrally simulated SNVs. We used the program *ms*
610 (Hudson, 2002) to neutrally simulate a set of 20,000 (independent, and therefore freely
611 recombining) variants 2,000 times using various models. All of these models included 216

612 haploid genomes representing YRI (matching the YRI sample size in 1000G dataset) and 2
613 haploid genomes representing each of the four archaic hominins. For every model, and for each
614 of the 2,000 runs we calculated the proportion of SNVs present in the YRI with minor allele-
615 count > 1 that were shared with at least one of the archaic hominins. Thereby, we obtained a
616 distribution of the proportion of ancient polymorphisms under each of the models.
617

618 We used a total of 36 models. The models varied by three parameters: N_e of Yoruba/AMH, N_e of
619 archaic hominins, and the time of divergence between the AMH and archaic hominin lineages.
620 For 27 of these models, N_e for humans was constant across time (ranging from 10,000 to
621 30,000 across models); for 9, it varied over time (Speidel et al., 2019) (**Figure 1B**). The N_e of
622 the archaic hominin lineage was constant over time in each model ranged from 1,000 to 3,000
623 across models; and the time of divergence between AMH and archaic hominins ranges from
624 500 to 700 kya across models (Bergström et al., 2021; Mafessoni et al., 2020; Meyer et al.,
625 2012; Prüfer et al., 2014). In the main text, we focused on the model with variable AMH- N_e
626 (Speidel et al., 2019), using the well-accepted archaic hominin N_e of 1,000 and a divergence
627 time of 700 kya; we referred to this as the “base model”. For each model, we assumed that
628 Denisovans diverged from Neanderthals ~400,000 years ago, the Altai Neanderthal lineage
629 separated from the Vindija-Chagyrskaya lineage ~130,000 years ago, and the Vindija lineage
630 separated from the Chagyrskaya lineage ~90,000 years ago (Bergström et al., 2021; Gravel et
631 al., 2011; Mafessoni et al., 2020; Meyer et al., 2012; Prüfer et al., 2014). Moreover, the
632 generation time was assumed to be 29 years (Fenner, 2005; Langergraber et al., 2012; Li and
633 Durbin, 2011).
634

635 To test if the empirical excess of ancient polymorphisms is more pronounced among low
636 derived allele frequency variants (due to recurrence), we repeated our analysis using the base
637 model (variable AMH N_e , Archaic $N_e = 1,000$, divergence time = 700 kya) for simulations,
638 dividing the empirical and simulated SNVs into derived allele frequency bins. To ensure there
639 are enough variants in each frequency bin, we used a larger set of randomly chosen YRI SNVs.
640 In particular, we used 300,000 SNVs that have a minor allele count > 1 and at least one variant
641 in LD ($r^2 > 0.9$) with them. Using the same method as described above, we identified the SNVs
642 that are shared with archaic hominins by common descent. We used the base model to simulate
643 1500 runs of 400,000 SNVs. Thereupon, we divided the empirical and simulated SNVs into 10
644 derived allele frequency bins of uniform length, and compared the simulated distribution of the
645 proportion of ancient polymorphisms with the observed proportion in each bin (**Figure 1—figure
646 supplement 1**). Moreover, to gauge whether recurrence at CpG sites leads to the excess of
647 ancient polymorphisms, we subsetted the 300,000 SNVs to retain only the 21,402 A↔T SNVs,
648 and calculated the proportion of ancient polymorphisms therein.
649

650 Additionally, we also performed another set of neutral simulations, this time with structure
651 introduced in the population ancestral to the archaic hominins and AMHs. This too was done
652 using Hudson’s *ms*. This was done using a constant size model with YRI/AMH $N_e = 14,474$
653 (Gravel et al., 2011), Archaic $N_e = 1,000$, and divergence time = 700 kya. The effective
654 population size for each of the subgroups in the population ancestral to both AMHs and archaic
655 hominins was set to 10,000. We define $m_{i,j}$ as the fraction of subgroup i that is formed by the
656 migrants of subgroup j in each generation, where $i \neq j$ and $i, j \in \{1, 2, 3\}$. For all i and j , where $i \neq j$,
657 we set $m_{i,j} = m$. The program *ms* takes this parameter in the form of $M = 4Nm$ (where $N =$
658 10,000 is the effective population size of each subgroup). We performed simulations for 10,000
659 different values of M chosen uniformly on the log scale from the range $(0.01, 100)$. This is akin
660 to running simulations using 10,000 different values of m in the range $(0.25 \times 10^{-7}, 0.25 \times 10^{-2})$.

661 For each m , 1000 variants were simulated. Thus, for each m , we calculated the percentage of
662 variants in Yoruba (with allele-count >1) wherein the derived allele was shared with the archaic
663 hominins. The proportion of allele-sharing in simulations equaled or exceeded the proportion
664 (13.7%) observed in real-life at approximately $m \leq 0.0075\%$.

665 **Identifying deletions in archaic genomes**

666 The identification of deletions in archaic hominins was predicated on the concept that a deletion
667 in an archaic hominin would correspond to a low read depth in the window of deletion in the
668 hominin's genome. We started with two main types of input files: 1) The VCF file for the 1000
669 Genomes Phase 3 dataset; and 2) A BAM file for each high-coverage archaic hominin.

670
671
672 The 1000G phase-3 VCF file was obtained from <https://www.internationalgenome.org/data>. This
673 includes 84.4 million variants from 2504 individuals across 26 populations.. The file was then
674 converted to a BED file (with three tab-separated columns representing chromosome numbers,
675 start positions, and end positions of deletions) using a script written in AWK. This VCF file was
676 filtered to retain only biallelic autosomal deletions. This amounted to 32,154 deletions. We
677 genotyped all these deletions in the four high-coverage archaic hominin genomes. (Note that in
678 the main text, we focused only on 4,863 deletions with both allele-count > 1 in YRI, CEU, and
679 CHB combined, and at least one SNV in LD.)

680
681 The sequence files for archaic genomes mapped to hg19 were obtained from
682 <https://www.eva.mpg.de/genetics/genome-projects.html?Fsize=0%2C%252%27A%3D0>. These
683 BAM files containing mapping information (such as the start and end coordinates of the part of
684 the genome to which a read maps) were converted to BED files. This was done using the
685 `bamToBed` command in the `bedtools` module (Quinlan and Hall, 2010). We then used the two
686 types of BED files to count the number of reads for each archaic genome that mapped to a
687 region of the genome that is polymorphically deleted in AMHs. In order to achieve this, we used
688 the `intersectBed` command with the `-c` option within the `bedtools` module. This command
689 counts the number of reads in an archaic genome that intersects with the region of the genome
690 harboring a deletion polymorphism in AMHs.

691
692 Next, for every archaic genome, we normalized the number of reads at each window of deletion
693 by the size of the window.

694

$$695 \quad \text{Normalized Read Depth} = \frac{\# \text{ of reads intersecting the window of deletion}}{\text{Size of the window of deletion}}$$

696

697 For each archaic genome, we wanted to calculate the Z-scores of the normalized read depths
698 across all windows of deletion, and classify a window as a deletion if the normalized read depth
699 was below a certain threshold. To prevent outliers from affecting measures of central tendency
700 and spread, and therefore the Z-score threshold, we use the modified Z-score to classify a
701 region as deleted or non-deleted in an archaic genome. The modified Z-score uses median (as
702 opposed to mean) and median absolute deviation (as opposed to standard deviation) to
703 calculate the Z-score. For a given archaic genome, the modified Z-score of the normalized read
704 depth at the i^{th} window of deletion is given by:

705

$$706 \quad \text{Mod}Z_i = \frac{r_i - \text{Median}(R)}{\text{MedianAbsoluteDeviation}(R)}$$

707

708 where r_i denotes the normalized read depth at a given window of deletion, and R denotes the
709 random variable representing normalized read depth.

710
711 (Iglewicz and Hoaglin, 1993) have suggested that a threshold ± 3.5 is reasonable for outlier
712 detection using modified z-scores. Nevertheless, for our purposes, we used a more
713 conservative threshold of -5, which we deemed more-appropriate based on spot-checking. For
714 example, if the modified Z-score (of the normalized read-depth) at a window of deletion was less
715 than -5 in the Vindija Neanderthal, that window was classified as deleted in the Vindija
716 Neanderthal. The distributions of these modified Z-scores across windows of deletions for the
717 four high-coverage archaic genomes are illustrated in **figure 2—figure supplement 1**. All
718 calculations downstream of obtaining the raw numbers of reads from archaic genomes
719 intersecting with the windows of deletion, were performed using a script that we wrote in R. The
720 read-depth analysis was done using all 32,154 AHM deletions (results for the status of these
721 deletions in the four high-coverage archaic genomes are available in **Supplementary File 1**).

722
723

724 **Identifying SNVs that are in linkage disequilibrium with deletion polymorphisms in** 725 **AMHs:**

726 We subsetted the 1000G phase-3 VCF files (there is a separate file for each chromosome)
727 obtained from <https://www.internationalgenome.org/data> to retain individuals only from CEU
728 ($n=103$), CHB ($n=99$), and YRI ($n=108$) populations. This filter was applied using the --keep
729 option in the module VCFtools (Danecek et al., 2011). All variants that had a minor allele count
730 of less than 2 were eliminated using the --mac filter in VCFtools. We used the resulting VCF
731 files to identify SNVs in LD with each of the autosomal biallelic deletions (with minor allele-count
732 > 1 in YRI, CEU, and CHB combined) within a 50 kb radius of the deletion. We did this using the
733 --hap-r2-positions and --min-r2 0.9 flags in VCFtools. For each autosomal biallelic deletion
734 with allele-count > 1 , this gave us a list of SNVs in LD with the deletion with $r^2 > 0.9$, if such
735 SNVs existed. At least one such SNV in LD existed for 4,863 deletions. We called this set of
736 deletion the “deletion dataset” and based all our downstream analysis on it.

737

738 It is important to describe why we only focused on deletions with identifiable variants (most of
739 them SNVs) in LD with them. We can only eliminate potentially introgressed deletions by
740 checking whether at least one of the SNVs in LD with a deletion is already known to be
741 introgressed. Moreover, we can confirm whether a deletion shared between archaic hominins
742 and AMHs is identical by descent (thereby eliminating recurrence) by checking whether the
743 same SNVs accompany the deletion in AMHs and archaic hominins. This filtering would not be
744 possible if our deletions were not flanked by variants in LD with them. A shortcoming of this
745 approach is that it fails to capture balanced deletions that are not in LD with at least one SNV.

746

747 **Eliminating instances of recurrence and introgression**

748 We found that 575 human polymorphic deletions are also present in at least one archaic
749 hominin genome. In order to ensure that we do further analysis only on deletions that are
750 shared with archaic hominins by common descent, we wanted to eliminate shared deletions that
751 were recurrent or introgressed.

752

753 We removed recurrent deletions (those emerging in AMHs and archaic hominins independently)
754 by retaining only those shared deletions for which at least one allele in LD with the deletion in
755 AMHs was also present in at least one archaic genome that harbored the deletion. To do this,
756 we needed to know whether variants in LD with deletions are present or absent in the archaic
757 genomes. We started with two types of inputs: 1) VCF files for each of the archaic genomes

758 (mapped to hg19) and 2) a file containing all the variants (SNVs) in LD with polymorphic
759 deletions in AMHs. We filtered the VCF files to include only the SNVs in LD with shared
760 deletions. This was done using the --positions flag in VCFtools. The presence or absence of
761 every variant in LD was then determined using the vcf files for the archaic hominins. The
762 procedure was implemented using an AWK script. 53 shared deletions were classified as
763 “recurrent” using this approach.

764
765 In order to eliminate introgressed shared deletions, we used the results published by (Taskent
766 et al., 2020). In their study, the authors had identified introgressed haplotypes in Eurasians
767 using the S^* statistic. They had also published a list of S^* -significant SNVs that characterize
768 introgressed haplotypes. We stamped out the shared deletions that were both absent in Yoruba
769 and for which at least one allele in LD was among the S^* -significant variants listed in the study
770 mentioned above. We thus eliminated 92 deletions that were likely introgressed from archaic
771 hominins into AMHs.

772
773
774

775 **Age of deletions and allele frequency trajectories**

776 We estimate the ages of the deletions in the *deletion dataset* using two methods: 1) Human
777 Genome Dating database (<https://human.genome.dating/download/index>); and 2) RELATE
778 (Speidel et al., 2019).

779
780 The Human Genome Dating database (<https://human.genome.dating/download/index>) hosts
781 age estimates for over 45 million single nucleotide variants (SNVs) (Albers and McVean, 2020).
782 This database reports multiple age estimates for each SNV. We used the median age estimate
783 calculated using the joint clock. Since this database only includes age estimates for SNVs (and
784 not for deletions), we could only date a deletion if the dating database contained the age
785 estimate for at least one of the variants in LD ($r^2 > 0.9$) with the deletion. If age estimates were
786 available for only one variant in LD, the same age estimate was assigned to the deletion. If age
787 estimates were available for more than one variant in LD with the deletion, we used the highest
788 age estimate, which may be inaccurate in certain cases.

789
790 Relate is a method that estimates genome-wide genealogies and can be used to infer the age of
791 a variant (Speidel et al., 2019). We used Relate to infer the ages of the deletions in the *deletion*
792 *dataset*. To this end, we used previously inferred genome-wide genealogies for samples of the
793 SGDP dataset (Mallick et al., 2016; Speidel et al., 2021), available from
794 <https://www.dropbox.com/sh/2gjyxe3kqzh932o/AAQcipCHnySgEB873t9EQjNa?dl=0>. For each
795 deletion, we used SNVs in LD where the derived allele was tagging the deletion at an r^2
796 exceeding 0.9 and calculated the mean age of such SNVs to date each deletion.

797
798 To quantify allele frequency variation, we computed the ratio of lineages carrying the derived
799 allele by the total number of lineages remaining at 5,000 years and 50,000 years before
800 present, but only if the number of lineages remaining at 50,000 years exceeded 10% of the
801 present-day sample size. We then standardized the allele frequency change stratified by
802 present-day allele frequency, by calculating the mean and standard deviation given present-day
803 frequency. Finally, we squared this standardized allele frequency change to obtain our statistic
804 χ^2 , which is expected to have a Chi-squared distribution with one degree of freedom under
805 neutrality, and smaller values for more stable trajectories. This approach was inspired by (Edge
806 and Coop, 2019), who used a similar approach to quantify polygenic positive selection using
807 genealogies.

808

809 **Beta measure**

810 We used a recent and robust measure of balancing selection (Siewert and Voight, 2017), std^2 ,
811 to investigate whether ancient deletions are enriched for targets of balancing selection. A high
812 std^2 for a variant is indicative of balancing selection.

813
814 In our study, we estimated std^2 for the deletions in the deletion dataset using SNVs in LD with
815 them. We did this for the CEU, CHB, and YRI population separately. The std^2 scores for
816 SNVs are publicly available (<https://github.com/ksiewert/BetaScan>) for the CEU, CHB, and YRI
817 populations. For each deletion in each of these populations, we obtained the std^2 values for
818 variants in LD with deletions whenever they were available. For a given deletion, when the
819 std^2 values were available for more than one variant in LD with the deletion, we used two
820 approaches to estimate the std^2 for the deletions. In the first approach, we used the highest
821 std^2 among the LD variants as the estimate for the std^2 for the deletion. We call this
822 BETAMAX. In the second approach, we focused on the std^2 values for the SNVs that were in
823 LD with the deletion with the highest r^2 value. If multiple SNVs were in LD with the deletion with
824 the highest r^2 value, we used the std^2 value of the SNV that was closest to the deletion among
825 these SNVs as the estimate for std^2 for the deletion. We call this BETAPRIME. We performed
826 this process for YRI, CEU, and CHB populations separately to arrive at std^2 estimates for
827 deletions in our deletion dataset in each of these three populations. Using both BETAMAX and
828 BETAPRIME gave us similar trends across populations. These values are available in
829 **Supplementary File 2.**

830
831

832 **Ascribing phenotypic relevance to deletion**

833 We used two criteria to ascribe phenotypic relevance to deletions: 1) intersection of the deletion
834 with at least one exon; and 2) association of a SNV in LD with the deletion with a GWAS trait.

835
836 In order to identify deletions that intersect with exons, we started with the genome annotation
837 file download from [https://hgdownload-](https://hgdownload-test.gi.ucsc.edu/goldenPath/hg19/bigZips/genes/hg19.refGene.gtf.gz)
838 [test.gi.ucsc.edu/goldenPath/hg19/bigZips/genes/hg19.refGene.gtf.gz](https://hgdownload-test.gi.ucsc.edu/goldenPath/hg19/bigZips/genes/hg19.refGene.gtf.gz). Using an AWK script, this
839 GTF file was then converted to a BED file containing five columns: 1) annotation's chromosome
840 number; 2) annotation's start position; 3) annotation's end position; 4) gene name; and 5) type
841 of feature. Only the rows wherein the type of feature was "exon" and the chromosome number
842 was between 1 and 22 were retained. All repeated entries (rows) were eliminated. The resulting
843 file contained only columns 1) to column 4). We then used a BED file containing information
844 about the AMH deletions in our deletion dataset and the BED file mentioned above to identify
845 deletions spanning exons. This was done using the intersectBed option with -wa and -wb flags
846 in the BEDtools module. On the resulting file, we used the groupby tool with the "-o freqdesc"
847 flag in the BEDtools module in order to obtain a file containing the names of the genes (and the
848 number of exons within each intersecting gene) that overlap the deletions. 243 (5%) of the
849 4,863 deletions in the deletion dataset were exonic.

850
851 The second method to ascribe phenotypic relevance to deletions was to use results from
852 previously published Genome Wide Association Studies (GWAS). We used a publicly available
853 catalog of GWAS results based on the UK BioBank data (<http://www.nealelab.is/uk-biobank/>). In
854 particular, we used data for 4,113 traits. For each trait, we used data that produced results using
855 both sexes. For continuous traits, we used the raw version of the data, as opposed to the
856 inverse rank normalized version. For each trait, only those SNVs were retained that were
857 associated with the phenotype with a p-value less than 10^{-8} . We are making available a
858 consolidated table with all statistically significant associations from this dataset
859 (https://figshare.com/articles/dataset/Table_S3_for_Aqil_et_al_2022/19606192). We hope this

860 will make it easier for the scientific community to use GWAS results than the currently available
861 datasets which store associations with each trait in a different table. Then, for each of the 4,863
862 deletions in our deletion dataset, we checked if any of the SNVs in LD were among the SNVs
863 that were significantly ($p < 10^{-8}$) associated with a phenotype. We then obtained the phenotype
864 that was associated with one of the SNVs in LD with the lowest p-value, and ascribed it to the
865 deletion. Thus, 433 (8.9%) of the 4,863 deletions in the deletion dataset had phenotypic
866 associations.

867

868 **Enrichment analysis for ancient deletions**

869 We performed enrichment analyses for phenotypic relevance and length among ancient
870 deletions using variants with a pooled frequency $> 5\%$ in YRI, CEU, and CHB combined. First,
871 we investigate whether a higher proportion of ancient deletions, relative to non-ancient
872 deletions, have phenotypic relevance. To this end, we defined phenotypic relevance in three
873 ways: 1) GWAS association, 2) exonic overlap, and 3) both GWAS association and exonic
874 overlap. For each definition, we first calculated the observed proportions of phenotypic deletions
875 among both ancient and non-ancient categories in turn. Then we shuffled the “ancient” and
876 “non-ancient” labels among the deletions in 10,000 permutations, calculating the proportion of
877 phenotypic deletions among both ancient and non-ancient labels for each permutation. Using
878 the number of permutations in which the difference in proportions of phenotypic deletions was
879 more extreme than the difference in observed proportions, we obtained an empirical p-value for
880 phenotypic enrichment among ancient deletions.

881

882 We also investigated whether certain phenotypic categories are overrepresented in ancient
883 deletions relative to non-ancient deletions. For this, we used the UKBioBank traits associated
884 with the deletions. In total, 1,675 traits were associated with the deletions in the deletion
885 dataset. We manually placed each of these traits into one of 18 categories such that any
886 deletion could be associated with one or more phenotypic categories (**Supplementary File 3**).
887 Only deletions with a pooled frequency $> 5\%$ in YRI, CEU, and CHB combined were retained for
888 analysis. For each phenotypic category, we obtained the proportion of deletions associated with
889 that category among ancient and non-ancient deletions. We then shuffled the “ancient” and
890 “non-ancient” labels in 10,000 permutations. Just as above, we used the number of
891 permutations in which the difference in proportions of deletions associated with a phenotypic
892 category was more extreme than the difference in observed proportions, and we obtained an
893 empirical p-value.

894

895 Now, we turn to length enrichment. We calculated the 2nd, 5th, 10th, 20th, 30th, 40th, 50th,
896 60th, 70th, 80th, 90th, 95th, and 98th length percentiles for deletions in ancient and non-ancient
897 categories. We calculated the differences in corresponding percentiles in ancient and non-
898 ancient deletions. Again, we shuffled the “ancient” and “non-ancient” labels into 10,000 random
899 permutations, calculating the differences in corresponding percentiles in ancient versus non-
900 ancient deletions for each permutation. This gave us empirical p-values for differences in the
901 length of ancient and non-ancient deletions at various percentiles.

902

903

904

905

906

907

908

909

909 **Simulations to identify signatures of overdominance**

910 We set out to identify signatures associated with a locus that has evolved under
911 overdominance.
912 Methodologically, we approached the problem of separating overdominance from neutrality on
913 two fronts: (i) The trajectory of the allele-frequency of the mutation conditioning on the age of
914 the mutation, and (ii) the patterns of neutral polymorphisms around the so-called focal mutation
915 which has evolved either under overdominance or neutrally (with the same age of the mutation).
916 Given that a mutation under overdominance (heterozygote advantage) may be at intermediate
917 frequency, we also (iii) studied simple coalescent simulations conditioning on the existence of at
918 least one SNV within a certain frequency range (e.g. 50%-60%). The age of a mutation is crucial
919 for the study of balancing selection. We considered two values for the age of a mutation: (i)
920 10,000 generation and (ii) 40,000 generations old. In the first case, the age of the mutation
921 corresponds to $10,000 \times 29 \text{ years} = 290,000 \text{ years}$ old mutation. In the second case, the
922 mutation is 1,160,000 years old. Conditioning on the age of the mutation, we generated allele
923 frequency trajectories of the mutation, *i.e.*, the frequency of the mutation at each time point from
924 its onset until the present-day. For the overdominance scenario, we used the software
925 *trajdemognpops*, implemented using tools from *ms* and *mssel* (kindly provided by R.R. Hudson),
926 to generate trajectories of a mutation under overdominance. The dominance coefficient is
927 characterized by a large value (here, $h = 10$) in order to assign a benefit for the heterozygote.
928 Thus, the fitness for genotypes at a biallelic locus is by:

929
930 **Aa:** $(1 + sh)$, where s is the selection coefficient and h the dominance coefficient. Here, $s =$
931 0.005 and $h = 10$, *i.e.*, the heterozygote Aa has a fitness value 1.05. If we had considered
932 codominance (as opposed to overdominance), we would have set $h = 0.5$.

933 **AA:** $1 + s$, thus, the fitness for the AA is 1.005

934 **aa:** 1.

935

936 Given the trajectory of a mutation, a population is split in two kinds of genotypes. The
937 haplotypes carrying the mutation (or the derived allele) and the haplotypes that carry the
938 ancestral allele. Each *neutral* (*i.e.* passenger mutation) can change “population” (genetic
939 background) by recombination. Therefore, the coalescent in this case is described as a
940 structured coalescent of two populations (a population with the derived allele and another
941 population with the ancestral allele) that communicate between themselves via recombination.
942 The size of each population is determined at each time point by the trajectory of the derived
943 allele.

944 In order to understand the effect of the age of the allele (and also to test the *mssel* code for
945 correctness), we performed coalescent simulations (using Hudson’s *ms*) conditioning on the
946 presence of at least one SNV at frequency within a given range. Since we are interested in the
947 mutations that are approximately at 50% frequency in the population, we conditioned on the
948 presence of derived mutations in the sample in 22-28 (out of 50) haplotypes. This set of
949 simulations are called *pseCoal*. For each simulated dataset, we calculated the relevant
950 population genetics statistics available through Comus (Papadantonakis et al., 2016) including
951 number of segregating sites, θ , π , Tajima’s D, ZnS, Fay and Wu’s H, dv_k , and dv_h . Then, we
952 conducted a summary of all these statistics using PCA. The goal is to understand whether the
953 different scenarios can get separated by using polymorphic patterns.

954

955

956

957 **Conceptual and methodological concerns:**

958

959 Human polymorphisms wherein derived allele is shared with archaic hominins are older than
960 700,000 years

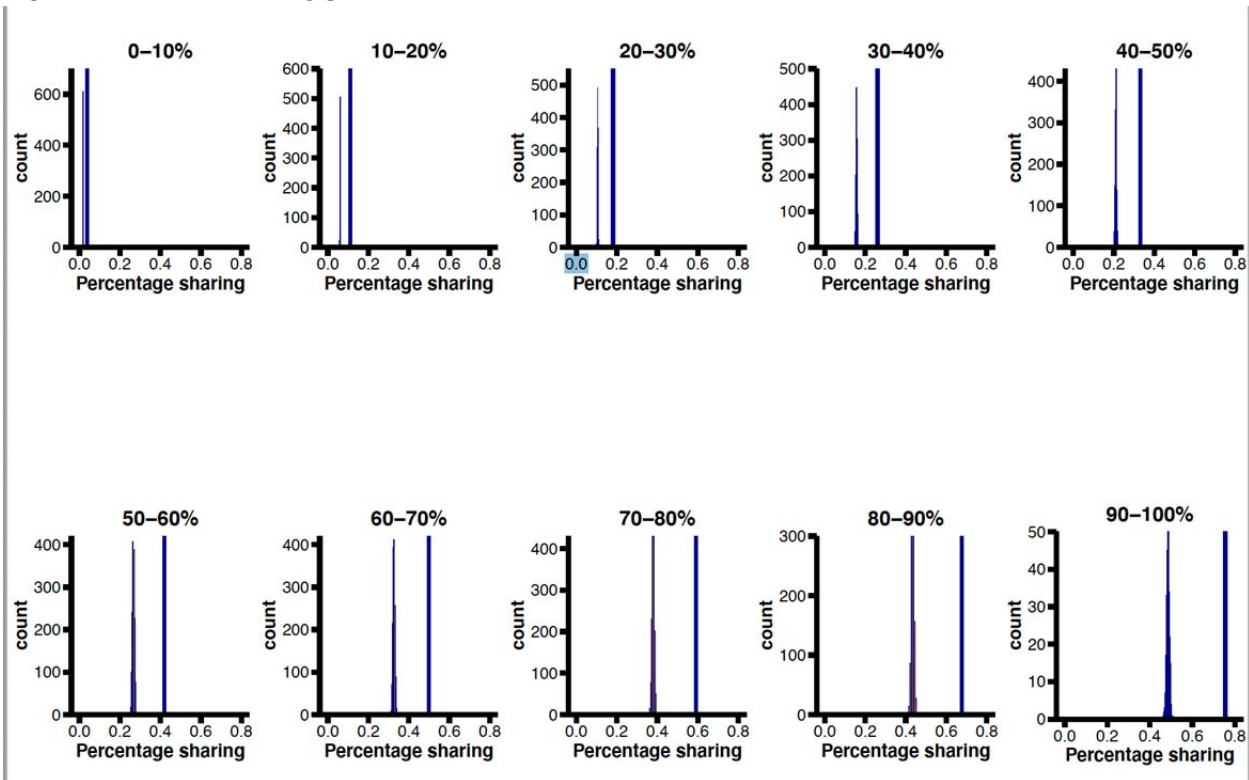
961 AMHs and archaic hominins are estimated to have diverged around ~700,000 years ago.
962 Therefore, if a human polymorphism has been maintained for more than ~700,000 years, it was
963 also present in the common ancestral population of AMH and archaic lineages. It follows that if
964 a polymorphism (the presence of both ancestral and derived alleles in a population) is present in
965 AMHs and archaic hominins, then (barring recurrence and introgression), by parsimony, the
966 polymorphism was also present in their common ancestral population (**Figure 1A**). Thus, a
967 polymorphism that is shared by common descent between AMHs and archaic hominins, has
968 been maintained for over ~700,000 years. Moreover, because the ancestral allele is fixed in
969 chimpanzees by definition, AMH polymorphisms wherein the derived allele is fixed in archaic
970 hominins were also present (in the polymorphic state) in the common ancestral population of
971 archaic hominins and AMHs. In essence, AMH polymorphisms for which archaic hominins carry
972 the derived allele (fixed or polymorphic) have been maintained for more than 700,000 years.
973

974
975 Why use SNVs (as opposed to deletion polymorphisms) for comparison of the real-life
976 proportion of ancient polymorphisms against neutrally simulated SNVs.
977 It is worth explaining why we used SNVs, instead of deletions – the class of variants that we
978 are interested in – for comparing observed versus simulated (under neutral conditions)
979 proportions of ancient polymorphisms. Deletions are not suitable for such a comparison
980 because, in general, they are targets of strong negative selection (Conrad et al., 2006;
981 Kondrashov, 2017; Lin et al., 2015; Lin and Gokcumen, 2019). Thus, negative selection would
982 have purged a large proportion of deletions that emerged in the common ancestral population of
983 AMHs and archaic hominins. It follows that a smaller proportion of AMH polymorphic deletions
984 than expected under neutral conditions will be shared by common descent with archaic
985 hominins. Even if balancing selection were inflating the proportion of deletions that are shared
986 with archaic hominins, it would not be observable due to the opposite deflationary effect of
987 negative selection. Since negative selection is not as strong a force in the evolution of SNVs as
988 it is for deletions, this problem would not be as pronounced if we used SNVs instead of
989 deletions for testing this premise. Hence, our choice of SNVs for this analysis. In the main text,,
990 we have expanded on the idea that deletions are targets of stronger negative selection.
991

992 The vast majority of human deletions are derived relative to chimpanzees
993 In order to identify deletion polymorphisms that have been maintained in the human lineage for
994 over 700,000 years, we focused on deletions that were present (either polymorphically or fixed)
995 in the four high coverage archaic genomes. This technique would work only for deletions that
996 were derived in humans, relative to chimps. However, variants that have been called as
997 deletions (wherein deletion is the alternative allele) in the 1000 Genomes project may, in fact,
998 be human-specific insertions, such that the reference allele (non-deletion) is derived. To
999 investigate how common this situation is among the 4,863 deletions in our dataset, we lifted
1000 over the the coordinates of the deletions from hg19 on to the Chimpanzee reference panTro3
1001 using the LiftOver tool in UCSC Genome Browser (Kent et al., 2002). If the liftover for a deletion
1002 fails on account of the window being completely or partially deleted in the Chimpanzee
1003 reference, it is indicative of the region being a human-specific insertion. The liftover failed for
1004 this reason for only 184 (3.8%) of the 4,863 deletions. Therefore, the vast majority of deletions
1005 are, in fact, derived relative to chimps.
1006

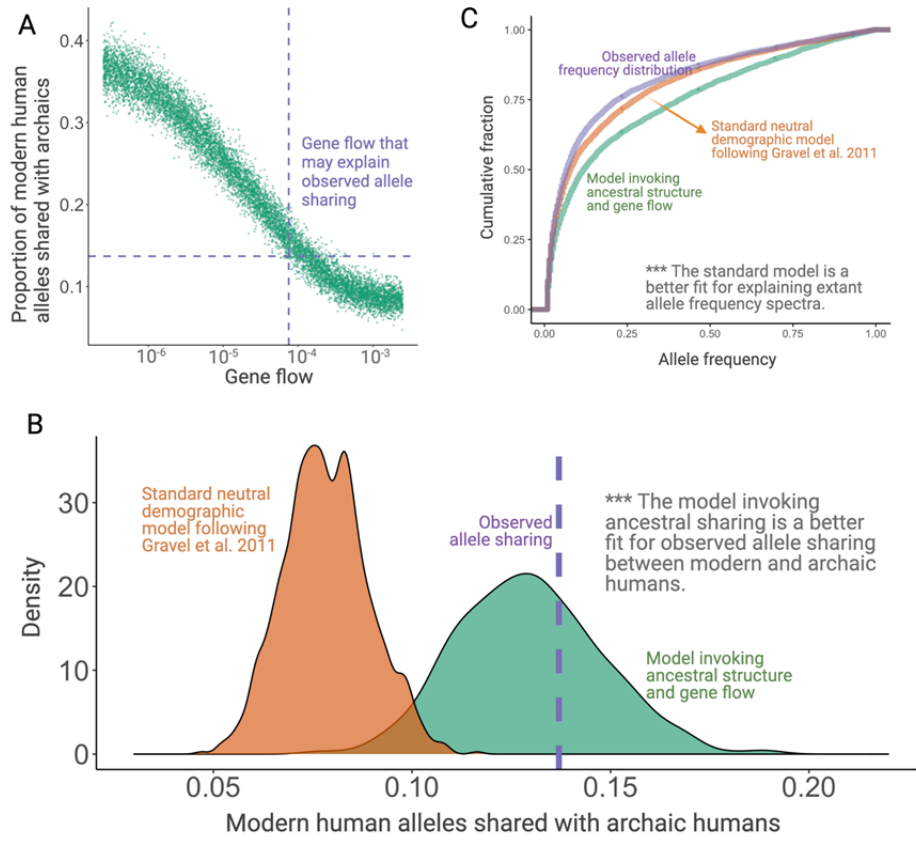
1007
1008
1009
1010
1011

1013 SUPPLEMENTARY FIGURES

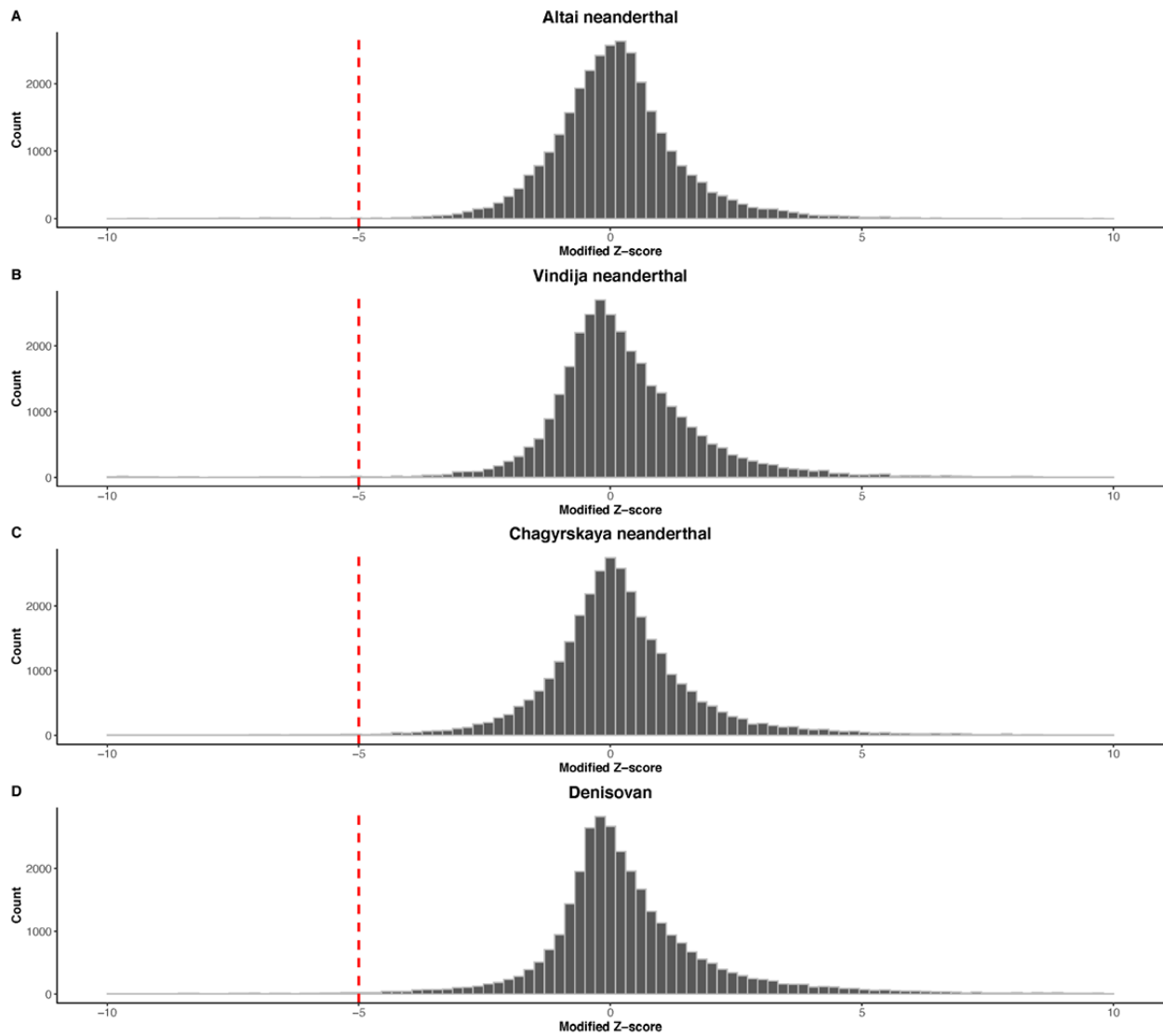


1014 **Figure 1—figure supplement 1. Proportion of ancient polymorphisms in observed data (YRI), relative to**
 1015 **neutral expectation (“base” model parameters) in various derived allele frequency bins. The vertical blue line**
 1016 **indicates the observed sharing, while the distributions are simulated expectations. The excess of ancient**
 1017 **polymorphisms in observed data becomes more pronounced at higher derived allele frequencies.**

1014
 1015
 1016
 1017
 1018
 1019
 1020

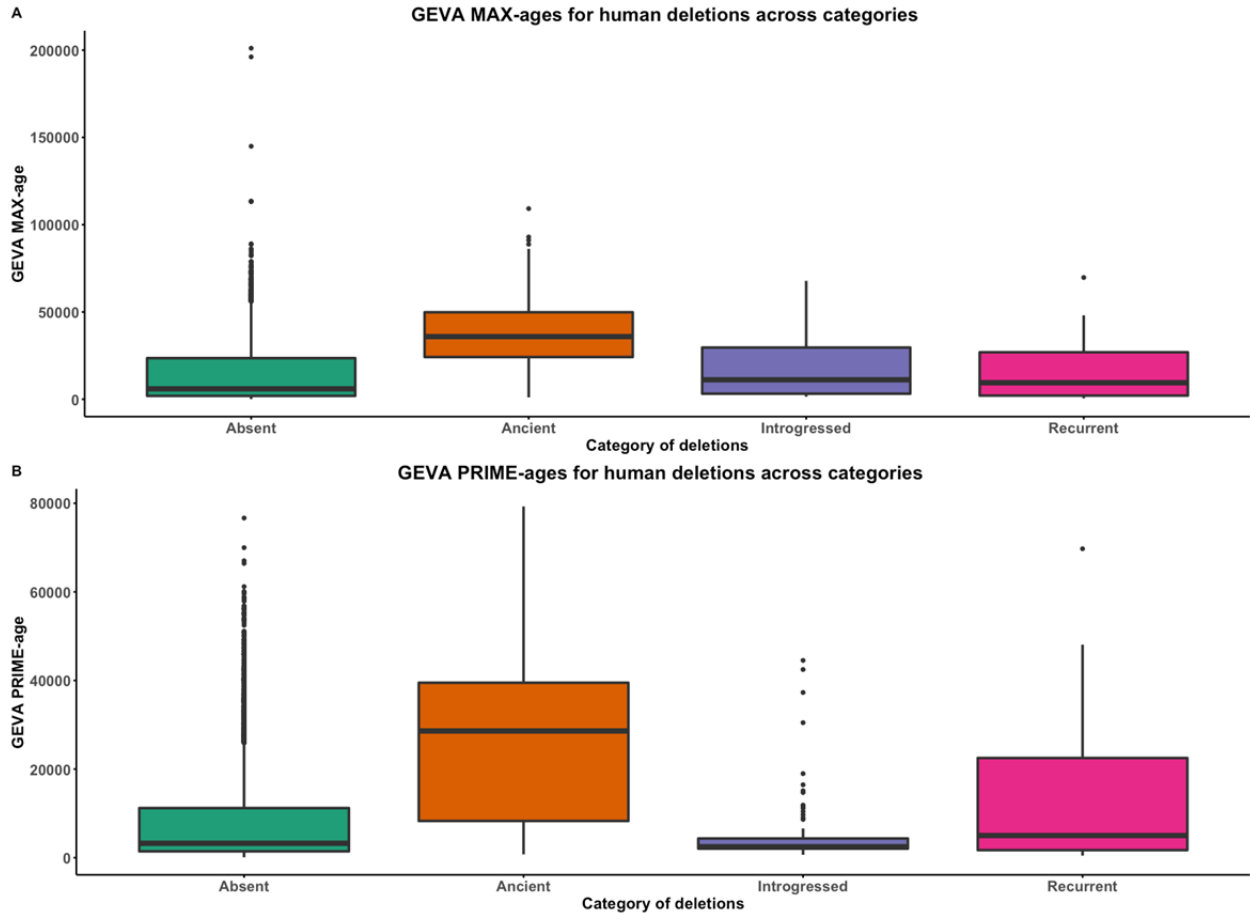


1021
 1022 **Figure 1—figure supplement 2. A)** Results from simulations invoking structure in the population that was ancestral
 1023 to both AMHs and archaic hominins. In this model, we have three latent subgroups in the ancestral populations. The
 1024 x-axis refers to the fraction of each subgroup that is formed by the migrants of each of the other subgroups in each
 1025 generation. **B.** Proportion of ancient polymorphisms in YRI. The purple line is the observed proportion of ancient
 1026 polymorphisms in Yoruba (YRI). The green and orange density plots indicate the distribution of the proportion of
 1027 ancient polymorphisms in neutral simulations with and without ancestral structure, respectively. We used Gravel et al.
 1028 parameters for these simulations. **C.** Comparison of the allele-frequency spectra of simulated SNVs with observed
 1029 SNVs. The purple, orange, and green lines represent allele frequency spectra in the YRI population using actual
 1030 SNVs, neutral simulations without ancestral structure, and neutral simulations invoking ancestral structure,
 1031 respectively.
 1032



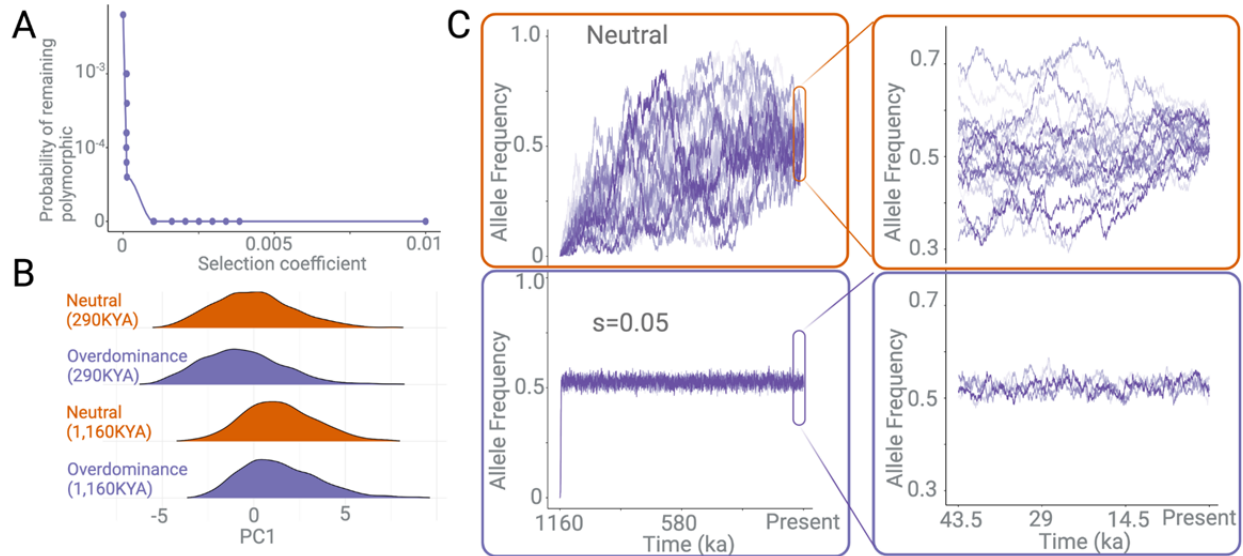
1033
 1034
 1035
 1036

Figure 2—figure supplement 1. Read-depth-based pipeline to identify deletions in archaic hominin genomes: Distribution of the modified Z-score of the read-depth across the 32,154 biallelic AMH deletions in the archaic genomes. **A.** Altai neanderthal. **B.** Vindija neanderthal. **C.** Chagyrskaya neanderthal. **D.** Denisovan.



1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045

Figure 3—figure supplement 1. GEVA ages of deletions across categories. Absent denotes polymorphic deletions in AMHs that are not present in any of the four high coverage archaic genomes. Introgressed refers to the shared deletions that were introgressed from archaic hominins into AMHs. Recurrent refers to the shared deletions that emerged independently in the AMH and archaic hominin lineages. Ancient refers to the AMH deletions that are shared with archaic hominins by common descent. **A.** GEVA PRIME-ages. **B.** GEVA MAX-ages. With both GEVA PRIME and GEVA MAX measures, we observe that ancient deletions are significantly older than absent, recurrent, and introgressed deletions. This implies that our pipeline to identify ancient deletions is sound.



1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058

Figure 7—figure supplement 1. **A.** The probability of a polymorphism persisting in the population for 1,000,000 years under different negative selection pressures. **B.** Density plots of the first principle component of multiple summary statistics based on variants simulated under neutral versus overdominance ($s=0.05$) scenarios. This is shown for two categories of variants: 1) those that emerged 290 kya, and 2) those that emerged 1,160 kya. There is no discernible difference between overdominance and neutrality within the time frame of these simulations. **C.** The allele frequency trajectories of variants over 1,000,000 years, under neutrally (top), versus under overdominance (bottom). The x-axis represents the time since the emergence of a variant in years, assuming a 29 year generation time. The right panel is a zoomed-in version of the same allele frequency trajectories in the last ~50 thousand years.

1059 **REFERENCES**

- 1060 1000 Genomes Project Consortium, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM,
1061 Korbel JO, Marchini JL, McCarthy S, McVean GA, Abecasis GR. 2015. A global reference
1062 for human genetic variation. *Nature* **526**:68–74.
- 1063 Abdul-Rahman F, Tranchina D, Gresham D. 2021. Fluctuating Environments Maintain Genetic
1064 Diversity through Neutral Fitness Effects and Balancing Selection. *Mol Biol Evol* **38**:4362–
1065 4375.
- 1066 Agarwal V, Kommaddi RP, Valli K, Ryder D, Hyde TM, Kleinman JE, Strobel HW, Ravindranath
1067 V. 2008. Drug metabolism in human brain: high levels of cytochrome P4503A43 in brain
1068 and metabolism of anti-anxiety drug alprazolam to its active metabolite. *PLoS One* **3**:e2337.
- 1069 Aho AV, Kernighan BW, Weinberger PJ. 1978. Awk, a Pattern Scanning and Processing
1070 Language. Bell Laboratories.
- 1071 Albers PK, McVean G. 2020. Dating genomic variants and shared ancestry in population-scale
1072 sequencing data. *PLoS Biol* **18**:e3000586.
- 1073 Alharbi AF, Sheng N, Nicol K, Strömberg N, Hollox EJ. 2022. Balancing selection at the human
1074 salivary agglutinin gene (DMBT1) driven by host-microbe interactions. *iScience* **25**:104189.
- 1075 Allison AC. 1954a. Protection afforded by sickle-cell trait against subtertian malarial infection.
1076 *Br Med J* **1**:290–294.
- 1077 Allison AC. 1954b. The distribution of the sickle-cell trait in East Africa and elsewhere, and its
1078 apparent relationship to the incidence of subtertian malaria. *Trans R Soc Trop Med Hyg*
1079 **48**:312–318.
- 1080 Andrés AM, Dennis MY, Kretzschmar WW, Cannons JL, Lee-Lin S-Q, Hurle B, NISC
1081 Comparative Sequencing Program, Schwartzberg PL, Williamson SH, Bustamante CD,
1082 Nielsen R, Clark AG, Green ED. 2010. Balancing selection maintains a form of ERAP2 that
1083 undergoes nonsense-mediated decay and affects antigen presentation. *PLoS Genet*
1084 **6**:e1001157.
- 1085 Barghi N, Hermisson J, Schlötterer C. 2020. Polygenic adaptation: a unifying framework to
1086 understand positive selection. *Nat Rev Genet* **21**:769–781.
- 1087 Benton ML, Abraham A, LaBella AL, Abbot P, Rokas A, Capra JA. 2021. The influence of
1088 evolutionary history on human health and disease. *Nat Rev Genet* **22**:269–283.
- 1089 Berg JJ, Coop G. 2014. A population genetic signal of polygenic adaptation. *PLoS Genet*
1090 **10**:e1004412.
- 1091 Bergström A, Stringer C, Hajdinjak M, Scerri EML, Skoglund P. 2021. Origins of modern human
1092 ancestry. *Nature* **590**:229–237.
- 1093 Bitarello BD, de Filippo C, Teixeira JC, Schmidt JM, Kleinert P, Meyer D, Andrés AM. 2018.
1094 Signatures of Long-Term Balancing Selection in Human Genomes. *Genome Biol Evol*
1095 **10**:939–955.
- 1096 Bromberg Y, Kahn PC, Rost B. 2013. Neutral and weakly nonneutral sequence variants may
1097 define individuality. *Proc Natl Acad Sci U S A* **110**:14255–14260.
- 1098 Conrad DF, Andrews TD, Carter NP, Hurles ME, Pritchard JK. 2006. A high-resolution survey of
1099 deletion polymorphism in the human genome. *Nat Genet* **38**:75–81.
- 1100 Conrad DF, Pinto D, Redon R, Feuk L, Gokcumen O, Zhang Y, Aerts J, Andrews TD, Barnes C,
1101 Campbell P, Fitzgerald T, Hu M, Ihm CH, Kristiansson K, Macarthur DG, Macdonald JR,
1102 Onyiah I, Pang AWC, Robson S, Stirrups K, Valsesia A, Walter K, Wei J, Wellcome Trust
1103 Case Control Consortium, Tyler-Smith C, Carter NP, Lee C, Scherer SW, Hurles ME. 2010.
1104 Origins and functional impact of copy number variation in the human genome. *Nature*
1105 **464**:704–712.
- 1106 Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE, Lunter G,
1107 Marth GT, Sherry ST, Others. 2011. 1000 Genomes Project Analysis Group. The variant
1108 call format and vcftools. *Bioinformatics* **27**:2156–2158.
- 1109 DeGiorgio M, Lohmueller KE, Nielsen R. 2014. A model-based approach for identifying

1110 signatures of ancient balancing selection in genetic data. *PLoS Genet* **10**:e1004561.
 1111 de Guzman Strong C, Conlan S, Deming CB, Cheng J, Sears KE, Segre JA. 2010. A milieu of
 1112 regulatory elements in the epidermal differentiation complex syntenic block: implications for
 1113 atopic dermatitis and psoriasis. *Hum Mol Genet* **19**:1453–1460.
 1114 Dudley JT, Kim Y, Liu L, Markov GJ, Gerold K, Chen R, Butte AJ, Kumar S. 2012. Human
 1115 genomic disease variants: a neutral evolutionary explanation. *Genome Res* **22**:1383–1394.
 1116 Edge MD, Coop G. 2019. Reconstructing the History of Polygenic Scores Using Coalescent
 1117 Trees. *Genetics* **211**:235–262.
 1118 Eyre-Walker A. 2010. Genetic architecture of a complex trait and its implications for fitness and
 1119 genome-wide association studies. *Proc Natl Acad Sci U S A* **107**:1752–1756.
 1120 Fenner JN. 2005. Cross-cultural estimation of the human generation interval for use in genetics-
 1121 based population divergence studies. *Am J Phys Anthropol* **128**:415–423.
 1122 Fijarczyk A, Babik W. 2015. Detecting balancing selection in genomes: limits and prospects. *Mol*
 1123 *Ecol* **24**:3529–3545.
 1124 Fisher RA. 1922. XXI.—On the Dominance Ratio. *Proceedings of the Royal Society of*
 1125 *Edinburgh* **42**:321–341.
 1126 Gibson G. 2018. Population genetics and GWAS: A primer. *PLoS Biol.*
 1127 Gokcumen O, Zhu Q, Mulder LCF, Iskow RC, Austermann C, Scharer CD, Raj T, Boss JM,
 1128 Sunyaev S, Price A, Stranger B, Simon V, Lee C. 2013. Balancing Selection on a
 1129 Regulatory Region Exhibiting Ancient Variation That Predates Human–Neandertal
 1130 Divergence. *PLoS Genet* **9**:e1003404.
 1131 Gravel S, Henn BM, Gutenkunst RN, Indap AR, Marth GT, Clark AG, Yu F, Gibbs RA, 1000
 1132 Genomes Project, Bustamante CD. 2011. Demographic history and rare allele sharing
 1133 among human populations. *Proc Natl Acad Sci U S A* **108**:11983–11988.
 1134 Hamid I, Korunes KL, Beleza S, Goldberg A. 2021. Rapid adaptation to malaria facilitated by
 1135 admixture in the human population of Cabo Verde. *Elife* **10**. doi:10.7554/eLife.63177
 1136 Hedrick PW. 2011. Population genetics of malaria resistance in humans. *Heredity* **107**:283–
 1137 304.
 1138 Hedrick PW. 2006. Genetic Polymorphism in Heterogeneous Environments: The Age of
 1139 Genomics. *Annu Rev Ecol Evol Syst* **37**:67–93.
 1140 Hudson RR. 2002. Generating samples under a Wright-Fisher neutral model of genetic
 1141 variation. *Bioinformatics*. doi:10.1093/bioinformatics/18.2.337
 1142 Iglewicz B, Hoaglin DC. 1993. How to detect and handle outliers. Asq Press.
 1143 Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, Haussler D. 2002. The
 1144 human genome browser at UCSC. *Genome Res* **12**:996–1006.
 1145 Kharitonov A, Chen Z, Sures I, Wang H, Schilling J, Ullrich A. 1997. A family of proteins that
 1146 inhibit signalling through tyrosine kinase receptors. *Nature* **386**:181–186.
 1147 Klunk J, Vilgalys TP, Demeure CE, Cheng X, Shiratori M, Madej J, Beau R, Elli D, Patino MI,
 1148 Redfern R, DeWitte SN, Gamble JA, Boldsen JL, Carmichael A, Varlik N, Eaton K, Grenier
 1149 J-C, Golding GB, Devault A, Rouillard J-M, Yotova V, Sindeaux R, Ye CJ, Bikaran M,
 1150 Dumaine A, Brinkworth JF, Missiakas D, Rouleau GA, Steinrücken M, Pizarro-Cerdá J,
 1151 Poinar HN, Barreiro LB. 2022. Evolution of immune genes is associated with the Black
 1152 Death. *Nature*. doi:10.1038/s41586-022-05349-x
 1153 Kondrashov AS. 2017. Crumbling Genome: The Impact of Deleterious Mutations on Humans.
 1154 John Wiley & Sons.
 1155 Lachance J, Tishkoff SA. 2013. Population Genomics of Human Adaptation. *Annu Rev Ecol*
 1156 *Evol Syst* **44**:123–143.
 1157 Langergraber KE, Prüfer K, Rowney C, Boesch C, Crockford C, Fawcett K, Inoue E, Inoue-
 1158 Muruyama M, Mitani JC, Muller MN, Robbins MM, Schubert G, Stoinski TS, Viola B, Watts
 1159 D, Wittig RM, Wrangham RW, Zuberbühler K, Pääbo S, Vigilant L. 2012. Generation times
 1160 in wild chimpanzees and gorillas suggest earlier divergence times in great ape and human

1161 evolution. *Proc Natl Acad Sci U S A* **109**:15716–15721.

1162 Lederberg J. 1999. JBS Haldane (1949) on infectious disease and evolution. *Genetics* **153**:1–3.

1163 Leffler EM, Gao Z, Pfeifer S, Ségurel L, Auton A, Venn O, Bowden R, Bontrop R, Wall JD, Sella

1164 G, Donnelly P, McVean G, Przeworski M. 2013. Multiple instances of ancient balancing

1165 selection shared between humans and chimpanzees. *Science* **339**:1578–1582.

1166 Lettre G. 2014. Rare and low-frequency variants in human common diseases and other

1167 complex traits. *J Med Genet* **51**:705–714.

1168 Levene H. 1953. Genetic Equilibrium When More Than One Ecological Niche is Available. *Am*

1169 *Nat* **87**:331–333.

1170 Li H, Durbin R. 2011. Inference of human population history from individual whole-genome

1171 sequences. *Nature* **475**:493–496.

1172 Lin Y-L, Gokcumen O. 2019. Fine-Scale Characterization of Genomic Structural Variation in the

1173 Human Genome Reveals Adaptive and Biomedically Relevant Hotspots. *Genome Biol Evol*

1174 **11**:1136–1151.

1175 Lin Y-L, Pavlidis P, Karakoc E, Ajay J, Gokcumen O. 2015. The evolution and functional impact

1176 of human deletion variants shared with archaic hominin genomes. *Mol Biol Evol* **32**:1008–

1177 1019.

1178 Loos RJF. 2020. 15 years of genome-wide association studies and no signs of slowing down.

1179 *Nat Commun* **11**:5900.

1180 Mafessoni F, Grote S, de Filippo C, Slon V, Kolobova KA, Viola B, Markin SV, Chintalapati M,

1181 Peyrégne S, Skov L, Skoglund P, Krivoschapkin AI, Derevianko AP, Meyer M, Kelso J, Peter

1182 B, Prüfer K, Pääbo S. 2020. A high-coverage Neandertal genome from Chagyrskaya Cave.

1183 *Proc Natl Acad Sci U S A* **117**:15132–15136.

1184 Mallick S, Li H, Lipson M, Mathieson I, Gymrek M, Racimo F, Zhao M, Chennagiri N, Nordenfelt

1185 S, Tandon A, Skoglund P, Lazaridis I, Sankararaman S, Fu Q, Rohland N, Renaud G,

1186 Erlich Y, Willems T, Gallo C, Spence JP, Song YS, Poletti G, Balloux F, Van Driem G, De

1187 Knijff P, Romero IG, Jha AR, Behar DM, Bravi CM, Capelli C, Hervig T, Moreno-Estrada A,

1188 Posukh OL, Balanovska E, Balanovsky O, Karachanak-Yankova S, Sahakyan H, Toncheva

1189 D, Yepiskoposyan L, Tyler-Smith C, Xue Y, Abdullah MS, Ruiz-Linares A, Beall CM, Di

1190 Rienzo A, Jeong C, Starikovskaya EB, Metspalu E, Parik J, Villems R, Henn BM,

1191 Hodoglugil U, Mahley R, Sajantila A, Stamatoyannopoulos G, Wee JTS, Khusainova R,

1192 Khusnutdinova E, Litvinov S, Ayodo G, Comas D, Hammer MF, Kivisild T, Klitz W, Winkler

1193 CA, Labuda D, Bamshad M, Jorde LB, Tishkoff SA, Watkins WS, Metspalu M, Dryomov S,

1194 Sukernik R, Singh L, Thangaraj K, Pääbo S, Kelso J, Patterson N, Reich D. 2016. The

1195 Simons Genome Diversity Project: 300 genomes from 142 diverse populations. *Nature*

1196 **538**:201–206.

1197 Mathieson S, Mathieson I. 2018. FADS1 and the Timing of Human Adaptation to Agriculture.

1198 *Mol Biol Evol* **35**:2957–2970.

1199 McArthur E, Rinker D, Capra JA. 2020. Quantifying the contribution of Neanderthal introgression

1200 to the heritability of complex traits. *bioRxiv*. doi:10.1101/2020.06.08.140087

1201 Mendoza-Revilla J, Chacón-Duque JC, Fuentes-Guajardo M, Ormond L, Wang K, Hurtado M,

1202 Villegas V, Granja V, Acuña-Alonzo V, Jaramillo C, Arias W, Lozano RB, Gómez-Valdés J,

1203 Villamil-Ramírez H, de Cerqueira CCS, Badillo Rivera KM, Nieves-Colón MA, Gignoux CR,

1204 Wojcik GL, Moreno-Estrada A, Hunemeier T, Ramallo V, Schuler-Faccini L, Gonzalez-José

1205 R, Bortolini M-C, Canizales-Quinteros S, Gallo C, Poletti G, Bedoya G, Rothhammer F,

1206 Balding D, Fumagalli M, Adhikari K, Ruiz-Linares A, Hellenthal G. 2021. Disentangling

1207 signatures of selection before and after European colonization in Latin Americans. *bioRxiv*.

1208 doi:10.1101/2021.11.15.467418

1209 Meyer M, Kircher M, Gansauge M-T, Li H, Racimo F, Mallick S, Schraiber JG, Jay F, Prüfer K,

1210 de Filippo C, Sudmant PH, Alkan C, Fu Q, Do R, Rohland N, Tandon A, Siebauer M, Green

1211 RE, Bryc K, Briggs AW, Stenzel U, Dabney J, Shendure J, Kitzman J, Hammer MF,

1212 Shunkov MV, Derevianko AP, Patterson N, Andrés AM, Eichler EE, Slatkin M, Reich D,
1213 Kelso J, Pääbo S. 2012. A High-Coverage Genome Sequence from an Archaic Denisovan
1214 Individual. *Science* **338**:222–226.

1215 Muller HJ. 1918. Genetic Variability, Twin Hybrids and Constant Hybrids, in a Case of Balanced
1216 Lethal Factors. *Genetics* **3**:422–499.

1217 Noonan JP, Coop G, Kudaravalli S, Smith D, Krause J, Alessi J, Chen F, Platt D, Pääbo S,
1218 Pritchard JK, Rubin EM. 2006. Sequencing and analysis of Neanderthal genomic DNA.
1219 *Science* **314**:1113–1118.

1220 Pajic P, Lin Y-L, Xu D, Gokcumen O. 2016. The psoriasis-associated deletion of late cornified
1221 envelope genes LCE3B and LCE3C has been maintained under balancing selection since
1222 Human Denisovan divergence. *BMC Evol Biol* **16**:265.

1223 Papadantonakis S, Poirazi P, Pavlidis P. 2016. CoMuS: Simulating coalescent histories and
1224 polymorphic data from multiple species. *Mol Ecol Resour*. doi:10.1111/1755-0998.12544

1225 Posth C, Wißing C, Kitagawa K, Pagani L, van Holstein L, Racimo F, Wehrberger K, Conard NJ,
1226 Kind CJ, Bocherens H, Krause J. 2017. Deeply divergent archaic mitochondrial genome
1227 provides lower time boundary for African gene flow into Neanderthals. *Nat Commun*
1228 **8**:16046.

1229 Pritchard JK, Pickrell JK, Coop G. 2010. The genetics of human adaptation: hard sweeps, soft
1230 sweeps, and polygenic adaptation. *Curr Biol* **20**:R208–15.

1231 Prüfer K, de Filippo C, Grote S, Mafessoni F, Korlević P, Hajdinjak M, Vernot B, Skov L, Hsieh
1232 P, Peyrégne S, Reher D, Hopfe C, Nagel S, Maricic T, Fu Q, Theunert C, Rogers R,
1233 Skoglund P, Chintalapati M, Dannemann M, Nelson BJ, Key FM, Rudan P, Kučan Ž, Gušić
1234 I, Golovanova LV, Doronichev VB, Patterson N, Reich D, Eichler EE, Slatkin M, Schierup
1235 MH, Andrés AM, Kelso J, Meyer M, Pääbo S. 2017. A high-coverage Neandertal genome
1236 from Vindija Cave in Croatia. *Science* **358**:655–658.

1237 Prüfer K, Racimo F, Patterson N, Jay F, Sankararaman S, Sawyer S, Heinze A, Renaud G,
1238 Sudmant PH, de Filippo C, Li H, Mallick S, Dannemann M, Fu Q, Kircher M, Kuhlwiilm M,
1239 Lachmann M, Meyer M, Ongyerth M, Siebauer M, Theunert C, Tandon A, Moorjani P,
1240 Pickrell J, Mullikin JC, Vohr SH, Green RE, Hellmann I, Johnson PLF, Blanche H, Cann H,
1241 Kitzman JO, Shendure J, Eichler EE, Lein ES, Bakken TE, Golovanova LV, Doronichev VB,
1242 Shunkov MV, Derevianko AP, Viola B, Slatkin M, Reich D, Kelso J, Pääbo S. 2014. The
1243 complete genome sequence of a Neanderthal from the Altai Mountains. *Nature* **505**:43–49.

1244 Qiu Q-W, Wu D-D, Yu L-H, Yan T-Z, Zhang W, Li Z-T, Liu Y-H, Zhang Y-P, Xu X-M. 2013.
1245 Evidence of recent natural selection on the Southeast Asian deletion (--SEA) causing α -
1246 thalassemia in South China. *BMC Evol Biol* **13**:63.

1247 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic
1248 features. *Bioinformatics* **26**:841–842.

1249 Saitou M, Gokcumen O. 2020. An Evolutionary Perspective on the Impact of Genomic Copy
1250 Number Variation on Human Health. *J Mol Evol* **88**:104–119.

1251 Saitou M, Masuda N, Gokcumen O. 2021a. Similarity-based analysis of allele frequency
1252 distribution among multiple populations identifies adaptive genomic structural variants. *Mol*
1253 *Biol Evol*. doi:10.1093/molbev/msab313

1254 Saitou M, Resendez S, Pradhan AJ, Wu F, Lie NC, Hall NJ, Zhu Q, Reinholdt L, Satta Y,
1255 Speidel L, Nakagome S, Hanchard NA, Churchill G, Lee C, Atilla-Gokcumen GE, Mu X,
1256 Gokcumen O. 2021b. Sex-specific phenotypic effects and evolutionary history of an ancient
1257 polymorphic deletion of the human growth hormone receptor. *Sci Adv* **7**:eabi4476.

1258 Sella G, Barton NH. 2019. Thinking About the Evolution of Complex Traits in the Era of
1259 Genome-Wide Association Studies. *Annu Rev Genomics Hum Genet* **20**:461–493.

1260 Siewert KM, Voight BF. 2020. BetaScan2: Standardized Statistics to Detect Balancing Selection
1261 Utilizing Substitution Data. *Genome Biol Evol* **12**:3873–3877.

1262 Siewert KM, Voight BF. 2017. Detecting Long-Term Balancing Selection Using Allele Frequency

1263 Correlation. *Mol Biol Evol* **34**:2996–3005.

1264 Smith Maynard J, Smith JM, Smith JM. 1998. *Evolutionary Genetics*. Oxford University Press.

1265 Soni V, Vos M, Eyre-Walker A. 2022. A new test suggests hundreds of amino acid
1266 polymorphisms in humans are subject to balancing selection. *PLoS Biol* **20**:e3001645.

1267 Speidel L, Cassidy L, Davies RW, Hellenthal G, Skoglund P, Myers SR. 2021. Inferring
1268 Population Histories for Ancient Genomes Using Genome-Wide Genealogies. *Mol Biol Evol*
1269 **38**:3497–3511.

1270 Speidel L, Forest M, Shi S, Myers SR. 2019. A method for genome-wide genealogy estimation
1271 for thousands of samples. *Nat Genet* **51**:1321–1329.

1272 Starr I, Seiffert-Sinha K, Sinha AA, Gokcumen O. 2021. Evolutionary context of psoriatic
1273 immune skin response. *Evolution, Medicine, and Public Health*. doi:10.1093/emph/eoab042

1274 Takahashi Y, Kawata M. 2013. A comprehensive test for negative frequency-dependent
1275 selection. *Popul Ecol* **55**:499–509.

1276 Taskent O, Lin YL, Patramanis I, Pavlidis P, Gokcumen O. 2020. Analysis of Haplotypic
1277 Variation and Deletion Polymorphisms Point to Multiple Archaic Introgression Events,
1278 Including from Altai Neanderthal Lineage. *Genetics* **215**:497–509.

1279 Taskent RO, Alioglu ND, Fer E, Melike Donertas H, Somel M, Gokcumen O. 2017. Variation
1280 and Functional Impact of Neanderthal Ancestry in Western Asia. *Genome Biol Evol*
1281 **9**:3516–3524.

1282 Van der Auwera GA, O'Connor BD. 2020. *Genomics in the Cloud: Using Docker, GATK, and*
1283 *WDL in Terra*. "O'Reilly Media, Inc."

1284 Vernot B, Akey JM. 2014. Resurrecting surviving Neanderthal lineages from modern human
1285 genomes. *Science* **343**:1017–1021.

1286 Wallace B. 1970. Genetic load: its biological and conceptual aspects. *Genetic load: its biological*
1287 *and conceptual aspects*.

1288 Wittmann MJ, Bergland AO, Feldman MW, Schmidt PS, Petrov DA. 2017. Seasonally
1289 fluctuating selection can maintain polymorphism at many loci via segregation lift. *Proc Natl*
1290 *Acad Sci U S A* **114**:E9932–E9941.

1291 Xue Y, Sun D, Daly A, Yang F, Zhou X, Zhao M, Huang N, Zerjal T, Lee C, Carter NP, Hurles
1292 ME, Tyler-Smith C. 2008. Adaptive evolution of UGT2B17 copy-number variation. *Am J*
1293 *Hum Genet* **83**:337–346.

1294 Zeng J, de Vlaming R, Wu Y, Robinson MR, Lloyd-Jones LR, Yengo L, Yap CX, Xue A,
1295 Sidorenko J, McRae AF, Powell JE, Montgomery GW, Metspalu A, Esko T, Gibson G, Wray
1296 NR, Visscher PM, Yang J. 2018. Signatures of negative selection in the genetic architecture
1297 of human complex traits. *Nat Genet* **50**:746–753.

1298