

SHARING OF WEAK SIGNALS OF POSITIVE SELECTION ACROSS THE GENOME

Nathan S. Harris and Alan R. Rogers

April 22, 2020

Abstract

Selection in humans often leaves subtle signatures at individual loci. Few studies have measured the extent to which these signals are shared among human populations. Here a new method is developed to compare weak signals of selection in aggregate across the genome using the 1000 Genomes Phase 3 Data. Results presented here show that selection producing weak selection serves to increase population differences around coding areas of the genome.

1 Introduction and background

Until relatively recently, studies of natural selection in humans focused on classic selective sweeps that have large effects on isolated regions of the genome (Sabeti et al., 2002; Voight et al., 2006). In a classic selective sweep, a new beneficial mutation appears in one person and spreads through a population (Smith and Haigh, 1974). When an allele is sweeping through the population, surrounding DNA from the original haplotype on which the mutation occurred tends to “hitchhike” with the selected allele. This results in linkage disequilibrium (LD), a nonrandom association of alleles at two or more loci (Lewontin and Kojima, 1960). The blocks around selected loci are longer and contain less diversity the greater the strength of selection. Mutation and recombination reintroduce variation into blocks of LD (Lande, 1975), and given sufficient generations following the sweep, LD blocks around a locus are broken apart. The extent to which this has occurred depends on the local mutation and recombination rates and the amount of time that has passed. This “signal” is used by a variety of methods to detect natural selection (Booker et al., 2017; Haas and Payseur, 2016; Vitti et al., 2013).

As research continues it has become apparent that the most common forms of selection in humans are those that have smaller effects on LD around individual loci such as polygenic adaptation (Daub et al., 2013; Hernandez et al.,

2011; Pritchard et al., 2010) and selection on existing variation (Harris et al., 2018; Schrider and Kern, 2017). These forms of selection are unlikely to generate significant signals using statistics designed for classic sweeps, although they may account for some fraction of those that fail to reach significance. Furthermore, while classic selective sweeps are usually geographically local and therefore tend to increase population differences (Fagny et al., 2014; Vitti et al., 2013), much less is known about how often weak signals of selection are shared between populations. Weak signals of selection are more likely to be shared between populations because the types of weak selection that produce them are slower, and alleles that arose in a common ancestor are more likely to still be polymorphic. In the context of this paper, a “weak” signal refers to signals of selection that do not reach significance or an Integrated Haplotype Score ($|iHS|$) greater than two. This term is therefore relative to the statistic being calculated rather than reflecting a rigid category of selective pressure.

Theoretically, populations ought to share more signals of selection if they are closely related for the following reasons:

1. Beneficial mutations are more likely to become lost or fixed the more time that has elapsed since mutation. If two populations are distantly related, beneficial mutations are likely to be either fixed or lost in one or both of the populations. If two populations are closely related, signals from completed sweeps in their common ancestor are more likely to be preserved and detectable in each. Furthermore, the same beneficial mutation may still be polymorphic and increasing in frequency in both populations, producing a shared signal.
2. Closely related populations often share similar environments. If each population experiences a beneficial mutation near the same locus, both populations may show evidence of selection in the same region.
3. Neutral mechanisms can produce spurious signals of selection by chance. Closely related populations are expected to share such signals because a larger portion of their population history is shared.

This theoretical expectation has been supported empirically. Pickrell et al. (2009) show that populations within the same continent are more likely to share signals of selection. Similarly, Johnson and Voight (2018) found that regions of the genome with high concentrations of large $|iHS|$ scores were more likely to overlap between populations if those populations are closely related.

Here, methods traditionally used to detect classic selective sweeps are implemented to characterize genome-wide patterns of shared weak signals of selection. While methods have been developed to detect subtler signatures of selection (Field et al., 2016; Schrider and Kern, 2016), some of the methods for detecting classic selective sweeps still have some utility for studying population differences. Hard selective sweeps may be an uncommon mechanism of adaptation in humans (Coop et al., 2009; Hernandez et al., 2011; Schrider and Kern, 2017), but they are an important case. Many large adaptive changes have

76 occurred in the recent past via classic selective sweeps (Fagny et al., 2014; Vitti
77 et al., 2013). These signals are known to be recent because ancient adaptation
78 that occurred via this mechanism has already driven the causative variant to
79 fixation and the signal has been obscured by subsequent recombination and
80 mutation. It may be possible to detect older instances of selection with some of
81 the same methods because moderately beneficial alleles increase in frequency
82 much slower and the resulting signal persists longer. However, because $|iHS|$
83 is standardized, and most large significant signals are the result of selection
84 within the last 20,000 years (Voight et al., 2006), ancient signals are unlikely to
85 produce significant signals of selection at individual loci. For this reason, this
86 research investigates genome-wide patterns of weak signals of selection rather
87 than identifying particular loci under selection.

88 **2 Results**

89 **2.1 Weak signals of selection**

90 To get an idea of the relevant strength of selection, a catalog of significant $|iHS|$
91 signals for the 1000 Genomes Phase 3 data were obtained from Johnson and
92 Voight (2018). These signals were binned by their size to search for any clear
93 cutoffs in the size (in base pairs) of selection signals from $|iHS|$ (Figure 1). Most
94 significant $|iHS|$ signals were at least 100kb long. A model was adapted from
95 Gillespie (2004) to determine the relationship between the selection coefficient
96 (s), the recombination rate (r), and the size of linkage disequilibrium blocks
97 around a selected allele at an intermediate frequency of 0.5 (Figure 2). The
98 strongest signals of classic sweeps will therefore commonly have a selection
99 coefficient of 0.01 or greater. Here we hoped to exclude the majority of these
100 signals by removing sites with $|iHS|$ values greater than two. The remaining
101 sites should disproportionately be from loci with selection coefficients smaller
102 than 0.01.

103 **2.2 The Integrated Haplotype Score ($|iHS|$)**

104 To measure signals of selection, $|iHS|$ was calculated for samples from the 1000
105 Genomes Project. $|iHS|$ identifies regions under selection by comparing the
106 difference in LD between carriers of the reference and alternate alleles. Large
107 $|iHS|$ scores indicate a substantial difference in LD. Correlation of sample $|iHS|$
108 scores of two populations was calculated for each pair of samples. This correla-
109 tion was calculated separately for genic and nongenic portions of the genome.
110 Unlike genetic drift, selection affects specific loci and linked variation rather
111 than the entire genome. Genic regions should more commonly be the target
112 of selection because they are more often functional (Barreiro et al., 2008; Coop
113 et al., 2009). The difference between genic and nongenic correlations at a given
114 value of genetic distance is likely the result of selection. While selection is
115 known to occur in some noncoding regions (Forni et al., 2014; Hernandez et al.,

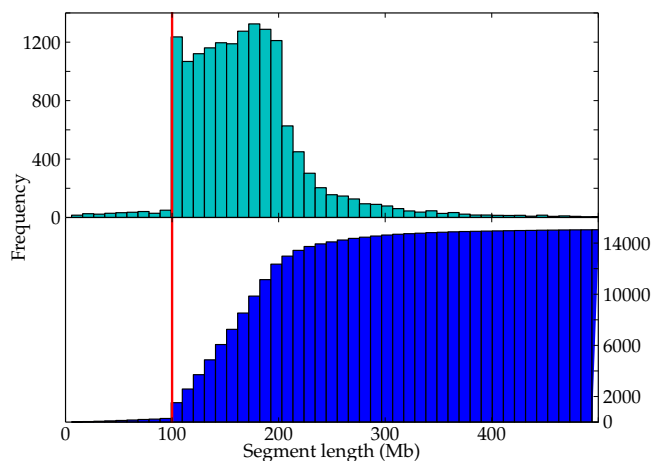


Figure 1: Normal and cumulative histogram of significant $|iHS|$ region sizes in 26 1000 Genomes Phase 3 data. A sharp cliff occurs at 100kb, implying most cases of classic selective sweeps considered to be significant leave signals grearelevantter than 100kb.

116 2011; Ponjavic et al., 2007), this should have the effect of making any observed
117 differences between genic and nongenic regions more conservative. In either
118 case, correlation is expected to be relatively large and positive when the two
119 populations have both experienced selection in the same areas of the genome.
120 This occurs not only because particular SNPs may be under selection, but also
121 because $|iHS|$ scores at linked sites near mutually selected loci should covary
122 as well.

123 To limit the effect of hard sweeps, loci with significant $|iHS|$ values (greater
124 than two) were removed from the calculation. To compensate for the effects
125 of linkage, regions were considered nongenic if they were at least 500kb from
126 genic regions. The resulting correlations were regressed with a Loess algorithm
127 against Nei's genetic distance for each pair of populations. Confidence inter-
128 vals were generated using a moving block bootstrap (Liu and Singh, 1992) with
129 a block size of 500kb. The outcome is shown in Figure 3.

130 The left edge of the graph refers to pairs of populations that are genetically
131 similar and tend also to be geographically close. Samples from the same geo-
132 graphic regions experience similar correlations in genic and nongenic sections
133 of their genome. As the genetic distance between samples increases, the cor-
134 relation between samples in both genic and nongenic regions decreases until
135 it approaches zero in the most genetically distant comparisons, Africans and
136 Eastern Asians. However, the correlation of $|iHS|$ in genic and nongenic re-
137 gions does not decrease at the same rate. The pairwise genic correlation is
138 consistently lower than the pairwise nongenic correlation for all sample pairs

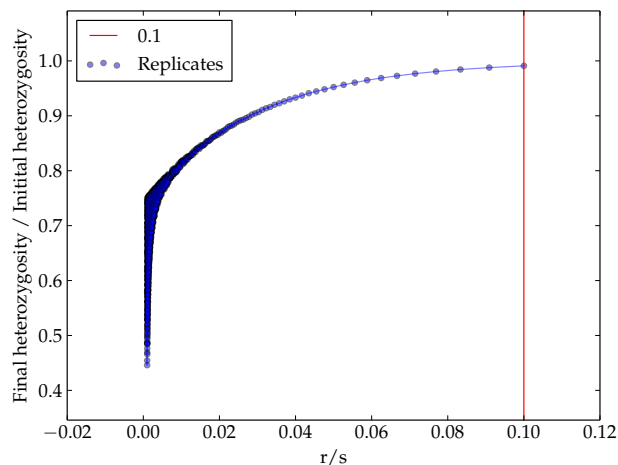


Figure 2: Determining the relationship between the strength of selection and LD block size. The recombination rate (r) is held constant while the selection coefficient (s) varies between replicates. Each replicate is allowed to run until the beneficial allele reaches a frequency of 0.5. The red line indicates the point at which sites are no longer in LD with the site under selection. With known r and block size, s can be determined.

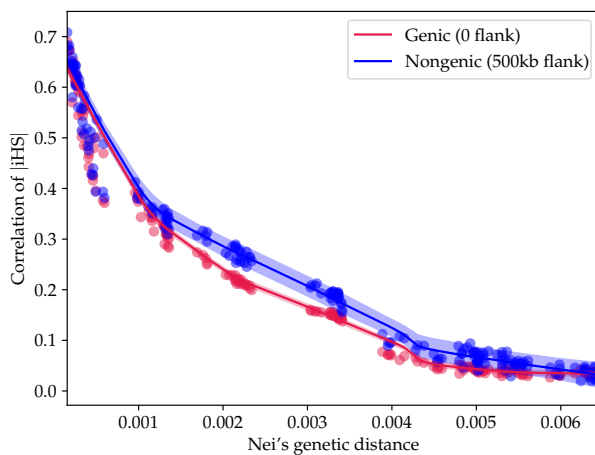


Figure 3: Correlation of $|iHS|$ scores for nongenic regions compared to genic regions.

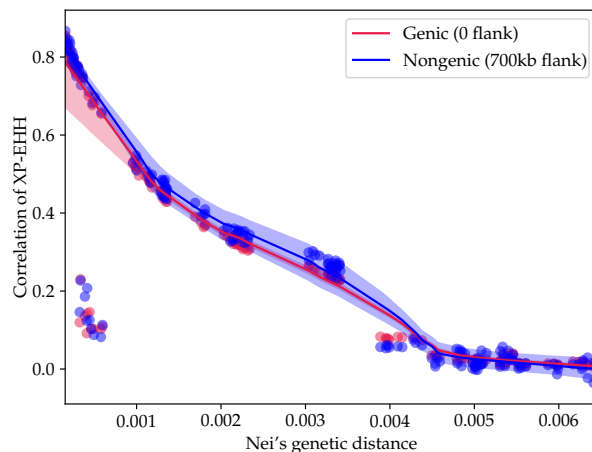


Figure 4: Correlation of XP-EHH scores for nongenic regions compared to genic regions.

139 except at large genetic distances.

140 **2.3 The Cross Population Extended Haplotype** 141 **Homozygosity (XP-EHH)**

142 The above analysis was repeated for XP-EHH. XP-EHH compares LD differ-
143 ences between two samples, identifying where selection has occurred in one
144 sample but not the other. XP-EHH was chosen as a supplementary analysis
145 to $|iHS|$ because it is more sensitive to selection in which the beneficial is near
146 fixation. Like $|iHS|$, XP-EHH correlation between populations decreases with
147 increasing genetic distance. Unlike iHS , XP-EHH genic and nongenic corre-
148 lations overlap considerably (see Figure 4). However, XP-EHH correlations
149 in either nongenic or genic regions were consistently larger than their $|iHS|$
150 counterparts.

151 **2.4 Simulations in SLiM**

152 Correlations tend to be lower in genic than in nongenic regions. I will ar-
153 gue that this implies that, even when selection is weak, its primary effect is
154 disruptive, tending to increase differences between populations. To establish
155 this point, simulations were conducted using Selection on Linked Mutations
156 (SLiM) to show that disruptive selection tends to reduce the correlation be-
157 tween $|iHS|$ signals. In each simulation, a single population splits into two.
158 The time of this split varies among simulations to model differences in genetic

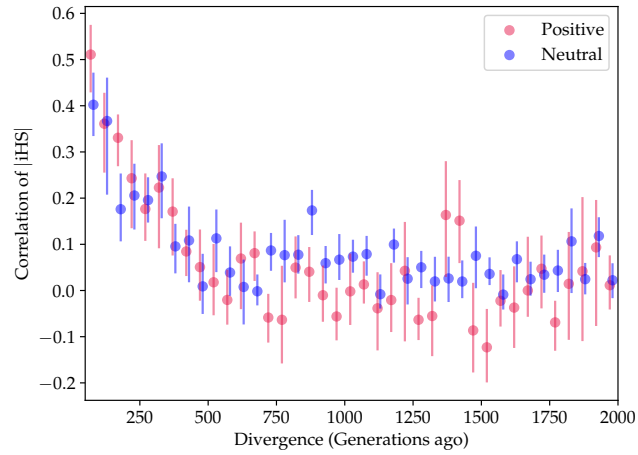


Figure 5: |iHS| correlation with positive selection occurring in one branch. Wilcoxon signed-rank test, $n=21$, $p=0.002$.

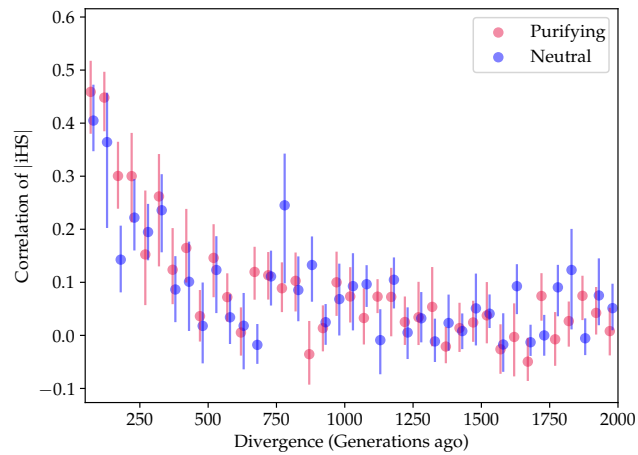


Figure 6: |iHS| correlation with purifying selection occurring in one branch. Wilcoxon signed-rank test, $n=21$, $p=0.357$.

159 distance. Following the split, populations experience one of three evolution-
160 ary scenarios: neutrality, positive selection, or purifying selection. All three
161 models experience neutral mutations. In both cases of selection, non-neutral
162 mutations occur in a single population following the split. Results were stan-
163 dardized against the neutral simulations with the same divergence time, and
164 the correlation analysis proceeds in the same manner as the real data.

165 The results of this process are shown in Figure 5 and 6. For any particu-
166 lar divergence time, confidence intervals of neutral and selection models over-
167 lapped. However, there is a difference between the scenarios when points are
168 considered together. In each case a Wilcoxon signed-rank test was conducted
169 to assess if the selection correlation could have been drawn from the same dis-
170 tribution as the neutral correlation. The correlation of $|iHS|$ in the presence
171 of positive selection was significantly lower than in the neutral model ($n=21$,
172 $p=0.002$). Previous work has suggested that $|iHS|$ is not sensitive to purifying
173 selection (Enard et al., 2014) and that conclusion was supported here. The cor-
174 relation of $|iHS|$ in the presence of purifying selection was indistinguishable
175 from the neutral model ($n=21$, $p=0.357$).

176 To get an idea of how far back selection can be detected, the basic model
177 of positive selection was repeated for increasing divergence times. In these
178 simulations, beneficial mutations ($s=0.01$) are introduced into one population
179 following the population split for 2,000 generations. At this point, no more
180 beneficial mutations are introduced and selection on existing beneficial muta-
181 tions is halted, allowing any existing signals to decay. When divergence times

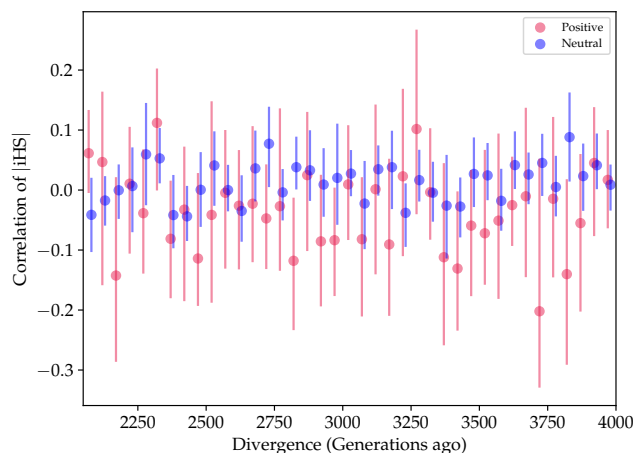


Figure 7: Correlation of $|iHS|$ for a set of divergence times ranging from 2,000 to 4,000 generations ago. While neutral regions have a correlation around zero, selection has the effect of decreasing correlation further.

182 range from 2,000 to 4,000 generations ago (about 50-100-kya) (Figure 7), corre-
183 lations are small relative to the recent divergence times discussed above, but
184 correlations in selected regions are still significantly smaller than in neutral re-
185 gions ($n=40$, $p=2.07e-4$). If divergence times are increased further (Figure 8), the
186 significant difference between selected and neutral regions disappears ($n=40$,
187 $p=0.83$). This strongly suggests that the correlation of $|iHS|$ is sensitive to in-
188 stances of selection that are much older than the effective range of traditional
189 use of $|iHS|$ (Voight et al., 2006).

190 Simulations of positive selection were repeated using a model of soft selec-
191 tive sweeps, in which a mutation becomes beneficial after it has already drifted
192 to a relatively high frequency. Soft sweeps introduce initial allele frequency as
193 an additional dimension to the simulations. While $|iHS|$ is sensitive to soft
194 sweeps with starting allele frequencies as large as 0.1 (Ferrer-Admetlla et al.,
195 2014), it is unclear to what extent soft selective sweeps will be detected by the
196 methods proposed above. To determine the relevant range of soft sweep pa-
197 rameters, simulations were run over a wide range of combinations of starting
198 allele frequencies and sweep frequency. Soft sweeps, in general, appear to have
199 the same effect as classic sweeps from de novo mutations, depressing the corre-
200 lation of $|iHS|$ in regions that have experienced selection. The largest allele fre-
201 quency for which soft sweeps caused a significant difference between selected
202 and neutral regions in our simulations was the 0.09-0.1 bin ($n=40$, $p=0.041$). The
203 effect size is small in both cases due to their proximity to the cutoff. Figure 9
204 shows the results of the simulations using the 0.09-0.1 bin. Figure 10 ($n=40$,

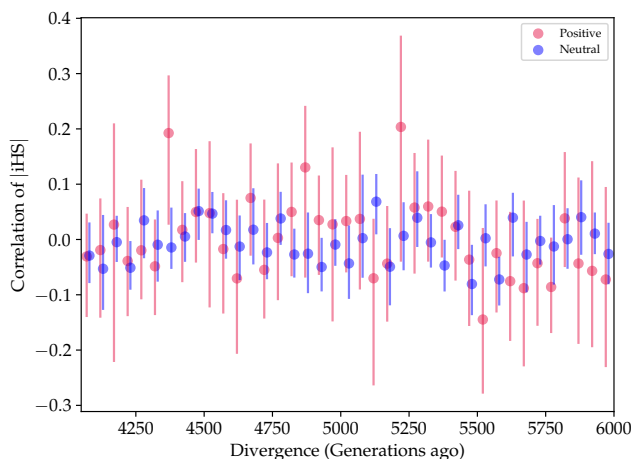


Figure 8: Correlation of $|iHS|$ for a set of divergence times ranging from 4,000 to 6,000 generations ago. No significant difference between neutral and selected regions is detectable.

205 $p=0.113$) shows the result of simulations in the first insignificant frequency bin
206 of 0.1-0.2.

207 3 Discussion

208 3.1 Weak signals differ between populations

209 The smaller correlation of $|iHS|$ in genic regions compared to nongenic regions
210 implies that populations share few signals of weak selection. This result sup-
211 ports the hypothesis that weak positive selection has a similar disruptive effect
212 as hard classic sweeps. This pattern is also present in the simulated data.

213 Simulations including beneficial mutations in one population resulted in
214 a depressed correlation of $|iHS|$ compared to neutral simulations. Models of
215 purifying selection did not have the same effect. This suggests that positive
216 selection rather than purifying selection across genic regions is more likely to
217 be driving the increase in population differences.

218 This effect is amplified when the distance from genic regions increases. If
219 nongenic regions are at least 1Mb from genic regions, confidence intervals for
220 correlation of nongenic regions increase, but the difference between genic cor-
221 relation and nongenic correlation of $|iHS|$ increases substantially (Figure 11).
222 The sample size as measured by the number of bases and number of regions
223 in the genome both decrease. This change implies two things. First, it sup-
224 ports the notion that nongenic regions are substantially affected by selection at
225 neighboring genic sites. Second, the pattern in the nongenic regions continues

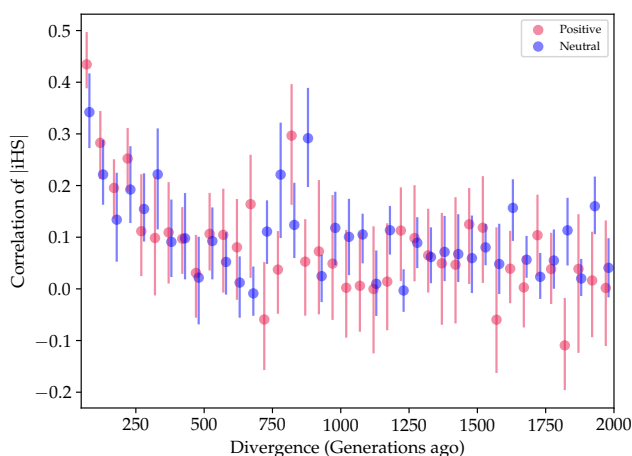


Figure 9: $|iHS|$ correlation with soft-sweeps occurring in one branch with a starting allele frequency between 0.09 and 0.1.

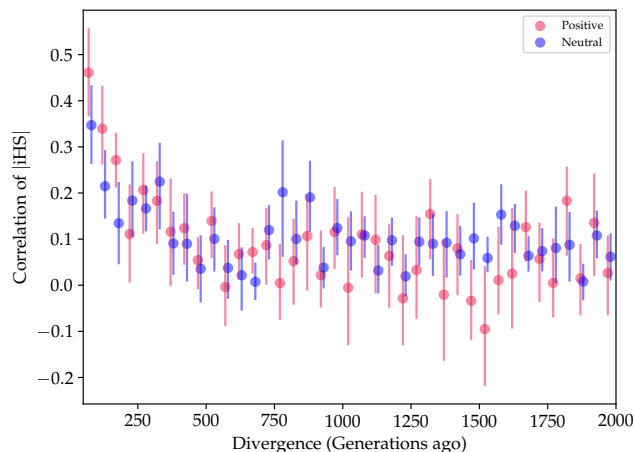


Figure 10: $|iHS|$ correlation with soft-sweeps occurring in one branch with a starting allele frequency between 0.1 and 0.2.

226 to change with increasing distance from genic regions, even to the point that
227 nongenic regions are reduced to a small portion of the genome. This result
228 could be used to support the hypothesis the majority of the genome is affected
229 by positive selection to some extent (e.g., Pouyet et al. (2018); Schrider and
230 Kern (2017)).

231 3.2 XP-EHH

232 The correlation of XP-EHH in genic regions was consistently lower than in non-
233 genic regions, but not significantly. The hypothesis that the patterns in genic
234 regions are due to population history rather than selection cannot be rejected.
235 In other words, the results are consistent with theoretical arguments 1, 2, and
236 3 enumerated above. However, the XP-EHH results are still informative when
237 compared to $|iHS|$. The difference in sensitivity between $|iHS|$ and XP-EHH is
238 visible in Figure 12. Pairwise XP-EHH correlation is larger in closely related
239 populations than pairwise $|iHS|$ correlation. This reflects the ability of XP-EHH
240 to detect differences in haplotype structure that are the result of older selection
241 (near fixation) or population history. This implies that the correlation methods
242 implemented here may be repeated for other methods of detecting selection.

243 Correlations of XP-EHH scores between African populations are substan-
244 tially small compared to population comparisons with similar genetic distance.
245 This values cluster in the lower left of Figure 12. This is likely due to the in-
246 creased similarity of these samples to the *reference* XP-EHH sample, Yoruba.

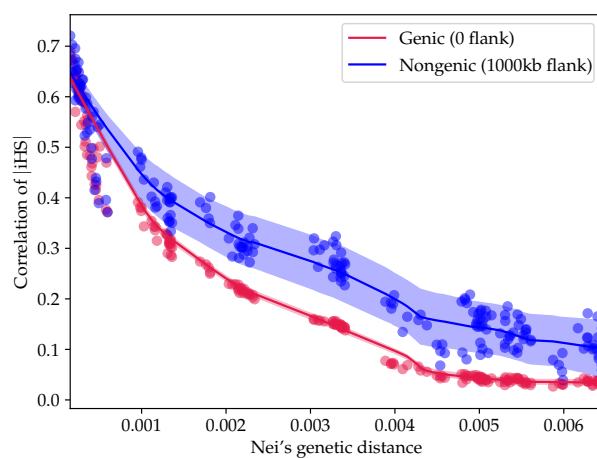


Figure 11: Genic regions compared to nongenic regions with a large flank size of 1Mb.

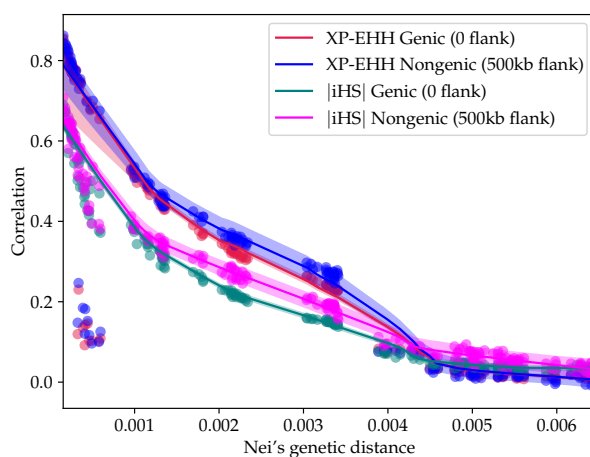


Figure 12: Correlation $|iHS|$ and correlation of XP-EHH scores. Genic regions with population labels can be found in Appendix A.

247 4 Methods

248 4.1 Pairwise correlation of |iHS|

249 To measure selection, the Integrated Haplotype Score (|iHS|) was calculated
250 for the 1000 Genomes Phase 3 data acquired from [ftp://ftp.1000ge-
251 nomes.ebi.ac.uk/vol1/ftp/release/20130502/](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/). Recently admixed populations were
252 excluded from this analysis, leaving 21 samples from Eurasian and African
253 populations. Data were divided by population sample and filtered to exclude
254 sites with a minor allele frequency less than 0.05.

255 |iHS| was calculated using the *selscan* package (Szpiech and Hernandez,
256 2014). |iHS| is sensitive to variation in allele frequency (Voight et al., 2006) and
257 local recombination rate (O'Reilly et al., 2008). This was compensated for in
258 two ways. First, when integrated across the chromosome, genetic map distance
259 was used rather than physical distance. Genetic maps for the 1000 genomes
260 were downloaded from the Pickrell lab ([https://github.com/joepickrell/
261 1000-genomes-genetic-maps](https://github.com/joepickrell/1000-genomes-genetic-maps)). The map positions for missing sites were im-
262 puted from neighboring sites. Second, when the ratio of scores is taken be-
263 tween the two allele types, the effects of recombination ought to disappear
264 because the recombination rate is the same for carriers of both alleles. To ver-
265 ify this process, |iHS| was first standardized using allele frequency bins and
266 then was regressed against recombination rate. A strong relationship between
267 frequency standardized |iHS| and recombination rate remained. Following the
268 example of Johnson and Voight (2018), |iHS| scores were restandardized using
269 46 frequency and 21 recombination bins. Frequency bins ranged from 0.05 to
270 0.95 and recombination bins were determined by grouping the data into per-
271 centiles.

272 Standardized |iHS| scores were split into genic and nongenic regions of
273 the genome for each sample using coordinates of known genes and gene pre-
274 dictions from the UCSC table browser (Karolchik et al., 2004). Dividing the
275 data into genic and nongenic regions allows us to distinguish between shared
276 sweeps in a common ancestor or independent sweeps after a common ancestor
277 from spurious signals of selection from common ancestry (theoretical points 1
278 and 2 from 3 above). Unlike genetic drift, selection affects specific loci and
279 linked variation rather than the entire genome. Genic regions should more
280 commonly be the target of selection because they are more often functional
281 (Barreiro et al., 2008; Coop et al., 2009). The difference between genic and non-
282 genic correlations at a given value of genetic distance can be attributed to se-
283 lection.

284 Regions were considered nongenic if they occurred outside of a flanking re-
285 gion around genes. This cutoff is used to compensate for the effects of linkage
286 (Slatkin, 2008; Wall and Pritchard, 2003). To determine the appropriate flank
287 size, correlation of |iHS| was calculated between populations for each subdivi-
288 sion of the genome and increasing flank size. The absolute value of iHS is
289 taken because the sign of iHS only indicates allele state. A beneficial allele will
290 produce both positive and negative scores with large magnitudes at neighbor-

291 ing sites. Correlation was limited to loci with $|iHS|$ values within two standard
 292 deviations of zero. This distinction was made to eliminate loci showing po-
 293 tential evidence for strong selective sweeps. The ideal flank size for nongenic
 294 regions was assessed based on sample size and the effect of linked genic sites
 295 (Table 1). As flanking regions become large, both sample size and the influence
 296 of genic regions on nongenic regions decrease. Within genic regions the size of
 297 the flanking region had little effect on $|iHS|$ correlations (Figure 13). There-
 298 fore, “genic” in this study is defined to mean a flank size of zero. For nongenic
 299 regions, a flank size of 500kb was used to balance between sample size and the
 300 effect of neighboring genic regions (Figure 14).

301 Once flank sizes were determined, genic and nongenic $|iHS|$ correlation
 302 matrices of the determined flank size were regressed against Nei’s genetic dis-
 303 tance using the Loess algorithm. Nei’s genetic distance (Nei, 1972) was cal-
 304 culated for each pair of population samples using the allele frequencies taken
 305 from the 1000 Genomes data.

306 A moving blocks bootstrap (Liu and Singh, 1992) was used to generate con-
 307 fidence intervals around the regressions. This method was chosen because loci
 308 near one another are used in each other’s calculation of $|iHS|$, implying that
 309 $|iHS|$ calculation of neighboring sites is not independent. The moving blocks
 310 bootstrap method compensates for this problem by sampling entire regions of
 311 the genome rather than individual loci. Blocks of 500 kilobases (kb) were sam-
 312 pled from the standardized $|iHS|$ output. This cutoff was chosen because most
 313 blocks of LD in the genome are smaller than 500kb (Slatkin, 2008; Wall and
 314 Pritchard, 2003). The number of blocks used in each bootstrap is equal to the
 315 number of blocks required to simulate the length of the real data. Correlation

Table 1: The flank sizes tested for genic and nongenic regions. Genic regions are given flanking regions and nongenic regions are considered to be anywhere not included.

Type	Genic flank size (kb)	Regions	Total bases
genic	0	21,531	1,281,434,774
	100	1,588	2,160,261,718
	200	7,458	2,374,098,102
	300	465	2,485,778,284
nongenic	300	445	256,313,467
	400	280	187,125,998
	500	193	142,103,389
	600	131	111,130,841
	700	97	89,398,183
	800	66	73,840,405
	900	47	63,052,677
	1,000	36	54,964,754

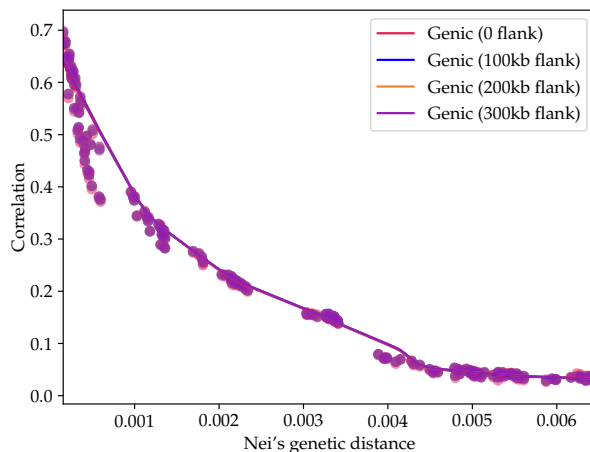


Figure 13: Difference in correlation of $|iHS|$ in genic regions given varying flank sizes.

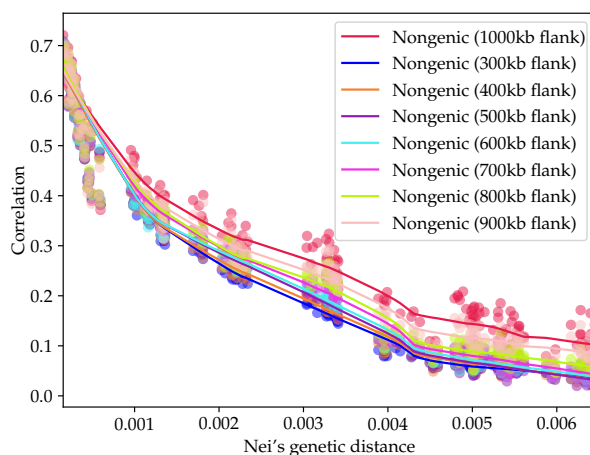


Figure 14: Difference in correlation of $|iHS|$ in nongenic regions given varying flank sizes.

316 of $|iHS|$ was calculated for each resampling. The model fit to the real data is
317 then applied to the resampling of the data. This process is repeated 1,000 times.
318 The inner 95% of these replicates become the confidence intervals for the real
319 data.

320 4.2 Pairwise correlation of XP-EHH

321 Tests performed on $|iHS|$ were repeated for XP-EHH. XP-EHH requires a sec-
322 ond population to serve as a comparison. Results of XP-EHH indicate where
323 selection has occurred in one population or the other, but not both. Each test
324 used the same second or *reference* population, the Yoruba. This will bias the
325 XP-EHH results against signals that populations share with the reference pop-
326 ulation. However, it allows for the pairwise comparison of nonreference pop-
327 ulations. XP-EHH was calculated using the *selscan* package (Szpiech and Her-
328 nandez, 2014). XP-EHH scores were standardized using 46 frequency and 21
329 recombination bins. Frequency bins ranged from 0.05 to 0.95 and recombina-
330 tion bins were determined by grouping the data into percentiles. A moving
331 blocks bootstrap (Liu and Singh, 1992) was used to generate confidence inter-
332 vals following the same methods described for $|iHS|$ above.

333 The correlation of XP-EHH was then calculated from the standardized val-
334 ues. Final standardized scores were filtered in both populations to exclude
335 sites that showed evidence for selection in Yoruba. The inclusion of these sites
336 would bias the results. Covariance between populations would be positive in
337 regions where Yoruba experienced selection that was relatively strong com-
338 pared to the populations being compared. In these regions, both populations
339 would have negative XP-EHH scores with relatively large magnitudes. This
340 increase in covariance would be artificial, rather than reflecting any difference
341 in the populations being compared. The pairwise XP-EHH correlations within
342 Africa clustered at lower correlations than other within continent comparisons.
343 This was visible in Figure 12. Figure 15 shows the same set of results with
344 Africa excluded. The trend in the data does not change in any statistically sig-
345 nificant way.

346 4.3 Simulation

347 Simulated data were generated using SLiM (Haller and Messer, 2018; Messer,
348 2013). SLiM is a forward-time simulator that allows researchers to model evo-
349 lutionary scenarios in the presence of linkage. For this work, three basic sim-
350 ulations were performed: a neutral model, a model of purifying selection, and
351 a model of positive selection. Instead of attempting to model human history,
352 a simple model was constructed in which a single population separates into
353 two populations at a prespecified time. This divergence time was adjusted to
354 values between 50 and 2000 generations. A constant population size of 10,000
355 was used throughout the simulation. The recombination rate was held con-
356 stant across the simulated genome to eliminate any possibility of an association
357 between linkage disequilibrium and recombination rate. Parameter values can
358 be found in Table 2.

359 In the neutral case, mutations have no effect. In both cases of selection,
360 selection occurs in one population following a population split. This creates a
361 scenario in which all loci under selection in one population were neutral in the
362 other. Beneficial or deleterious mutations constituted 5% of the total number

Table 2: Values used in the simulations

Model	Population size	Split time	μ	r	s	Starting p
Hard	10,000	50 to 2000	2.36e-8	1e-8	0.005 or -0.005	0.00005
Soft	10,000	50 to 2000	2.36e-8	1e-8	0.01	0.0001 to 0.5

363 of mutations in their respective simulations and have a selection coefficient of
364 0.005 and -0.005, respectively. $|iHS|$ is standardized to eliminate the effects of
365 allele frequency differences caused by drift. In real data, the entire genome is
366 standardized together, and results indicate how exceptional a particular site is
367 given its allele frequency. To replicate this effect in the simulated data, neutral
368 and non-neutral simulations with the same divergence time were standardized
369 jointly for allele frequency.

370 For each simulation, a moving blocks bootstrap was used to find confidence
371 intervals. A block size of 500 kb was used in the simulations to be consistent
372 with the analysis performed on the real data. Each simulation is independent
373 allowing the use of a sign-rank test. The Wilcoxon signed-rank test was per-
374 formed between the selection and neutral models.

375 Simulations of soft sweeps differed from hard sweeps models by using a
376 constant selective coefficient of 0.01, and varied beginning allele frequencies.
377 Soft sweeps start following the population split but occur at manually speci-
378 fied loci meeting the desired allele frequency. This occurs at user-specified in-

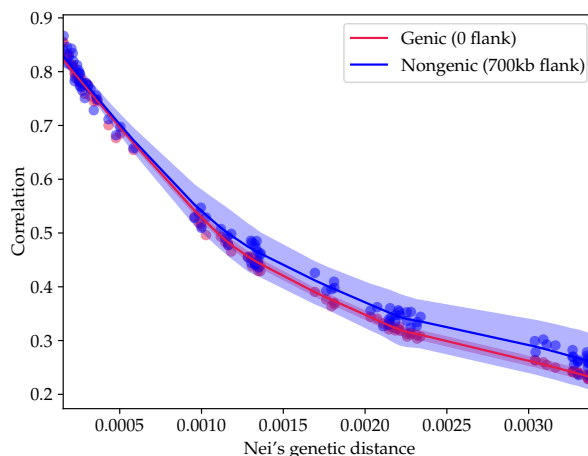


Figure 15: Correlation of XP-EHH scores omitting African samples.

379 tervals. In hard sweep models, most mutations, including the beneficial ones,
380 are lost to drift. Soft sweeps starting at higher allele frequencies are unlikely to
381 be lost. This presents a problem because the absolute number of soft sweeps
382 affects the result. Therefore these simulations were run varying the number of
383 introduced of soft sweeps until an allele frequency cutoff was observed.

384 References

- 385 Barreiro, L. B., Laval, G., Quach, H., Patin, E., and Quintana-Murci, L. (2008).
386 Natural selection has driven population differentiation in modern humans.
387 *Nature genetics*, 40:340.
- 388 Booker, T. R., Jackson, B. C., and Keightley, P. D. (2017). Detecting positive
389 selection in the genome. *BMC biology*, 15(1):98.
- 390 Coop, G., Pickrell, J. K., Novembre, J., Kudaravalli, S., Li, J., Absher, D., Myers,
391 R. M., Cavalli-Sforza, L. L., Feldman, M. W., and Pritchard, J. K. (2009). The
392 role of geography in human adaptation. *PLoS genetics*, 5(6):e1000500.
- 393 Daub, J. T., Hofer, T., Cutivet, E., Dupanloup, I., Quintana-Murci, L., Robinson-
394 Rechavi, M., and Excoffier, L. (2013). Evidence for polygenic adaptation to
395 pathogens in the human genome. *Molecular biology and evolution*, 30(7):1544–
396 1558.
- 397 Enard, D., Messer, P. W., and Petrov, D. A. (2014). Genome-wide signals of
398 positive selection in human evolution. *Genome research*, 24(6):885–895.
- 399 Fagny, M., Patin, E., Enard, D., Barreiro, L. B., Quintana-Murci, L., and Laval,
400 G. (2014). Exploring the occurrence of classic selective sweeps in humans
401 using whole-genome sequencing data sets. *Molecular biology and evolution*,
402 31(7):1850–1868.
- 403 Ferrer-Admetlla, A., Liang, M., Korneliussen, T., and Nielsen, R. (2014). On de-
404 tecting incomplete soft or hard selective sweeps using haplotype structure.
405 *Molecular biology and evolution*, 31(5):1275–1291.
- 406 Field, Y., Boyle, E. A., Telis, N., Gao, Z., Gaulton, K. J., Golan, D., Yengo, L.,
407 Rocheleau, G., Froguel, P., McCarthy, M. I., and Pritchard, J. K. (2016). Detec-
408 tion of human adaptation during the past 2000 years. *Science*, 354(6313):760–
409 764.
- 410 Forni, D., Cagliani, R., Tresoldi, C., Pozzoli, U., De Gioia, L., Filippi, G., Riva,
411 S., Menozzi, G., Colleoni, M., Biasin, M., Lo Caputo, S., Mazzotta, F., Comi,
412 G. P., Bresolin, N., Clerici, M., and Sironi, M. (2014). An evolutionary analy-
413 sis of antigen processing and presentation across different timescales reveals
414 pervasive selection. *PLoS genetics*, 10(3):e1004189.
- 415 Gillespie, J. H. (2004). *Population Genetics: A Concise Guide*. JHU Press.

- 416 Haasl, R. J. and Payseur, B. A. (2016). Fifteen years of genomewide scans for
417 selection: trends, lessons and unaddressed genetic sources of complication.
418 *Molecular ecology*, 25(1):5–23.
- 419 Haller, B. C. and Messer, P. W. (2018). SLiM 3: Forward genetic simulations
420 beyond the Wright-Fisher model. *bioRxiv*, page 418657.
- 421 Harris, R. B., Sackman, A., and Jensen, J. D. (2018). On the unfounded enthu-
422 siasm for soft selective sweeps II: Examining recent evidence from humans,
423 flies, and viruses. *PLoS genetics*, 14(12):e1007859.
- 424 Hernandez, R. D., Kelley, J. L., Elyashiv, E., Melton, S. C., Auton, A., McVean,
425 G., Sella, G., and Przeworski, M. (2011). Classic selective sweeps were rare
426 in recent human evolution. *Science*, 331(6019):920–924.
- 427 Johnson, K. E. and Voight, B. F. (2018). Patterns of shared signatures of recent
428 positive selection across human populations. *Nature Ecology & Evolution*,
429 2(4):713–720.
- 430 Karolchik, D., Hinrichs, A. S., Furey, T. S., Roskin, K. M., Sugnet, C. W., Haus-
431 sler, D., and Kent, W. J. (2004). The UCSC table browser data retrieval tool.
432 *Nucleic acids research*, 32(Database issue):D493–6.
- 433 Lande, R. (1975). The maintenance of genetic variability by mutation in a poly-
434 genic character with linked loci. *Genetical research*, 26(3):221–235.
- 435 Lewontin, R. C. and Kojima, K.-I. (1960). THE EVOLUTIONARY DYNAMICS
436 OF COMPLEX POLYMORPHISMS. *Evolution; international journal of organic*
437 *evolution*, 14(4):458–472.
- 438 Liu, R. Y. and Singh, K. (1992). Moving blocks jackknife and bootstrap capture
439 weak dependence. *Exploring the limits of bootstrap*, 225:248.
- 440 Messer, P. W. (2013). SLiM: simulating evolution with selection and linkage.
441 *Genetics*, 194(4):1037–1039.
- 442 Nei, M. (1972). Genetic distance between populations. *The American naturalist*,
443 106(949):283–292.
- 444 O'Reilly, P. F., Birney, E., and Balding, D. J. (2008). Confounding between re-
445 combination and selection, and the Ped/Pop method for detecting selection.
446 *Genome research*, 18(8):1304–1313.
- 447 Pickrell, J. K., Coop, G., Novembre, J., Kudaravalli, S., Li, J. Z., Absher, D.,
448 Srinivasan, B. S., Barsh, G. S., Myers, R. M., Feldman, M. W., and Pritchard,
449 J. K. (2009). Signals of recent positive selection in a worldwide sample of
450 human populations. *Genome research*, 19(5):826–837.
- 451 Ponjavic, J., Ponting, C. P., and Lunter, G. (2007). Functionality or transcrip-
452 tional noise? evidence for selection within long noncoding RNAs. *Genome*
453 *research*, 17(5):556–565.

- 454 Pouyet, F., Aeschbacher, S., Thiéry, A., and Excoffier, L. (2018). Background
455 selection and biased gene conversion affect more than 95% of the human
456 genome and bias demographic inferences. *eLife*, 7.
- 457 Pritchard, J. K., Pickrell, J. K., and Coop, G. (2010). The genetics of human
458 adaptation: hard sweeps, soft sweeps, and polygenic adaptation. *Current*
459 *biology: CB*, 20(4):R208–15.
- 460 Sabeti, P. C., Reich, D. E., Higgins, J. M., Levine, H. Z. P., Richter, D. J., Schaffner,
461 S. F., Gabriel, S. B., Platko, J. V., Patterson, N. J., McDonald, G. J., Ackerman,
462 H. C., Campbell, S. J., Altshuler, D., Cooper, R., Kwiatkowski, D., Ward, R.,
463 and Lander, E. S. (2002). Detecting recent positive selection in the human
464 genome from haplotype structure. *Nature*, 419(6909):832–837.
- 465 Schrider, D. R. and Kern, A. D. (2016). S/HIC: Robust identification of soft and
466 hard sweeps using machine learning. *PLoS genetics*, 12(3):e1005928.
- 467 Schrider, D. R. and Kern, A. D. (2017). Soft sweeps are the dominant mode of
468 adaptation in the human genome. *Molecular biology and evolution*, 34(8):1863–
469 1877.
- 470 Slatkin, M. (2008). Linkage disequilibrium—understanding the evolutionary
471 past and mapping the medical future. *Nature reviews. Genetics*, 9(6):477–485.
- 472 Smith, J. M. and Haigh, J. (1974). The hitch-hiking effect of a favourable gene.
473 *Genetical research*, 23(1):23–35.
- 474 Szpiech, Z. A. and Hernandez, R. D. (2014). Selscan: An efficient multithreaded
475 program to perform EHH-based scans for positive selection. *Molecular biol-*
476 *ogy and evolution*, 31(10):2824–2827.
- 477 Vitti, J. J., Grossman, S. R., and Sabeti, P. C. (2013). Detecting natural selection
478 in genomic data. *Annual review of genetics*, 47:97–120.
- 479 Voight, B. F., Kudaravalli, S., Wen, X., and Pritchard, J. K. (2006). A map of
480 recent positive selection in the human genome. *PLoS biology*, 4(3):0446–0458.
- 481 Wall, J. D. and Pritchard, J. K. (2003). Haplotype blocks and linkage disequilib-
482 rium in the human genome. *Nature reviews. Genetics*, 4(8):587–597.