

1 GENOMIC REGIONS LINKED TO SOFT  
2 SWEEPS APPROXIMATE NEUTRALITY  
3 WHEN INFERRING POPULATION  
4 HISTORY FROM SITE PATTERN  
5 FREQUENCIES

6 Nathan S. Harris and Alan R. Rogers

7 April 21, 2020

8 **ABSTRACT**

9 Recent studies have suggested that selection is widespread throughout the  
10 genome and largely uncompensated for in inferences of population history.  
11 To address this potential issue, we estimated site pattern frequencies for neu-  
12 tral and selection associated areas of the genome. There are notable differences  
13 in these frequencies between neutral regions and those affected by selection.  
14 However, these differences have relatively small effects when inferring popu-  
15 lation history.

16 **1 INTRODUCTION AND BACKGROUND**

17 In the past year, population geneticists have been debating the extent to which  
18 natural selection has shaped the human genome. Evidence suggests that soft  
19 sweeps (Harris et al., 2018; Schrider and Kern, 2017) and polygenic adaptation  
20 (Daub et al., 2013; Hernandez et al., 2011; Pritchard et al., 2010) are the pri-  
21 mary modes of selection in humans. This has led some researchers to suggest  
22 that most of the genome is in some way affected by selection either directly  
23 or indirectly through linkage with neighboring sites. This led Kern and Hahn  
24 (2018) to argue that the original lines of evidence that led to the neutral the-  
25 ory of evolution pioneered by Kimura (1983) do not hold up in the genomic  
26 era. Specifically, they believe inferred population histories are skewed because  
27 areas of the genome often perceived to be neutral are actually affected by selec-  
28 tion. In addition, they suggest genome-wide selection scans will have a high  
29 rate of false negatives if using a null distribution built from selected regions  
30 of the genome. On the other hand, Harris et al. (2018) and Jensen et al. (2018)  
31 suggest that many of these apparent signals of selection are false positives, and

32 others have suggested that selection has been less common but signals of se-  
33 lection are amplified by population history (Torres et al., 2018).

34 A simple way to test the idea proposed by Kern and Hahn (2018) is to  
35 reconstruct population history using different areas of the genomes. Schrider  
36 and Kern (2017) used a machine learning algorithm to assemble a list of re-  
37 gions in the genome inferred to be affected by selection but previously thought  
38 to be neutral, and those that are neutral. Here, relative site pattern frequencies  
39 (Rogers, 2019) are calculated in each of these subdivisions to measure the po-  
40 tential effects of selection. In what follows, we show that selection skews site  
41 pattern frequencies, but has little effect on the estimation of population history  
42 in our model.

## 43 2 RESULTS

44 **Selection affects population history.** A nucleotide site pattern is a particular  
45 arrangement of derived and ancestral alleles when a single haploid individ-  
46 ual is sampled from each sample population. The total number of distinct site  
47 patterns is therefore all combinations in which at least one sample, but not all,  
48 carries the derived allele. Three populations are discussed here, given the la-  
49 bels CEU ( $X$ ), JPT ( $Y$ ), and YRI ( $Z$ ). The possible site patterns are  $x$ ,  $y$ ,  $z$ ,  $xy$ ,  
50  $xz$ , and  $yz$ , the first three representing the cases in which the derived allele is  
51 found in only one population, and the rest representing when the derived al-  
52 lele is found in two populations. Estimations of common ancestry, divergence  
53 times, and admixture can be made using the relative frequency of of these site  
54 patterns (Rogers, 2019). The difference between the selection-affected site pat-  
55 terns and neutral site patterns was calculated. If neutral and selection-affected  
56 genomic regions have similar site pattern frequencies to neutral regions, it will  
57 produce similar estimates of population history, and the difference between  
58 them should not vary from zero significantly.

59 Figure 1 shows the difference in relative frequencies of site patterns be-  
60 tween affected and neutral regions for biallelic single nucleotide polymorphisms  
61 (SNP). Site pattern frequencies were calculated using the *sitapat* program from  
62 the *Legofit* package (Rogers, 2019). Confidence intervals were generated by  
63 using 1,000 bootstrap replicates generated from *sitapat*. Neutral and selection-  
64 affected regions differed significantly in the  $x$ ,  $y$ , and  $xy$  site patterns, with  
65 confidence intervals that do not overlap with zero. To investigate the possibil-  
66 ity that this pattern is driven primarily by one type of selection, hard sweeps or  
67 soft sweeps, affected regions were split accordingly. Hard sweeps differ from  
68 other distinctions substantially but the confidence intervals are large, likely  
69 due to the relatively small sample size of these regions (Table 1). Soft sweep  
70 regions are marginally more similar to neutral regions, but the large difference  
71 between affected and neutral regions in the  $x$ ,  $y$ , and  $xy$  patterns persists when  
72 only soft sweeps are considered.

73 **Inference of Population History.** Site pattern frequencies differ between  
74 neutral and selected affected regions of the genome. However, these differ-

75 ences are small and it is unclear how large their effect will be when inferring  
 76 population history. To test the effect these differences have, the site patterns  
 77 for neutral and soft-sweep affected regions were used to estimate ancestral di-  
 78 vergence times and population sizes using *Legofit* (Rogers, 2019). The demo-  
 79 graphic model used was taken from (Rogers et al., 2019). This model includes  
 80 admixture events from Neanderthals into Eurasia, ancient humans into Nean-  
 81 derthals, and from superarchaic hominins into Denisovans and the ancestor of  
 82 Neanderthals and Denisovans. Rogers et al. (2019) found that the exclusion of  
 83 these admixture events can strongly bias results.

84 Selection seems to affect the estimation of these parameters, but only to a

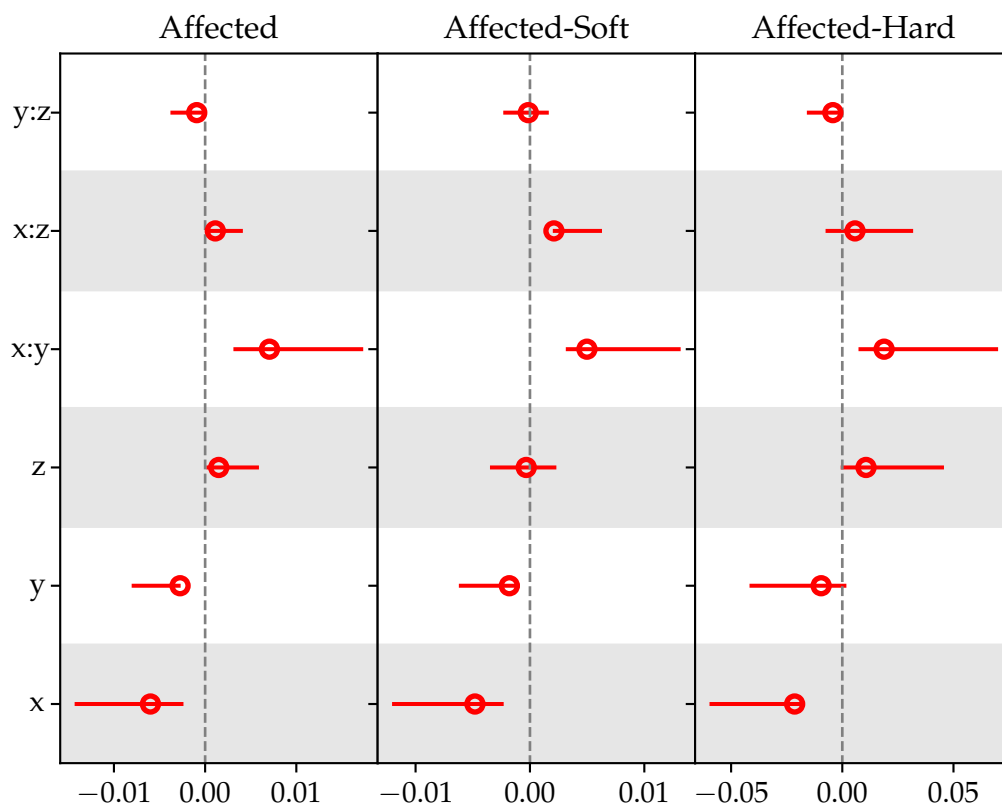


Figure 1: Difference in selection-affected and neutral relative site pattern frequencies. If all patterns rested at zero, the inferred population history would be identical. Each population is given an alphabetical label:  $x$  (CEU),  $y$  (JPT),  $z$  (YRI)

Samples	Neutral	Selection Affected	Soft-Sweep Affected	Hard-Sweep Affected
Humans only	285,833	619,136	431,105	4,460
Human and Neanderthal	251,953	543,098	378,700	4,104
Human, Neanderthal, and Denisova	261,996	563,523	392,625	4,275

Table 1: Number of sites tabulated in each run of *sitepat*

85 small extent. Figure 2 shows the percent differences between soft-affected and  
 86 neutral estimates of divergence times. Confidence intervals were generated by  
 87 taking differences of individual estimates in each of fifty bootstrap replicates.  
 88 If selection does not strongly affect quantitative estimates of demographic in-  
 89 ference, these differences should rest around zero.

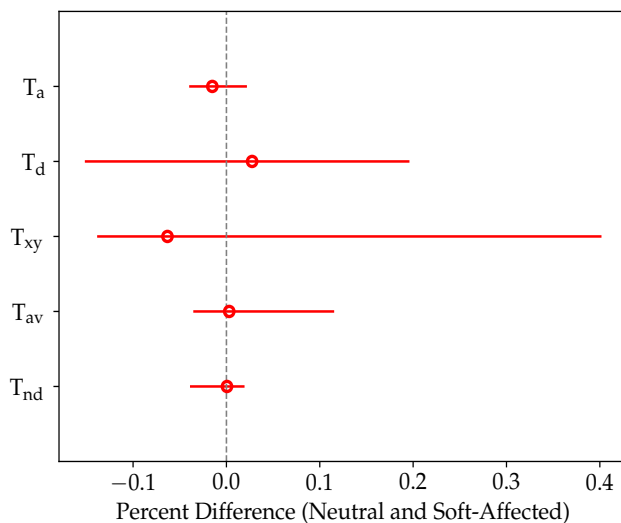


Figure 2: Percent difference between estimates of divergence times for soft-sweep associated regions and neutral regions for a model including Yorubans ( $x$ ), European (CEU) ( $y$ ), Neanderthals ( $n$ ) further split into the Vindija ( $v$ ) and Altai ( $a$ ), Denisovans ( $d$ ), and superarchaic hominins ( $s$ ).

90 Figure 3 shows the percent differences between soft-affected and neutral es-  
 91 timates of ancestral population sizes. In general, estimates of population size  
 92 were more likely to be disrupted by the effects of soft selective sweeps. How-

93 ever, these deviations are relatively small. The largest difference in population  
94 size without confidence intervals overlapping with zero is approximately one  
95 or two percent of approximately 40,000 (Table 2).

Table 2: Estimates of divergence times and population sizes for a model including Yorubans ( $x$ ), Europeans (CEU) ( $y$ ), Neanderthals ( $n$ ), Denisovans ( $d$ ), and superarchaic hominins ( $s$ ). “T” indicates divergence time and “2N” is the diploid effective population size.

Parameter	Neutral	Soft Selection Affected
$T_{xynds}$	67,049.1	52,396.2
$T_{nd}$	24,906.3	24,921.1
$T_{av}$	15,570.6	15,617.5
$T_{xy}$	3,864.6	3,628.0
$T_d$	1,855.0	1,906.8
$T_a$	4,958.1	4,883.7
$T_v$	2,339.7	2,315.4
$2N_{av}$	31,673.6	24,333.9
$2N_n$	7,372.7	7,414.5
$2N_{nd}$	2,111.9	2,063.7
$2N_{xy}$	50,259.1	50,639.3
$2N_{xynd}$	42,152.5	42,826.4
$2N_s$	81,794.2	50,639.7

### 96 3 DISCUSSION

97 The difference in site pattern frequencies is in agreement with the results of  
98 Schrider and Kern (2017). The site patterns for the neutral regions imply a  
99 different population history than the selection affected regions. The neutral  
100 and selection-affected regions are close for each site pattern except the  $x$  and  
101  $xy$  patterns. The  $x$  and  $y$  site patterns are relatively over-represented in the  
102 neutral case, while the  $xy$  site pattern is relatively underrepresented. The se-  
103 lection affected site patterns would therefore overestimate the length of the  
104 ancestral Eurasian branch, and underestimate the individual Eurasian popula-  
105 tion branches.

106 The significant differences in site pattern frequencies do not translate to  
107 large differences in estimations of population history in our model. In the Eu-  
108 ropean model, only one estimated parameter, the size of the ancestral popula-  
109 tion of humans, Neanderthals, and Denisovans ( $2N_{xynd}$ ), shows a significant  
110 difference between neutral and soft sweep affected regions. However, in this  
111 case the difference is more than an order of magnitude smaller than the esti-  
112 mated parameter itself.

113 The results here are consistent with large portions of the genome being un-

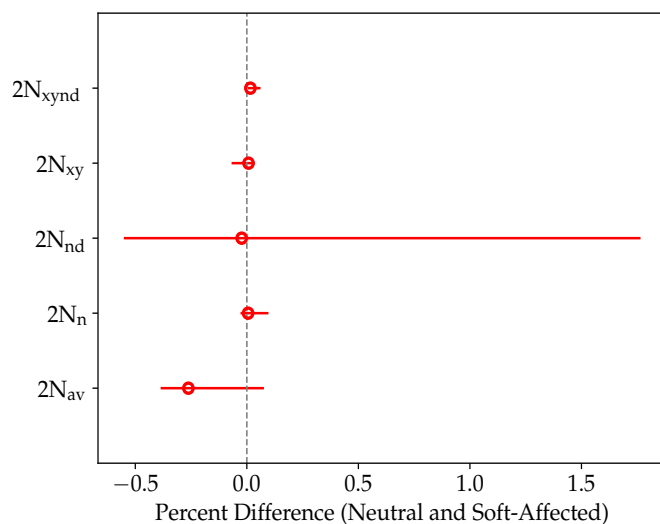


Figure 3: Percent difference between estimates of ancestral effective population size for soft-sweep associated regions and neutral regions for a model including Yorubans ( $x$ ), European (CEU) ( $y$ ), Neanderthals ( $n$ ) further split into the Vindija ( $v$ ) and Altai ( $a$ ), Denisovans ( $d$ ), and superarchaic hominins ( $s$ ).

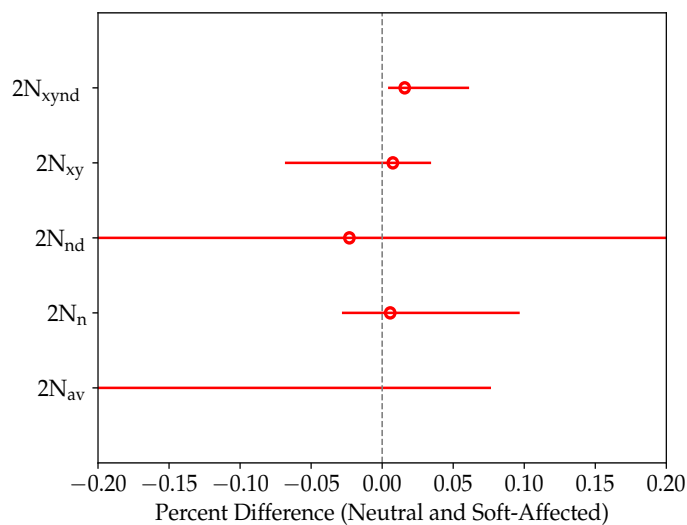


Figure 4: Percent difference shown in 3 zoomed in to compensate for the visual bias caused by the large amount of error around estimates of  $2N_{xynd}$ .

114 der selection. However, differences in site pattern frequencies generated by  
115 soft selection do not appear to have substantial effects on the estimation of  
116 parameters of population history, at least when using *Legofit*. This result may  
117 have a larger effect at finer scales. For instance, the significant difference in  $T_{av}$ ,  
118 the divergence time of the Altai and Vindija Neanderthals, could be as large as  
119 5,000 generations or 125,000 years. This may not too large when considering  
120 recent human population history and admixture between humans and other  
121 hominins, but could be problematic in reconstructing a fine scale history of  
122 Neanderthal populations. Further work exploring different time scales and  
123 population histories should be done before any generalization of these results  
124 is made. Nonetheless, soft sweeps do not cause meaningful disturbances in  
125 estimates of population history at the scale studied here.

## 126 4 METHODS

127 **Simons Data.** Simons Genome Diversity Mallick et al. (2016) data for Japan  
128 (JPT), Yoruba (YRI), and Europe (England and France) was acquired from <https://www.ebi.ac.uk/ena/data/view/PRJEB9586>. Sites with a map quality or geno-  
129 type quality below thirty were excluded. Sites that were fixed across all human  
130 populations were removed, because they cannot differentiate human popula-  
131 tions. Indels, SNPs within seven bases of an indel, and low quality SNPs (filter  
132 level equal to zero) were removed. The Central Europeans from Utah (CEU)  
133 sample is not represented in the SGDP. Instead, individuals from England and  
134 France were used as a proxy for CEU. These samples were used because CEU  
135 is present in the analysis and results of Schrider and Kern (2017), but absent  
136 from the Simons data.  
137

138 **Subdivisions.** Schrider and Kern (2017) subdivide the genome into neu-  
139 tral and selection linked, soft-sweep affected, and hard-sweep affected regions.  
140 These results were obtained from [https://github.com/kern-lab/shIC/tree/](https://github.com/kern-lab/shIC/tree/master/humanScanResults)  
141 [master/humanScanResults](https://github.com/kern-lab/shIC/tree/master/humanScanResults). The selection linked regions are those that are  
142 commonly assumed to be neutral in the literature, but showed evidence of  
143 being linked to and affected by selection in their machine learning analysis.  
144 The selection-affected and neutral regions were used here to divide data re-  
145 spectively. Regions inferred to be under selection are not studied here because  
146 selected regions are already expected to produce different estimates of popu-  
147 lation history from neutral regions. Here we are concerned with regions of the  
148 genome that could be mistaken for neutral regions and skew population his-  
149 tories in the literature. Each sample has its own subdivision that reflects the  
150 population and selective history of the population it belongs to. The neutral  
151 distinction made here refers to sites where all three populations are considered  
152 neutral. For a site to be included in the “selection-affected” for the purpose  
153 of this analysis, at least one population needs to be represented in the original  
154 “selection-affected” distinction.

155 **Site Pattern Frequencies.** The program *sitapat* from the *Legofit* (Rogers,  
156 2019) package was used to calculate relative site pattern frequencies (SPF) (Rogers

157 et al., 2017). The ancestral allele is determined by using reference alleles of  
158 chimpanzees and gorillas. The analysis was limited to sites where the chimpanzee and gorilla were fixed for the same allele, and the human samples were  
159 polymorphic. One thousand bootstrap replicates were generated for each combination of selection type and set of populations. The difference in site pattern  
160 frequencies between selection affected regions and neutral regions was then taken, with confidence intervals generated from differences in individual bootstrap  
161 replicates.  
162  
163

164  
165 **Legofit.** Site pattern frequencies were used to estimate population history parameters using *Legofit*. The model of population history includes admixture  
166 between Eurasians and Neandethals, and between superarchaics and Denisovans and the ancestor of Denisovans and Neanderthals. Site pattern frequencies  
167 are generated using high coverage Neanderthal genomes from the Vindija (Prüfer et al., 2017) and Altai (Prüfer et al., 2014) Neanderthals, and a high  
168 coverage Denisovan genome from Siberia (Meyer et al., 2012). The Yoruban samples from the SGDP represent Africa. The final human population was a  
169 mixture of English and French individuals serving as a proxy for CEU. One thousand bootstrap replicates were generated for each model of population  
170 history. Confidence intervals for these differences were generated by taking the inner ninety-five percent of the differences between individual bootstrap  
171 replicate estimations. This process was conducted for neutral and soft-selection affected regions separately.  
172  
173  
174  
175  
176  
177  
178

## 179 References

- 180 Daub, J. T., Hofer, T., Cutivet, E., Dupanloup, I., Quintana-Murci, L., Robinson-Rechavi, M., and Excoffier, L. (2013). Evidence for polygenic adaptation to pathogens in the human genome. *Molecular biology and evolution*, 30(7):1544–1558.  
181  
182  
183
- 184 Harris, R. B., Sackman, A., and Jensen, J. D. (2018). On the unfounded enthusiasm for soft selective sweeps II: Examining recent evidence from humans, flies, and viruses. *PLoS genetics*, 14(12):e1007859.  
185  
186
- 187 Hernandez, R. D., Kelley, J. L., Elyashiv, E., Melton, S. C., Auton, A., McVean, G., Sella, G., and Przeworski, M. (2011). Classic selective sweeps were rare in recent human evolution. *Science*, 331(6019):920–924.  
188  
189
- 190 Jensen, J. D., Payseur, B. A., Stephan, W., Aquadro, C. F., Lynch, M., Charlesworth, D., and Charlesworth, B. (2018). The importance of the neutral theory in 1968 and 50 years on: a response to kern & hahn 2018. *Evolution; international journal of organic evolution*.  
191  
192  
193
- 194 Kern, A. D. and Hahn, M. W. (2018). The neutral theory in light of natural selection. *Molecular biology and evolution*, 35(6):1366–1371.  
195



- 196 Kimura, M. (1983). *The Neutral Theory of Molecular Evolution*. Cambridge Uni-  
197 versity Press.
- 198 Mallick, S., Li, H., Lipson, M., Mathieson, I., Gymrek, M., Racimo, F., Zhao,  
199 M., Chennagiri, N., Nordenfelt, S., Tandon, A., Skoglund, P., Lazaridis, I.,  
200 Sankararaman, S., Fu, Q., Rohland, N., Renaud, G., Erlich, Y., Willems, T.,  
201 Gallo, C., Spence, J. P., Song, Y. S., Poletti, G., Balloux, F., van Driem, G.,  
202 de Knijff, P., Romero, I. G., Jha, A. R., Behar, D. M., Bravi, C. M., Capelli, C.,  
203 Hervig, T., Moreno-Estrada, A., Posukh, O. L., Balanovska, E., Balanovsky,  
204 O., Karachanak-Yankova, S., Sahakyan, H., Toncheva, D., Yepiskoposyan,  
205 L., Tyler-Smith, C., Xue, Y., Abdullah, M. S., Ruiz-Linares, A., Beall, C. M.,  
206 Di Rienzo, A., Jeong, C., Starikovskaya, E. B., Metspalu, E., Parik, J., Villems,  
207 R., Henn, B. M., Hodoglugil, U., Mahley, R., Sajantila, A., Stamatoyannopou-  
208 los, G., Wee, J. T. S., Khusainova, R., Khusnutdinova, E., Litvinov, S., Ayodo,  
209 G., Comas, D., Hammer, M. F., Kivisild, T., Klitz, W., Winkler, C. A., Labuda,  
210 D., Bamshad, M., Jorde, L. B., Tishkoff, S. A., Watkins, W. S., Metspalu, M.,  
211 Dryomov, S., Sukernik, R., Singh, L., Thangaraj, K., Pääbo, S., Kelso, J., Pat-  
212 terson, N., and Reich, D. (2016). The simons genome diversity project: 300  
213 genomes from 142 diverse populations. *Nature*, 538(7624):201–206.
- 214 Meyer, M., Kircher, M., Gansauge, M.-T., Li, H., Racimo, F., Mallick, S.,  
215 Schraiber, J. G., Jay, F., Prüfer, K., de Filippo, C., Sudmant, P. H., Alkan, C.,  
216 Fu, Q., Do, R., Rohland, N., Tandon, A., Siebauer, M., Green, R. E., Bryc,  
217 K., Briggs, A. W., Stenzel, U., Dabney, J., Shendure, J., Kitzman, J., Ham-  
218 mer, M. F., Shunkov, M. V., Derevianko, A. P., Patterson, N., Andrés, A. M.,  
219 Eichler, E. E., Slatkin, M., Reich, D., Kelso, J., and Pääbo, S. (2012). A high-  
220 coverage genome sequence from an archaic denisovan individual. *Science*,  
221 338(6104):222–226.
- 222 Pritchard, J. K., Pickrell, J. K., and Coop, G. (2010). The genetics of human  
223 adaptation: hard sweeps, soft sweeps, and polygenic adaptation. *Current*  
224 *biology: CB*, 20(4):R208–15.
- 225 Prüfer, K., de Filippo, C., Grote, S., Mafessoni, F., Korlević, P., Hajdinjak, M.,  
226 Vernot, B., Skov, L., Hsieh, P., Peyrégne, S., Reher, D., Hopfe, C., Nagel, S.,  
227 Maricic, T., Fu, Q., Theunert, C., Rogers, R., Skoglund, P., Chintalapati, M.,  
228 Dannemann, M., Nelson, B. J., Key, F. M., Rudan, P., Kućan, Ž., Gušić, I.,  
229 Golovanova, L. V., Doronichev, V. B., Patterson, N., Reich, D., Eichler, E. E.,  
230 Slatkin, M., Schierup, M. H., Andrés, A. M., Kelso, J., Meyer, M., and Pääbo,  
231 S. (2017). A high-coverage neandertal genome from vindija cave in croatia.  
232 *Science*, 358(6363):655–658.
- 233 Prüfer, K., Racimo, F., Patterson, N., Jay, F., Sankararaman, S., Sawyer, S.,  
234 Heinze, A., Renaud, G., Sudmant, P. H., de Filippo, C., Li, H., Mallick, S.,  
235 Dannemann, M., Fu, Q., Kircher, M., Kuhlwillm, M., Lachmann, M., Meyer,  
236 M., Ongyerth, M., Siebauer, M., Theunert, C., Tandon, A., Moorjani, P., Pick-  
237 rell, J., Mullikin, J. C., Vohr, S. H., Green, R. E., Hellmann, I., Johnson, P. L. F.,  
238 Blanche, H., Cann, H., Kitzman, J. O., Shendure, J., Eichler, E. E., Lein, E. S.,

- 239 Bakken, T. E., Golovanova, L. V., Doronichev, V. B., Shunkov, M. V., Dere-  
240 vianko, A. P., Viola, B., Slatkin, M., Reich, D., Kelso, J., and Pääbo, S. (2014).  
241 The complete genome sequence of a neanderthal from the altai mountains.  
242 *Nature*, 505(7481):43–49.
- 243 Rogers, A. R. (2019). Legofit: Estimating population history from genetic data.
- 244 Rogers, A. R., Bohlender, R. J., and Huff, C. D. (2017). Early history of ne-  
245 anderthals and denisovans. *Proceedings of the National Academy of Sciences*,  
246 114(37):201706426.
- 247 Rogers, A. R., Harris, N. S., and Achenbach, A. A. (2019). Neanderthal-  
248 Denisovan ancestors interbred with a distantly-related hominin.
- 249 Schrider, D. R. and Kern, A. D. (2017). Soft sweeps are the dominant mode of  
250 adaptation in the human genome. *Molecular biology and evolution*, 34(8):1863–  
251 1877.
- 252 Torres, R., Szpiech, Z. A., and Hernandez, R. D. (2018). Human demographic  
253 history has amplified the effects of background selection across the genome.  
254 *PLoS genetics*, 14(6):e1007387.