

1 *Type: Article*

2

## 3 **selscan 2.0: scanning for sweeps in unphased data**

4

5 Zachary A. Szpiech<sup>1,2,\*</sup>

6

7 <sup>1</sup> *Department of Biology, Pennsylvania State University, University Park, PA 16802, USA*

8 <sup>2</sup> *Institute for Computational and Data Sciences, Pennsylvania State University, University Park,*

9 *PA 16802, USA*

10 *\* Correspondence: [szpiech@psu.edu](mailto:szpiech@psu.edu)*

11

### 12 **Abstract**

13 Haplotype-based scans to identify recent and ongoing positive selection have become  
14 commonplace in evolutionary genomics studies of numerous species across the tree of life.  
15 However, the most widely adopted approaches require phased haplotypes to compute the key  
16 statistics. Here we release a major update to the selscan software that re-defines popular  
17 haplotype-based statistics for use with unphased “multi-locus genotype” data. We provide  
18 unphased implementations of iHS, nSL, XP-EHH, and XP-nSL and evaluate their performance  
19 across a range of important parameters in a generic demographic history. We also show that  
20 these implementations often outperform a naïve application of the original statistics to unphased  
21 data, and that they perform with minimal to no reduction in power compared to the original  
22 statistics when phase is perfectly known. Source code and executables are available at  
23 <https://www.github.com/szpiech/selscan>.

### 24 **1 Introduction**

25 Haplotype-based summary statistics—such as iHS (Voight, et al. 2006), nSL (Ferrer-  
26 Admetlla, et al. 2014), XP-EHH (Sabeti, et al. 2007), and XP-nSL (Szpiech, et al. 2021)—have  
27 become commonplace in evolutionary genomics studies to identify recent and ongoing positive  
28 selection in populations (e.g., Colonna, et al. 2014; Zoledziwska, et al. 2015; Nedelec, et al.  
29 2016; Crawford, et al. 2017; Meier, et al. 2018; Lu, et al. 2019; Zhang, et al. 2020; Salmon, et al.  
30 2021). When an adaptive allele sweeps through a population, it leaves a characteristic pattern

31 of long high-frequency haplotypes and low genetic diversity in the vicinity of the allele. These  
32 statistics aim to capture these signals by summarizing the decay of haplotype homozygosity as  
33 a function of distance from a putatively selected region, either within a single population (iHS  
34 and nSL) or between two populations (XP-EHH and XP-nSL).

35         These haplotype-based statistics are powerful for detecting recent positive selection  
36 (Colonna, et al. 2014; Zoledziewska, et al. 2015; Nedelec, et al. 2016; Crawford, et al. 2017;  
37 Meier, et al. 2018; Lu, et al. 2019; Zhang, et al. 2020; Salmon, et al. 2021), and the two-  
38 population versions can even out-perform pairwise  $F_{st}$  scans on a large swath of the parameter  
39 space (Szpiech, et al. 2021). Furthermore, haplotype-based methods have also been shown to  
40 be robust to background selection (Fagny, et al. 2014; Schrider 2020). However, each of these  
41 statistics presumes that haplotype phase is known or well-estimated.

42         As the generation of genomic sequencing data for non-model organisms is becoming  
43 routine (Ellegren 2014), there are many great opportunities for studying recent adaptation  
44 across the tree of life (e.g., Campagna and Toews (2022)). However, often these  
45 organisms/populations do not have a well-characterized demographic history or recombination  
46 rate map, two pieces of information which are important inputs for statistical phasing methods  
47 (Delaneau, et al. 2013; Browning, et al. 2021).

48         Recent work has shown that converting haplotype data into “multi-locus genotype” data  
49 is an effective approach for using haplotype-based selection statistics such as G12, LASSI, and  
50 saltiLASSI (Harris, et al. 2018; Harris and DeGiorgio 2020; DeGiorgio and Szpiech 2022) in  
51 unphased data. Recognizing this, we have reformulated the iHS, nSL, XP-EHH, and XP-nSL  
52 statistics to use multi-locus genotypes and provided an easy-to-use implementation in selscan  
53 2.0 (Szpiech and Hernandez 2014). We evaluate the performance of these unphased statistics  
54 under various generic demographic models and compare against the original statistics applied  
55 to simulated datasets when phase is either known or unknown.

## 56 2 New Approaches

57 When the --unphased flag is set in selscan v2.0+, biallelic genotype data is collapsed  
58 into multi-locus genotype data by representing the genotype as either 0, 1, or 2—the number of  
59 derived alleles observed. In this case, selscan v2.0+ will then compute iHS, nSL, XP-EHH, and  
60 XP-nSL as described below. We follow the notation conventions of Szpiech and Hernandez  
61 (2014).

### 62 2.1 Extended Haplotype Homozygosity

63 In a sample of  $n$  diploid individuals, let  $\mathcal{C}$  denote the set of all possible genotypes at  
64 locus  $x_0$ . For multi-locus genotypes,  $\mathcal{C} := \{0,1,2\}$ , representing the total counts of a derived  
65 allele. Let  $\mathcal{C}(x_i)$  be the set of all unique haplotypes extending from site  $x_0$  to site  $x_i$  either  
66 upstream or downstream of  $x_0$ . If  $x_1$  is a site immediately adjacent to  $x_0$ , then  $\mathcal{C}(x_1) :=$   
67  $\{00,01,02,10,11,12,20,21,22\}$ , representing all possible two-site multi-locus genotypes. We can  
68 then compute the extended haplotype homozygosity (EHH) of a set of multi-locus genotypes as

$$69 \text{EHH}(x_i) = \sum_{h \in \mathcal{C}(x_i)} \frac{\binom{n_h}{2}}{\binom{n}{2}},$$

70 where  $n_h$  is the number of observed haplotypes of type  $h$ .

71 If we wish to compute the EHH of a subset of observed haplotypes that all contain the  
72 same ‘core’ multi-locus genotype, let  $\mathcal{H}_c(x_i)$  be the partition of  $\mathcal{C}(x_i)$  containing genotype  $c \in \mathcal{C}$   
73 at  $x_0$ . For example, choosing a homozygous derived genotype ( $c = 2$ ) as the core,  $\mathcal{H}_2 :=$   
74  $\{20,21,22\}$ . Thus, we can compute the EHH of all individuals carrying a given genotype at site  $x_0$   
75 extending out to site  $x_i$  as

$$76 \text{EHH}_c(x_i) = \sum_{h \in \mathcal{H}_c(x_i)} \frac{\binom{n_h}{2}}{\binom{n_c}{2}},$$

77 where  $n_h$  is the number of observed haplotypes of type  $h$  and  $n_c$  is the number of observed  
 78 multi-locus genotypes with core genotype of  $c$ . Finally, we can compute the complement EHH of  
 79 a sample of multi-locus genotypes as

$$80 \quad cEHH_c(x_i) = \sum_{h \in \mathcal{C}(x_i) \setminus \mathcal{H}_c(x_i)} \frac{\binom{n_h}{2}}{\binom{n_{c'}}{2}},$$

81 where  $n_{c'}$  is the number of observed multi-locus genotypes with a core genotype of not  $c$ .

## 82 **2.2 iHS and nSL**

83 Unphased iHS and nSL are calculated using the equations above. First, we compute the  
 84 integrated haplotype homozygosity (iHH) for the homozygous ancestral ( $c = 0$ ) and derived ( $c =$   
 85 2) core genotypes as

$$86 \quad iHH_c = \sum_{i=1}^{|\mathcal{D}|} \frac{1}{2} (EHH_c(x_{i-1}) + EHH_c(x_i)) d(x_{i-1}, x_i) + \sum_{i=1}^{|\mathcal{U}|} \frac{1}{2} (EHH_c(x_{i-1}) + EHH_c(x_i)) d(x_{i-1}, x_i),$$

87 where  $\mathcal{D}$  is the set of downstream sites from the core locus and  $\mathcal{U}$  is the set of upstream sites.  
 88  $d(x_{i-1}, x_i)$  is a measure of genomic distance between to markers and is the genetic distance in  
 89 centimorgans or physical distance in basepairs for iHS (Voight, et al. 2006) or the number of  
 90 sites observed for nSL (Ferrer-Admetlla, et al. 2014). We similarly compute the complement  
 91 integrated haplotype homozygosity (ciHH) for both homozygous core genotypes as

$$92 \quad ciHH_c = \sum_{i=1}^{|\mathcal{D}|} \frac{1}{2} (cEHH_c(x_{i-1}) + cEHH_c(x_i)) d(x_{i-1}, x_i)$$

$$93 \quad + \sum_{i=1}^{|\mathcal{U}|} \frac{1}{2} (cEHH_c(x_{i-1}) + cEHH_c(x_i)) d(x_{i-1}, x_i).$$

94 The (unstandardized) unphased iHS is then calculated as

$$95 \quad iHS = \begin{cases} iHS_2, & \text{if } iHS_2 > iHS_0 \\ -iHS_0, & \text{otherwise} \end{cases}$$

96 where  $iHS_2 = \log_{10} \left( \frac{iHH_2}{ciHH_2} \right)$  and  $iHS_0 = \log_{10} \left( \frac{iHH_0}{ciHH_0} \right)$ . Conceptually, this is nearly identical to the  
97 phased version of iHS, where the log ratio of the integrated haplotype homozygosity is  
98 computed between all haplotypes carrying the ancestral allele at the core locus versus all  
99 haplotypes carrying the derived allele at the core locus. In this case, however, we compare the  
100 iHH of the haplotypes containing homozygous genotypes of one allele at the core locus to the  
101 iHH of the haplotypes containing all other genotypes at the core locus. Doing this for both  
102 homozygous derived and homozygous ancestral haplotypes separately, we then choose the  
103 most extreme value. We assign a positive sign for long low-diversity haplotypes containing the  
104 derived homozygous genotype at the core locus, and we assign a negative sign for long low-  
105 diversity haplotypes containing the ancestral homozygous genotype at the core locus.  
106 Unstandardized iHS scores are then normalized in frequency bins, as previously described  
107 (Voight, et al. 2006; Ferrer-Admetlla, et al. 2014). Unstandardized unphased nSL is computed  
108 similarly with the appropriate distance measure (see Ferrer-Admetlla, et al. (2014) where they  
109 show that nSL can be reformulated as iHS with a different distance measure). Large positive  
110 scores indicate long high-frequency haplotypes with a homozygous derived core genotype, and  
111 large negative scores indicate long high-frequency haplotypes with a homozygous ancestral  
112 core genotype. Clusters of extreme scores in both directions indicate evidence for a sweep.

### 113 **2.3 XP-EHH and XP-nSL**

114 Unphased XP-EHH and XP-nSL are calculated by comparing the iHH between  
115 populations  $A$  and  $B$ , using the entire sample in each population. iHH in a population  $P$  is  
116 computed as

$$117 \quad iHH_P = \sum_{i=1}^{|\mathcal{D}|} \frac{1}{2} (EHH(x_{i-1}) + EHH(x_i)) d(x_{i-1}, x_i) + \sum_{i=1}^{|\mathcal{U}|} \frac{1}{2} (EHH(x_{i-1}) + EHH(x_i)) d(x_{i-1}, x_i),$$

118 where the distance measure is given as centimorgans or basepairs for XP-EHH (Sabeti, et al.  
119 2007) and number of sites observed for XP-nSL (Szpiech, et al. 2021). The XP statistics

120 between population  $A$  and  $B$  are then computed as  $XP = \log_{10} \left( \frac{iHH_A}{iHH_B} \right)$  and are normalized  
121 genome wide. Large positive scores indicate long high-frequency haplotypes in population  $A$ ,  
122 and large negative scores indicate long high-frequency haplotypes in population  $B$ . Clusters of  
123 extreme scores in one direction indicate evidence for a sweep in that population.

### 124 **3 Results**

125 We find that the unphased versions of iHS and nSL generally have good power (Figures  
126 1, S1, and S2) to detect selection prior to fixation of the allele, with nSL generally outperforming  
127 iHS. In smaller populations (Figure 1C and 1D), power does suffer relative to larger populations  
128 (Figure 1A, 1B, 1E, 1F). We note that these statistics struggle to identify soft sweeps when the  
129 population is undergoing exponential growth (Figure 1E and 1F). Each of these statistics also  
130 have low false positive rates hovering around 1% (Table S1).

131 Similarly, we find that the unphased versions of XP-EHH and XP-nSL have good power  
132 as well (Figures 2, 3, and S3-S6). When the sweep takes place in the smaller of the two  
133 populations (Figure 2C and 2D), we see a similar decrease in power, likely related to the lower  
134 efficiency of selection in small populations. When one population is undergoing exponential  
135 growth (Figure 3) performance is generally quite good, likely the result of a larger effective  
136 selection coefficient in large populations. These two-population statistics generally outperform  
137 their single-population counterparts, especially for sweeps that have reached fixation recently.  
138 Each of these statistics also have low false positive rates hovering around 1% (Table S1).

139 Next, we turn to comparing the performance of these unphased statistics to their phased  
140 counterparts when they are used to analyze either phased data or unphased data. In Figures 4-  
141 6 and S7-S12, we plot the difference in power between the unphased statistics and the phased  
142 counterpart applied to data with phase known (red lines) or phase scrambled (blue lines).  
143 Where these lines are greater than or equal to 0 indicates that the unphased statistic performed  
144 as well as or better than the phased counterpart.

145 We find that iHS tends to underperform the traditional phased implementations, but nSL  
146 tends to perform as well as the phased versions (Figures 4, S7, and S8). Although we note  
147 noticeable drops in unphased nSL power for softer sweeps in exponential growth scenarios  
148 (Figures 4F, S7F, and S8F) and for sweeps near completion in small population sizes (Figures  
149 4E, S7E, and S8E).

150 When comparing the unphased versions of XP-EHH and XP-nSL, we find that they  
151 consistently perform as well or better than their phased counterparts (Figures 5, 6, S11, and  
152 S12), except in limited circumstances where phase is known and the sweep is fairly young  
153 (sweeping allele at 0.7 frequency) or the divergence time is further in the past.

## 154 **4 Discussion**

155 We introduce multi-locus genotype versions of four popular haplotype-based selection  
156 statistics—iHS (Voight, et al. 2006), nSL (Ferrer-Admetlla, et al. 2014), XP-EHH (Sabeti, et al.  
157 2007), and XP-nSL (Szpiech, et al. 2021)—that can be used when the phase of genotypes is  
158 unknown. Although phase would seem to be a critically important component of any haplotype-  
159 based method for detecting selection, here we show that, by collapsing haplotypes into derived  
160 allele counts (thus erasing phase information), we can achieve similar power to using this  
161 information. This follows other work that has shown similar patterns with other haplotype-based  
162 statistics for detecting selection (Harris, et al. 2018; Harris and DeGiorgio 2020; DeGiorgio and  
163 Szpiech 2022). Importantly, this approach now opens up the application of several popular  
164 haplotype-based selection statistics (based on extended haplotype homozygosity) to more  
165 species where phase information is challenging to know or infer.

166 For ease of use of these new unphased versions of iHS, nSL, XP-EHH, and XP-nSL, we  
167 implement these updates in the latest v2.0 update of the program selscan (Szpiech and  
168 Hernandez 2014), with source code and pre-compiled binaries available at  
169 <https://www.github.com/szpiech/selscan>.

## 170 **5 Methods**

### 171 **5.1 Simulations**

172 We evaluate the performance of the phased and unphased versions of iHS, nSL, XP-  
173 EHH, and XP-nSL under a generic two-population divergence model using the coalescent  
174 simulation program discoal (Kern and Schrider 2016). We explore five versions of this generic  
175 model and name them Demo 1 through Demo 5 (Table 1). Let  $N_0$  and  $N_1$  be the effective  
176 population sizes of Population 0 and Population 1 after the split from their ancestral population  
177 (of size  $N_A$ ). For Demo 1, we keep a constant population size post-split and let  $N_0 = N_1 =$   
178 10,000. For Demo 2, we keep a constant population size post-split and let  $N_0 = 2N_1 = 10,000$ .  
179 For Demo 3, we keep a constant population size post-split and let  $2N_0 = N_1 = 10,000$ . For  
180 Demo 4, we initially set  $N_0 = N_1 = 10,000$  and let  $N_0$  grow stepwise exponentially every 50  
181 generations starting at 2,000 generations ago until  $N_0 = 5N_1 = 50,000$ . For Demo 5, we initially  
182 set  $N_0 = N_1 = 10,000$  and let  $N_1$  grow stepwise exponentially every 50 generations starting at  
183 2,000 generations ago until  $5N_0 = N_1 = 50,000$ .

184 For each demographic history we vary the population divergence time  $t_d \in$   
185 {2000, 4000, 8000} generations ago. For non-neutral simulations, we simulate a sweep in  
186 Population 0 in the middle of the simulated region across a range of selection coefficients  $s \in$   
187 {0.005, 0.01, 0.02}. We vary the frequency at which the adaptive allele starts sweeping as  $e \in$   
188 {0, 0.01, 0.02, 0.05, 0.10}, where  $e = 0$  indicates a hard sweep and  $e > 0$  indicates a soft sweep,  
189 and we also vary the frequency of the selected allele at time of sampling  $f \in \{0.7, 0.8, 0.9, 1.0\}$   
190 as well as  $g \in \{50, 100\}$  representing fixation of the sweeping allele  $g$  generations ago. For all  
191 simulations we set the genome length to be  $L = 500,000$  basepairs, the ancestral effective  
192 population size to be  $N_A = 10,000$ , the per site per generation mutation rate at  $\mu = 2.35 \times 10^{-8}$ ,  
193 and the per site per generation recombination rate at  $r = 1.2 \times 10^{-8}$ . For neutral simulations, we  
194 simulate 1,000 replicates for each parameter set, and for non-neutral simulations we simulate

195 100 replicates for each parameter set. We sample 200 haplotypes, randomly paired together to  
196 form 100 diploid individuals, from each population for analysis. These data sets represent the  
197 case where phase is known perfectly. We also create a set of “unphased” data sets from these  
198 phased data sets by swapping the alleles of each heterozygote to the opposing haplotype with  
199 probability 0.5.

200 As iHS and nSL are single population statistics, we only analyze Demo 1, Demo 3, and  
201 Demo 4 with these statistics, as Demo 2 and Demo 5 have a constant size history identical to  
202 Demo 1 for Population 0, where the sweeps are simulated. For XP-EHH and XP-nSL we  
203 analyze all five demographic histories.

204 For all simulations, we compute the relevant statistics (--ihs, --nsl, --xpehh, or --xpnsi)  
205 with selscan v2.0 using the --trunc-ok flag. We set --unphased when computing the unphased  
206 versions of these statistics, and we do not set it when computing the original phased versions.  
207 For iHS and XP-EHH, we also use the --pmap flag to use physical distance instead of a  
208 recombination map.

## 209 **5.2 Power and False Positive Rate**

210 Here we evaluate the power and false positive rate for the unphased version of iHS,  
211 nSL, XP-EHH, and XP-nSL. For comparison, we also compute the power for the original phased  
212 versions of these statistics in two different ways. We compute the phased statistics for a set of  
213 simulated datasets where perfect phase is known, and we compute them again for a set of  
214 simulated datasets where we destroy phase information (see section 5.1). As the unphased  
215 statistics collapse genotypes into derived allele counts, there is no functional difference between  
216 these two datasets for these statistics. We compute power in the same way for each statistic  
217 regardless of underlying dataset analyzed as described below.

218 To compute power for iHS and nSL, we follow the approach of Voight, et al. (2006). For  
219 these statistics, each non-neutral replicate is individually normalized jointly with all neutral  
220 replicates with matching demographic history in 1% allele frequency bins. Because extreme

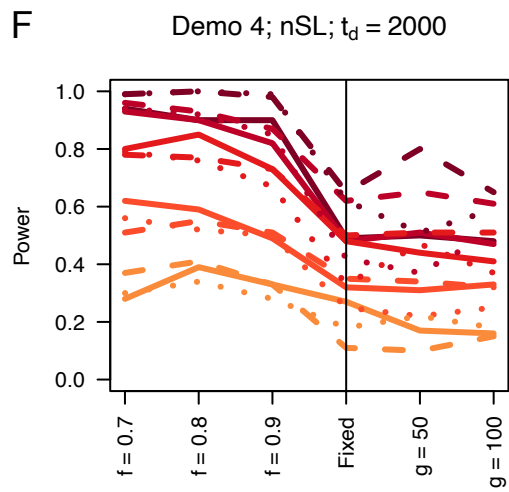
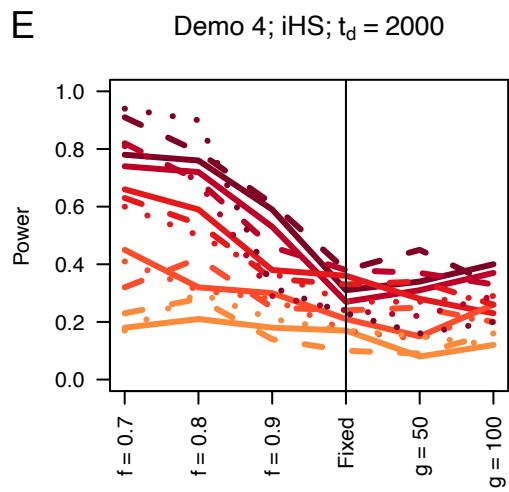
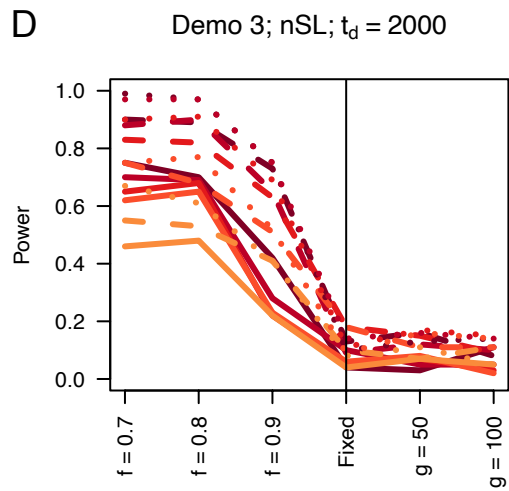
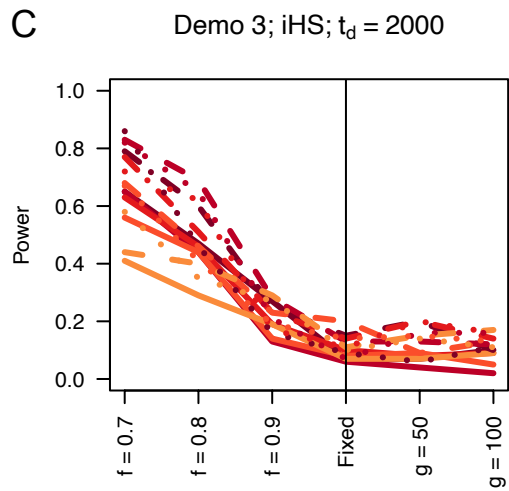
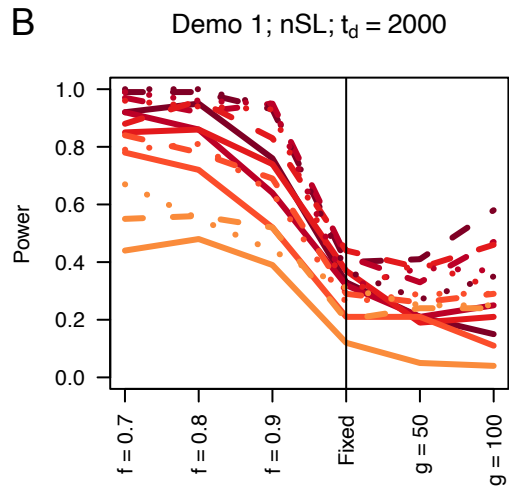
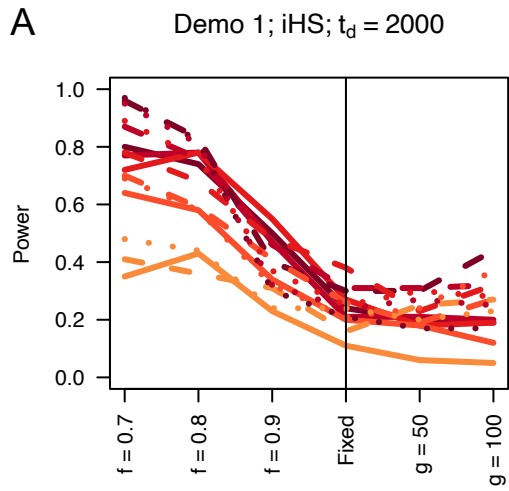
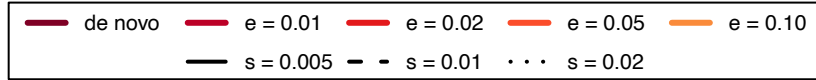
221 values of the statistic are likely to be clustered along the genome (Voight, et al. 2006), we then  
222 compute the proportion of extreme scores ( $|iHS| > 2$  or  $|nSL| > 2$ ) within 100kbp non-  
223 overlapping windows. We then bin these windows into 10 quantile bins based on the number of  
224 scores observed in each window and call the top 1% of these windows as putatively under  
225 selection. We calculate the proportion of non-neutral replicates that fall in this top 1% as the  
226 power. To compute the false positive rate, we compute the proportion of neutral simulations that  
227 fall within the top 1%.

228         To compute power for XP-EHH and XP-nSL, we follow the approach of (Szpiech, et al.  
229 2021). For these statistics, each non-neutral replicate is individually normalized jointly with all  
230 matching neutral replicates. Because extreme values of the statistic are likely to be clustered  
231 along the genome (Szpiech, et al. 2021), we then compute the proportion of extreme scores  
232 ( $XP-EHH > 2$  or  $XP-nSL > 2$ ) within 100kbp non-overlapping windows. We then bin these  
233 windows into 10 quantile bins based on the number of scores observed in each window and call  
234 the top 1% of these windows as putatively under selection. We calculate the proportion of non-  
235 neutral replicates that fall in this top 1% as the power. To compute the false positive rate, we  
236 compute the proportion of neutral simulations that fall within the top 1%.

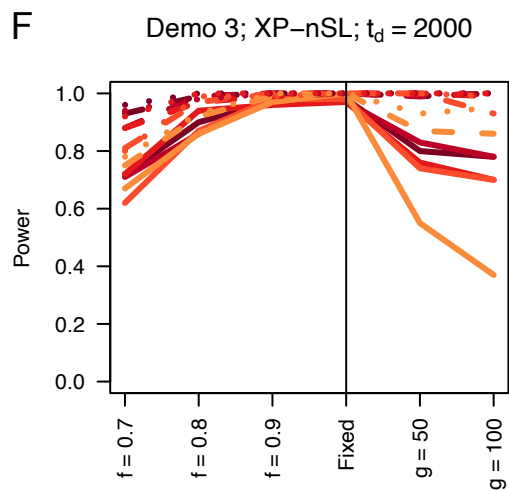
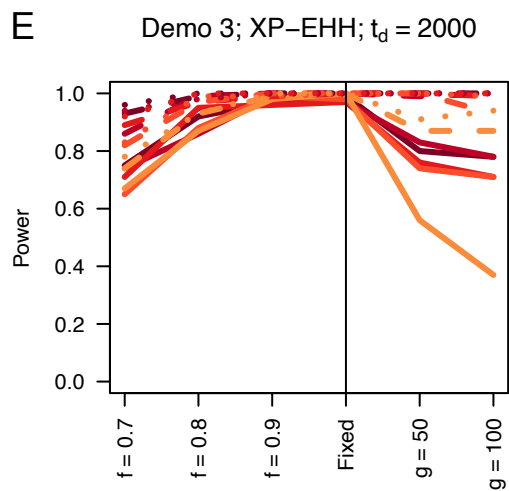
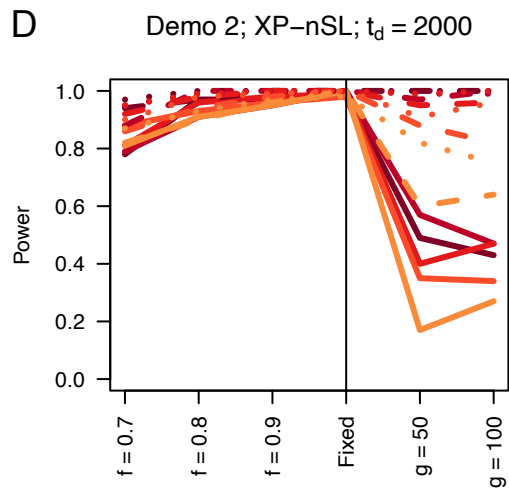
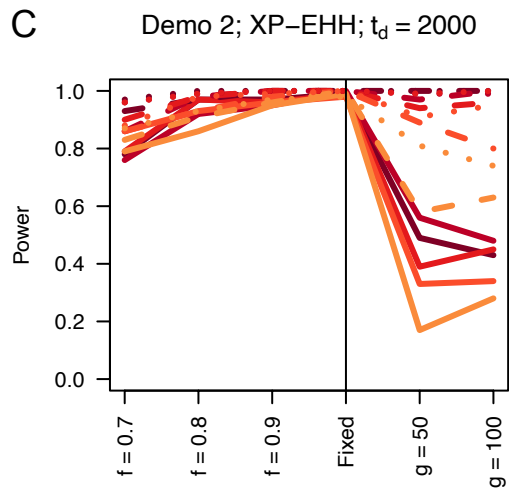
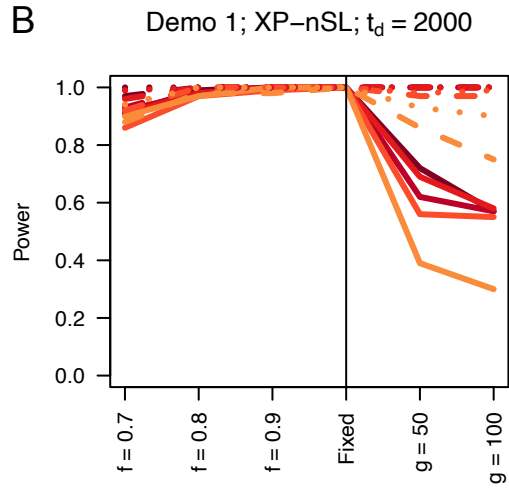
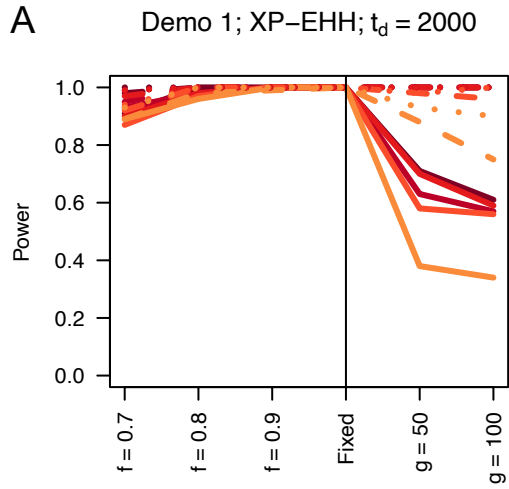
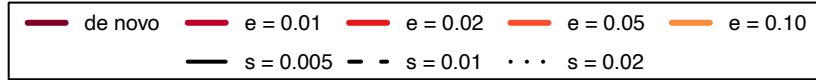
## 237 **6 Acknowledgements**

238         This work was supported by the National Institute of General Medical Sciences of the  
239 National Institutes of Health under Award Number R35GM146926 and by start-up funds from  
240 the Pennsylvania State University's Department of Biology. Computations for this research were  
241 performed using the Pennsylvania State University's Institute for Computational Data Sciences'  
242 Roar supercomputer.

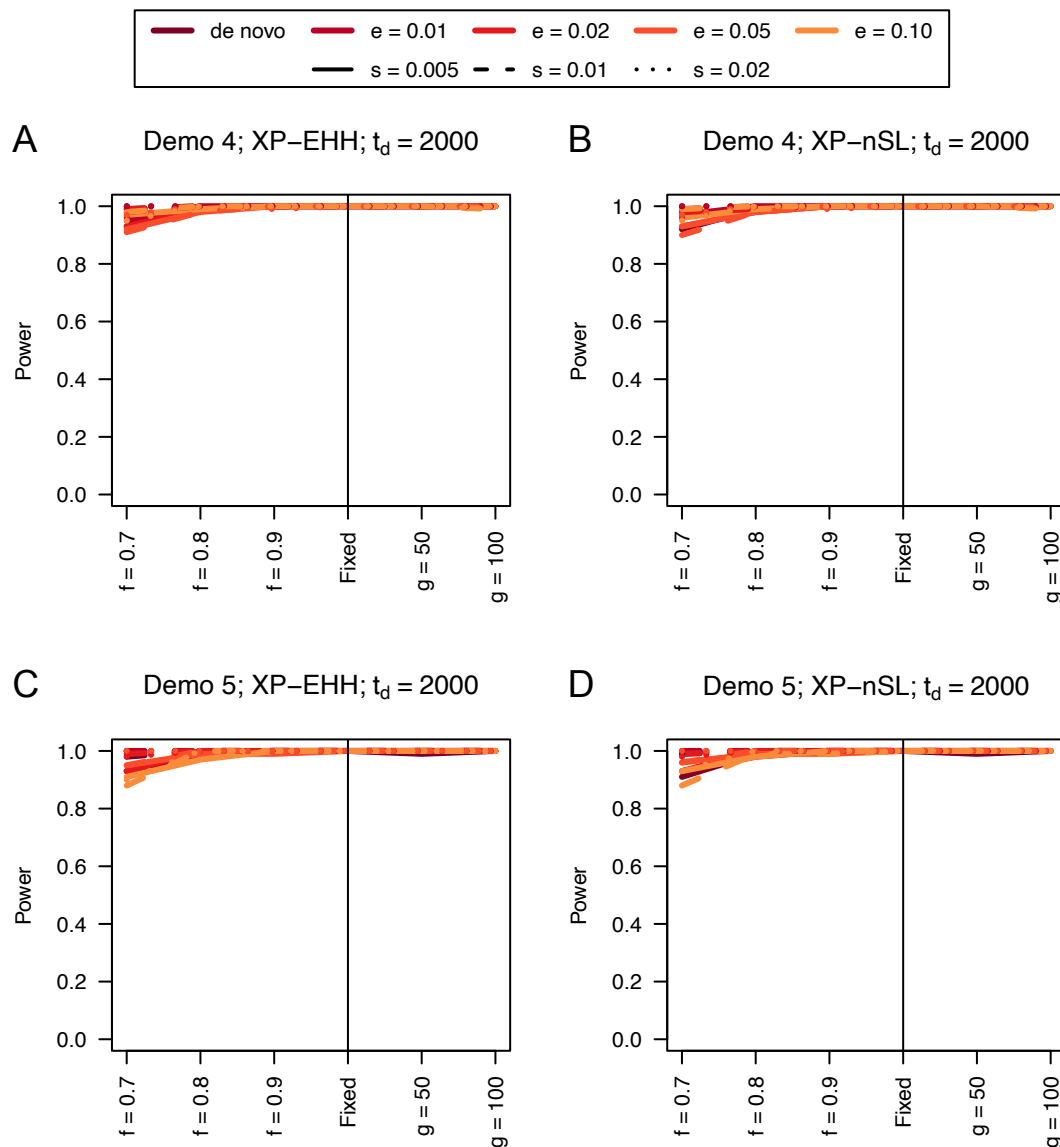
243



245 **Figure 1.** Power curves for unphased implementations of iHS (A, C, and E) and nSL (B, D, and  
246 F) under demographic histories Demo 1 (A and B), Demo 3 (C and D), and Demo 4 (E and F).  $s$   
247 is the selection coefficient,  $f$  is the frequency of the adaptive allele at time of sampling,  $g$  is the  
248 number of generations at time of sampling since fixation,  $e$  is the frequency at which selection  
249 began, and  $t_d = 2000$  is the time in generations since the two populations diverged.

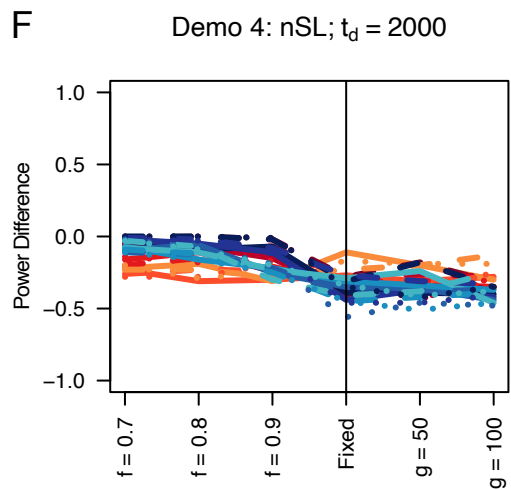
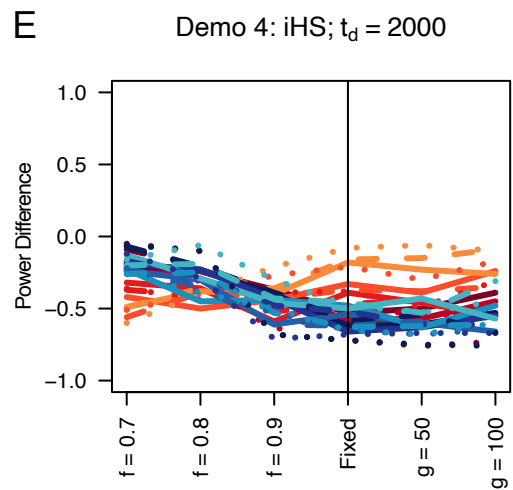
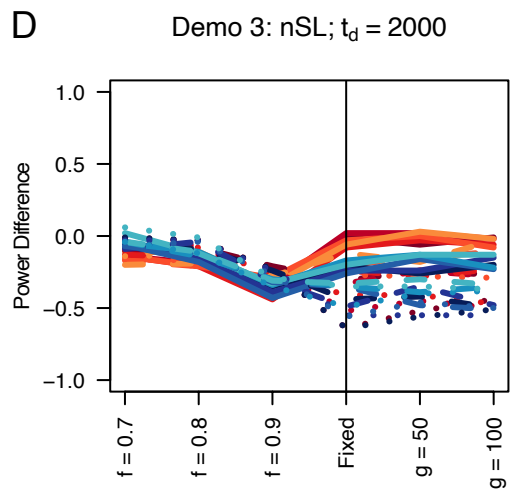
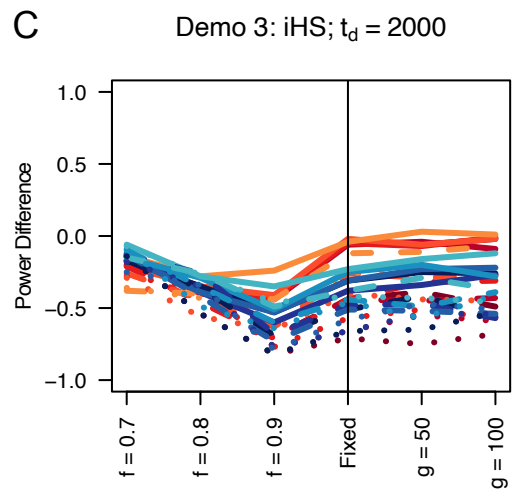
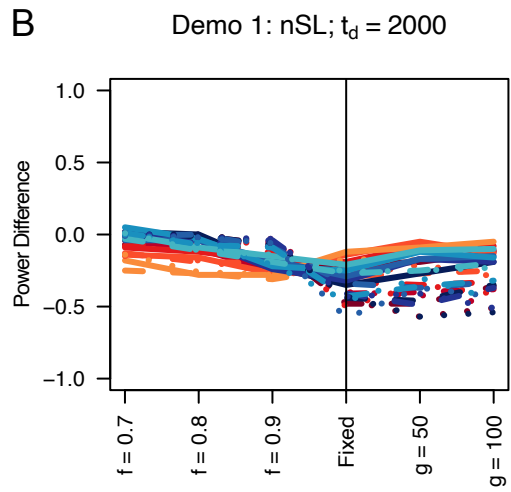
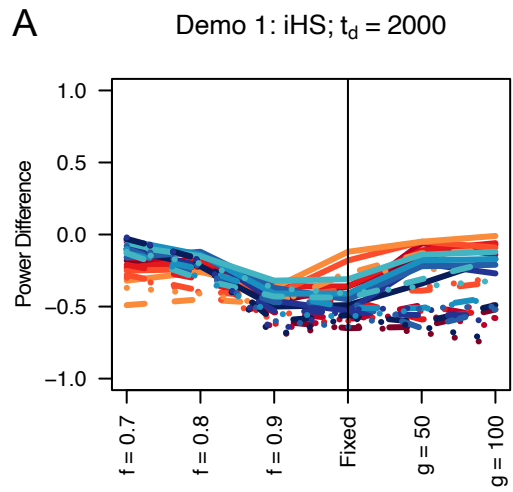
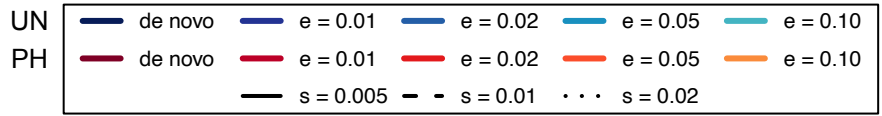


251 **Figure 2.** Power curves for unphased implementations of XP-EHH (A, C, and E) and XP-nSL  
 252 (B, D, and F) under demographic histories Demo 1 (A and B), Demo 2 (C and D), and Demo 3  
 253 (E and F).  $s$  is the selection coefficient,  $f$  is the frequency of the adaptive allele at time of  
 254 sampling,  $g$  is the number of generations at time of sampling since fixation,  $e$  is the frequency at  
 255 which selection began, and  $t_d = 2000$  is the time in generations since the two populations  
 256 diverged.

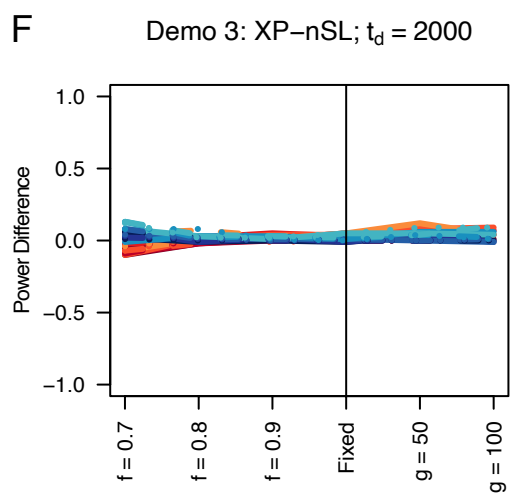
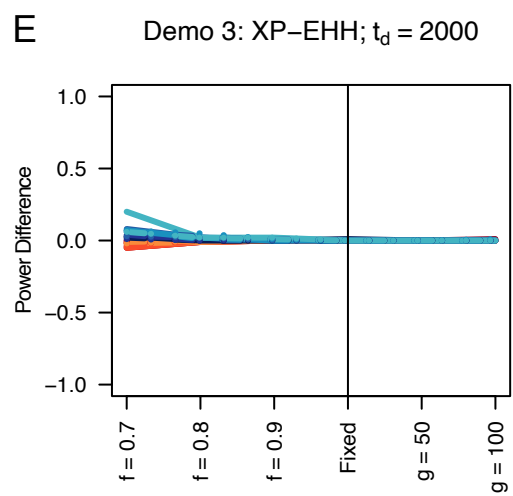
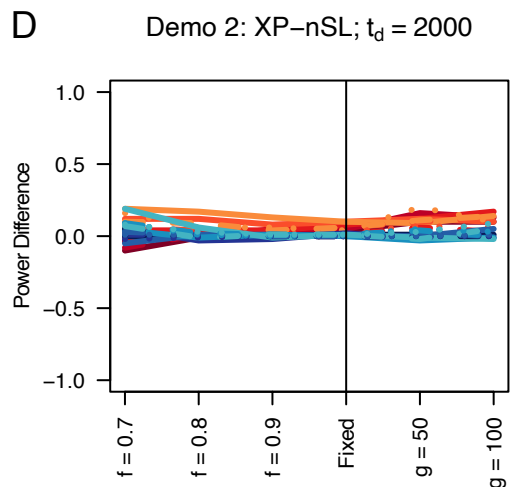
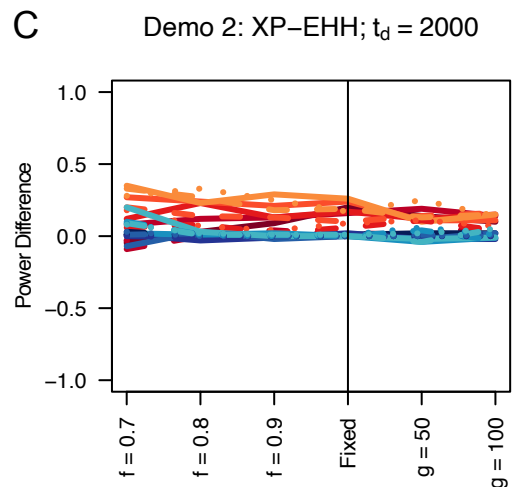
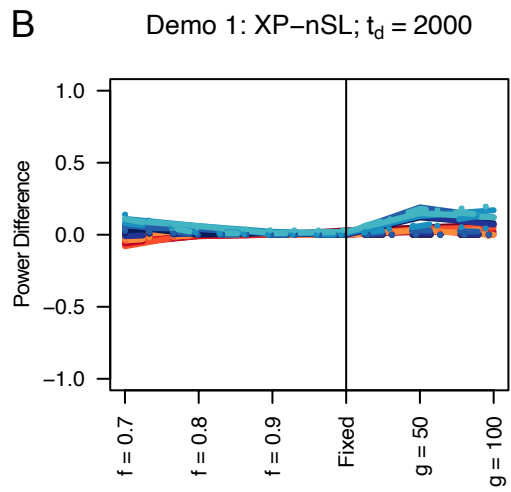
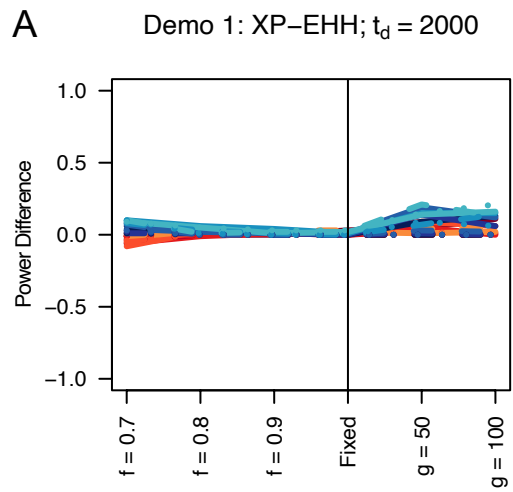
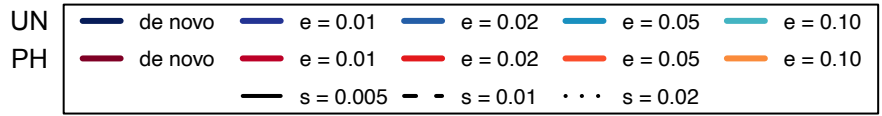


257  
 258 **Figure 3.** Power curves for unphased implementations of XP-EHH (A and C) and XP-nSL (B  
 259 and D) under demographic histories Demo 4 (A and B), and Demo 5 (C and D).  $s$  is the

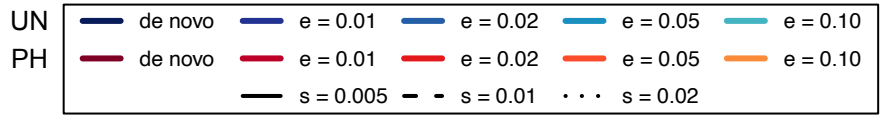
260 selection coefficient,  $f$  is the frequency of the adaptive allele at time of sampling,  $g$  is the  
261 number of generations at time of sampling since fixation,  $e$  is the frequency at which selection  
262 began, and  $t_d = 2000$  is the time in generations since the two populations diverged.



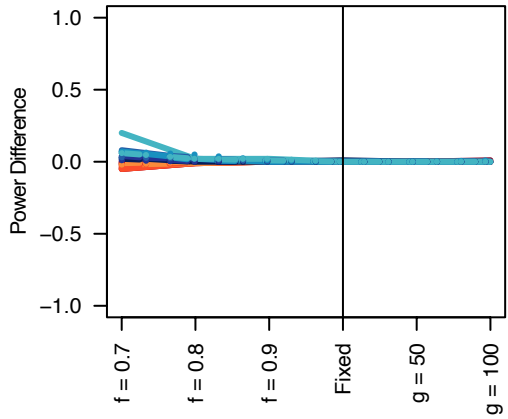
264 **Figure 4.** Power difference between unphased implementations of iHS (A, C, and E) and nSL  
265 (B, D, and F) and phased implementations. Blue curves represent the power difference between  
266 the unphased and phased statistics when applied to unphased data (UN). Red curves represent  
267 the power difference between the unphased and phased statistics when applied to perfectly  
268 phased data (PH). Values greater than 0 indicate the unphased statistic had higher power.  
269 Applied to demographic histories Demo 1 (A and B), Demo 3 (C and D), and Demo 4 (E and F).  
270  $s$  is the selection coefficient,  $f$  is the frequency of the adaptive allele at time of sampling,  $g$  is  
271 the number of generations at time of sampling since fixation,  $e$  is the frequency at which  
272 selection began, and  $t_d = 2000$  is the time in generations since the two populations diverged.



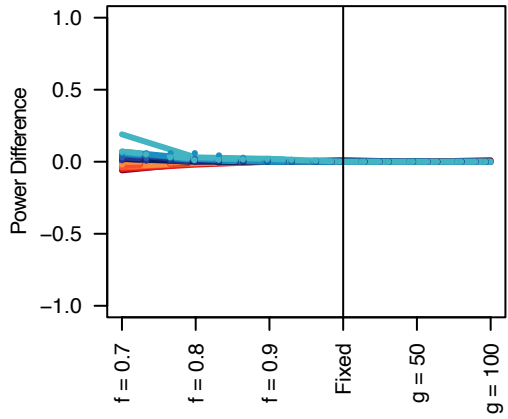
274 **Figure 5.** Power difference between unphased implementations of XP-EHH (A, C, and E) and  
275 XP-nSL (B, D, and F) and phased implementations. Blue curves represent the power difference  
276 between the unphased and phased statistics when applied to unphased data (UN). Red curves  
277 represent the power difference between the unphased and phased statistics when applied to  
278 perfectly phased data (PH). Values greater than 0 indicate the unphased statistic had higher  
279 power. Applied to demographic histories Demo 1 (A and B), Demo 2 (C and D), and Demo 3 (E  
280 and F).  $s$  is the selection coefficient,  $f$  is the frequency of the adaptive allele at time of sampling,  
281  $g$  is the number of generations at time of sampling since fixation,  $e$  is the frequency at which  
282 selection began, and  $t_d = 2000$  is the time in generations since the two populations diverged.



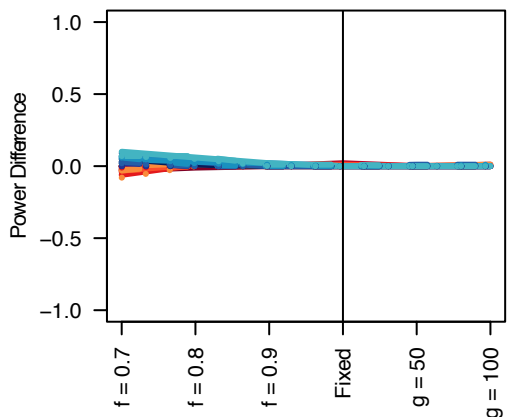
**A** Demo 4: XP-EHH;  $t_d = 2000$



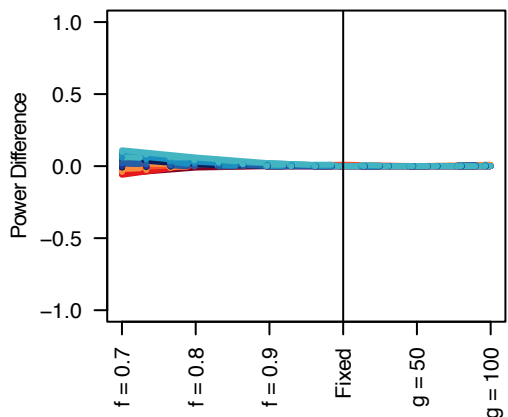
**B** Demo 4: XP-nSL;  $t_d = 2000$



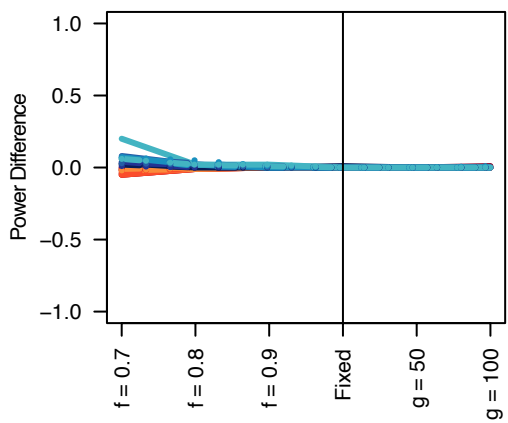
**C** Demo 5: XP-EHH;  $t_d =$



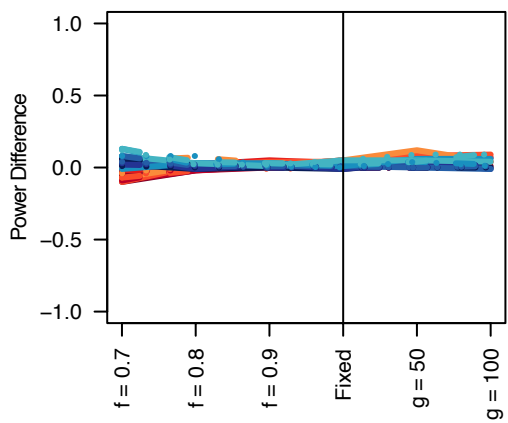
**D** Demo 5: XP-nSL;  $t_d = 2000$



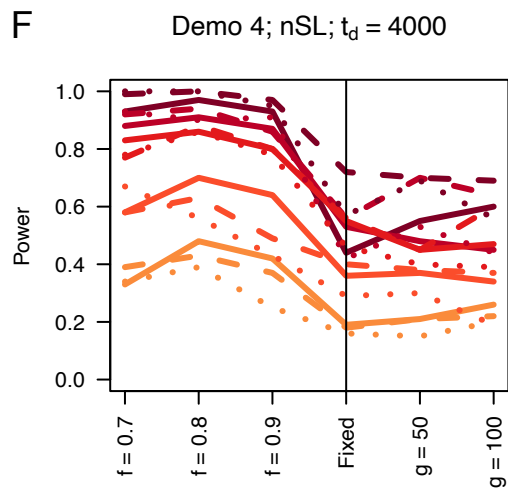
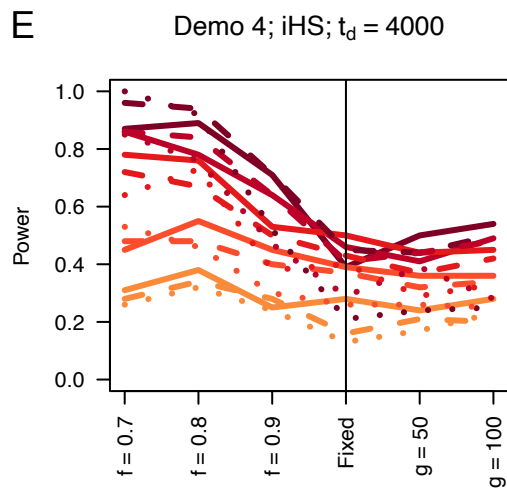
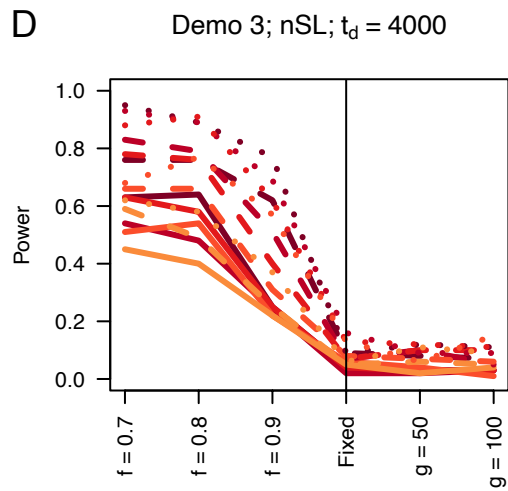
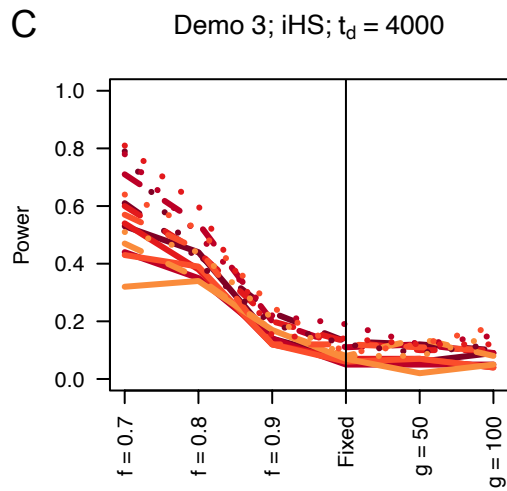
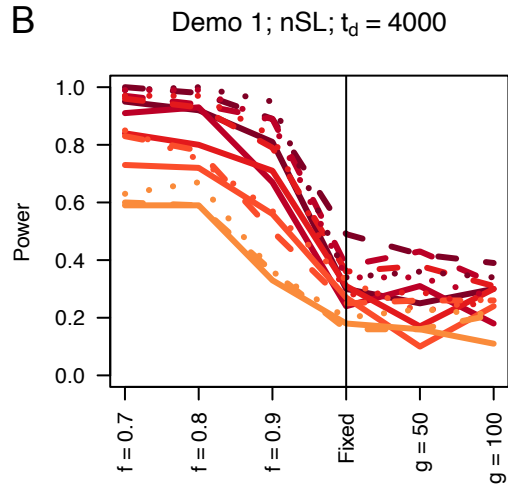
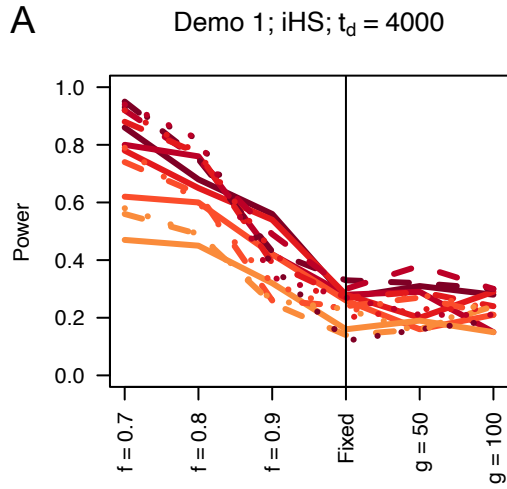
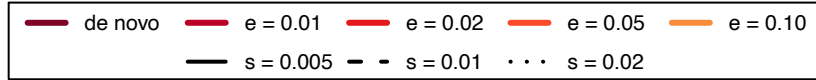
**E** Demo 3: XP-EHH;  $t_d = 2000$



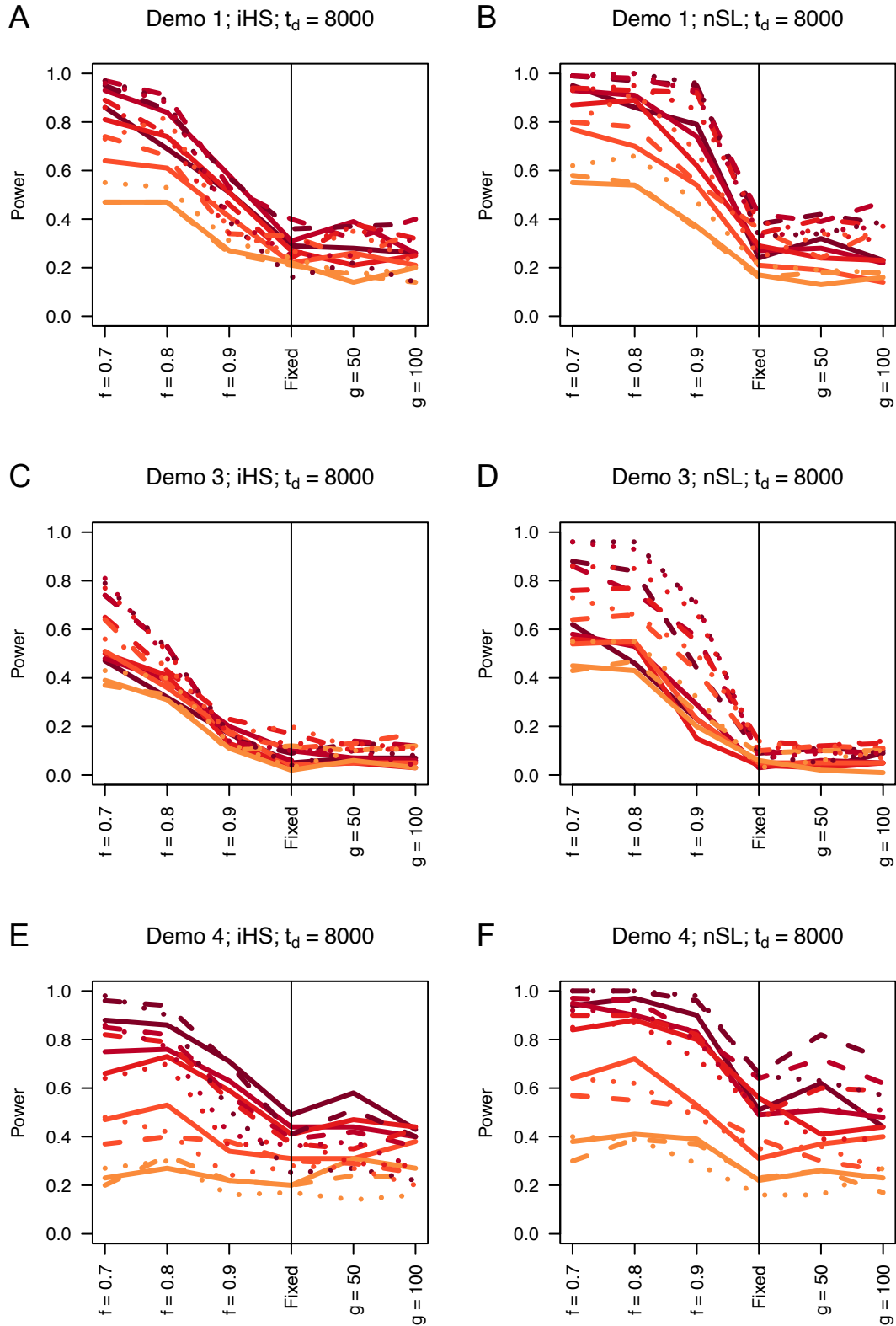
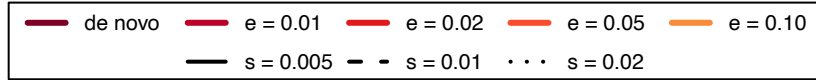
**F** Demo 3: XP-nSL;  $t_d = 2000$



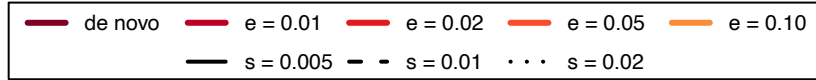
284 **Figure 6.** Power difference between unphased implementations of XP-EHH (A and C) and XP-  
285 nSL (B and D) and phased implementations. Blue curves represent the power difference  
286 between the unphased and phased statistics when applied to unphased data (UN). Red curves  
287 represent the power difference between the unphased and phased statistics when applied to  
288 perfectly phased data (PH). Values greater than 0 indicate the unphased statistic had higher  
289 power. Applied to demographic histories Demo 4 (A and B), and Demo 5 (C and D).  $s$  is the  
290 selection coefficient,  $f$  is the frequency of the adaptive allele at time of sampling,  $g$  is the  
291 number of generations at time of sampling since fixation,  $e$  is the frequency at which selection  
292 began, and  $t_d = 2000$  is the time in generations since the two populations diverged.



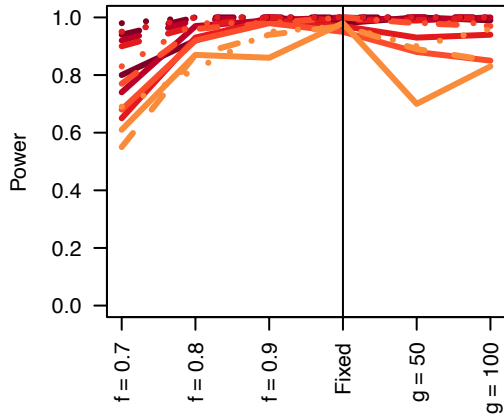
294 **Figure S1.** Power curves for unphased implementations of iHS (A, C, and E) and nSL (B, D,  
295 and F) under demographic histories Demo 1 (A and B), Demo 3 (C and D), and Demo 4 (E and  
296 F).  $s$  is the selection coefficient,  $f$  is the frequency of the adaptive allele at time of sampling,  $g$  is  
297 the number of generations at time of sampling since fixation,  $e$  is the frequency at which  
298 selection began, and  $t_d = 4000$  is the time in generations since the two populations diverged.



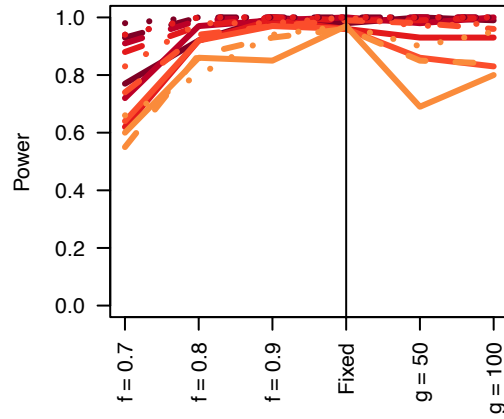
300 **Figure S2.** Power curves for unphased implementations of iHS (A, C, and E) and nSL (B, D,  
301 and F) under demographic histories Demo 1 (A and B), Demo 3 (C and D), and Demo 4 (E and  
302 F).  $s$  is the selection coefficient,  $f$  is the frequency of the adaptive allele at time of sampling,  $g$  is  
303 the number of generations at time of sampling since fixation,  $e$  is the frequency at which  
304 selection began, and  $t_d = 8000$  is the time in generations since the two populations diverged.



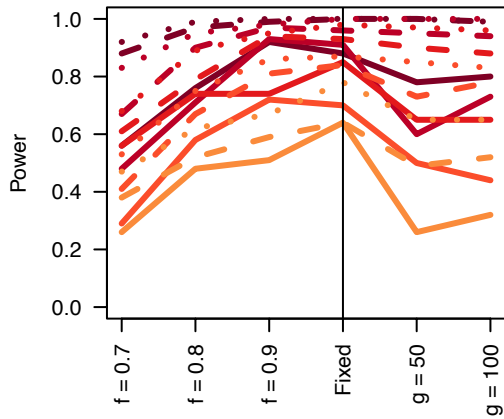
**A** Demo 1; XP-EHH;  $t_d = 4000$



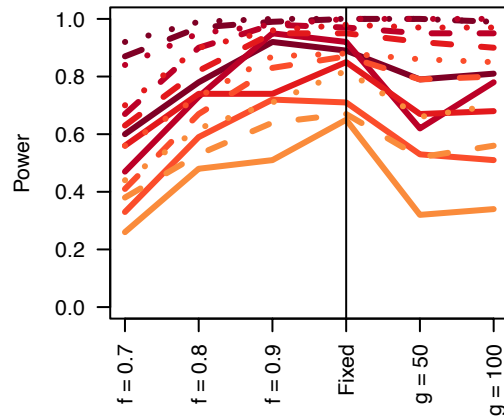
**B** Demo 1; XP-nSL;  $t_d = 4000$



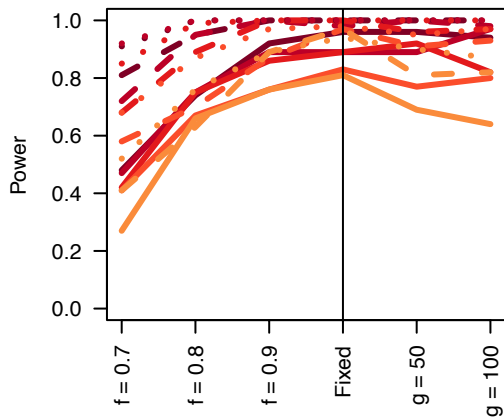
**C** Demo 2; XP-EHH;  $t_d = 4000$



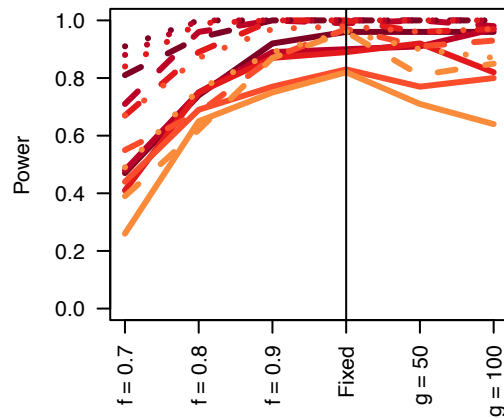
**D** Demo 2; XP-nSL;  $t_d = 4000$



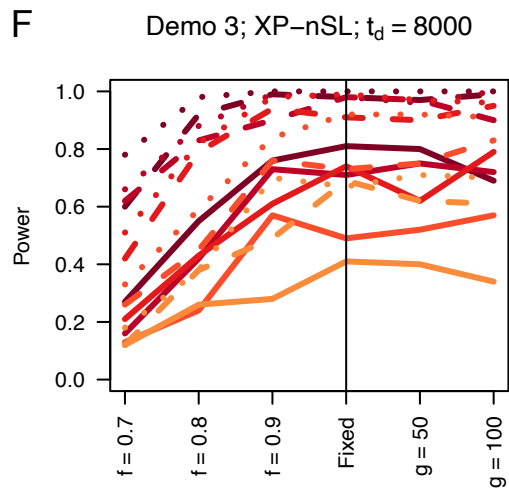
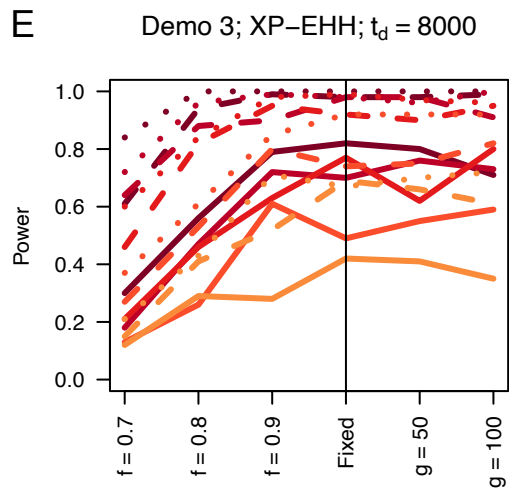
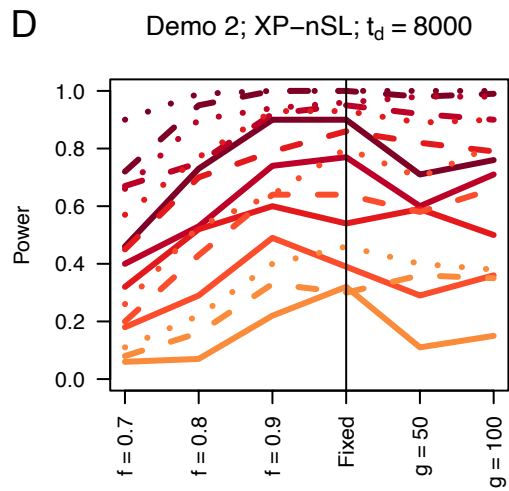
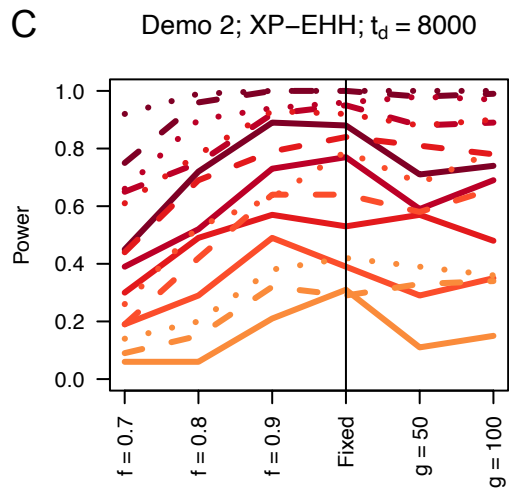
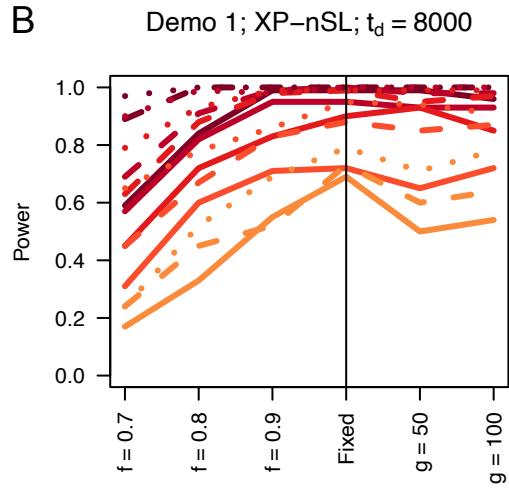
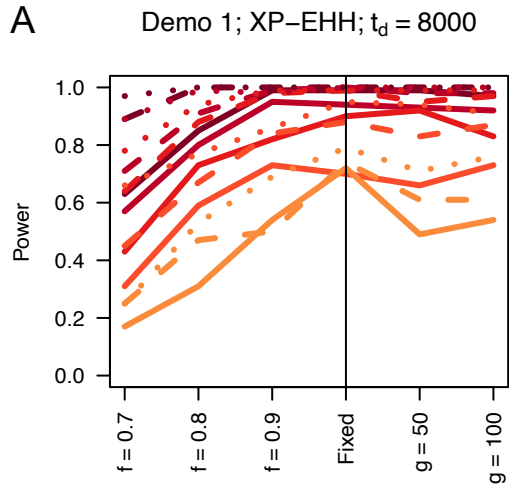
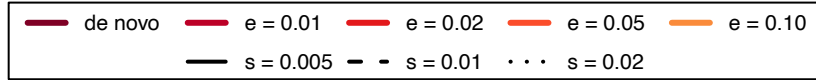
**E** Demo 3; XP-EHH;  $t_d = 4000$



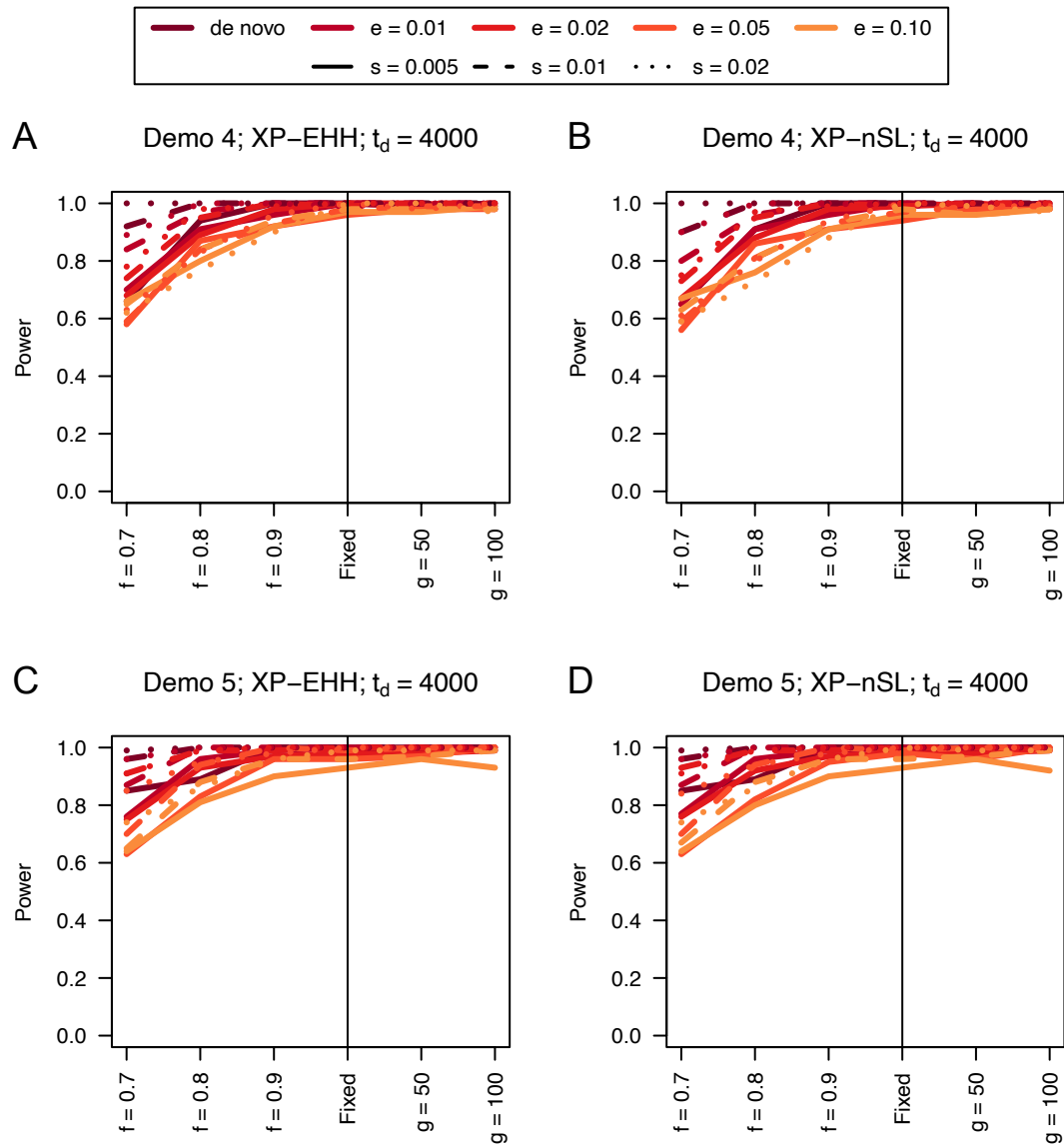
**F** Demo 3; XP-nSL;  $t_d = 4000$



306 **Figure S3.** Power difference between unphased implementations of XP-EHH (A, C, and E) and  
307 XP-nSL (B, D, and F) and phased implementations. Blue curves represent the power difference  
308 between the unphased and phased statistics when applied to unphased data (UN). Red curves  
309 represent the power difference between the unphased and phased statistics when applied to  
310 perfectly phased data (PH). Values greater than 0 indicate the unphased statistic had higher  
311 power. Applied to demographic histories Demo 1 (A and B), Demo 2 (C and D), and Demo 3 (E  
312 and F).  $s$  is the selection coefficient,  $f$  is the frequency of the adaptive allele at time of sampling,  
313  $g$  is the number of generations at time of sampling since fixation,  $e$  is the frequency at which  
314 selection began, and  $t_d = 4000$  is the time in generations since the two populations diverged.



316 **Figure S4.** Power difference between unphased implementations of XP-EHH (A, C, and E) and  
317 XP-nSL (B, D, and F) and phased implementations. Blue curves represent the power difference  
318 between the unphased and phased statistics when applied to unphased data (UN). Red curves  
319 represent the power difference between the unphased and phased statistics when applied to  
320 perfectly phased data (PH). Values greater than 0 indicate the unphased statistic had higher  
321 power. Applied to demographic histories Demo 1 (A and B), Demo 2 (C and D), and Demo 3 (E  
322 and F).  $s$  is the selection coefficient,  $f$  is the frequency of the adaptive allele at time of sampling,  
323  $g$  is the number of generations at time of sampling since fixation,  $e$  is the frequency at which  
324 selection began, and  $t_d = 8000$  is the time in generations since the two populations diverged.



325

326

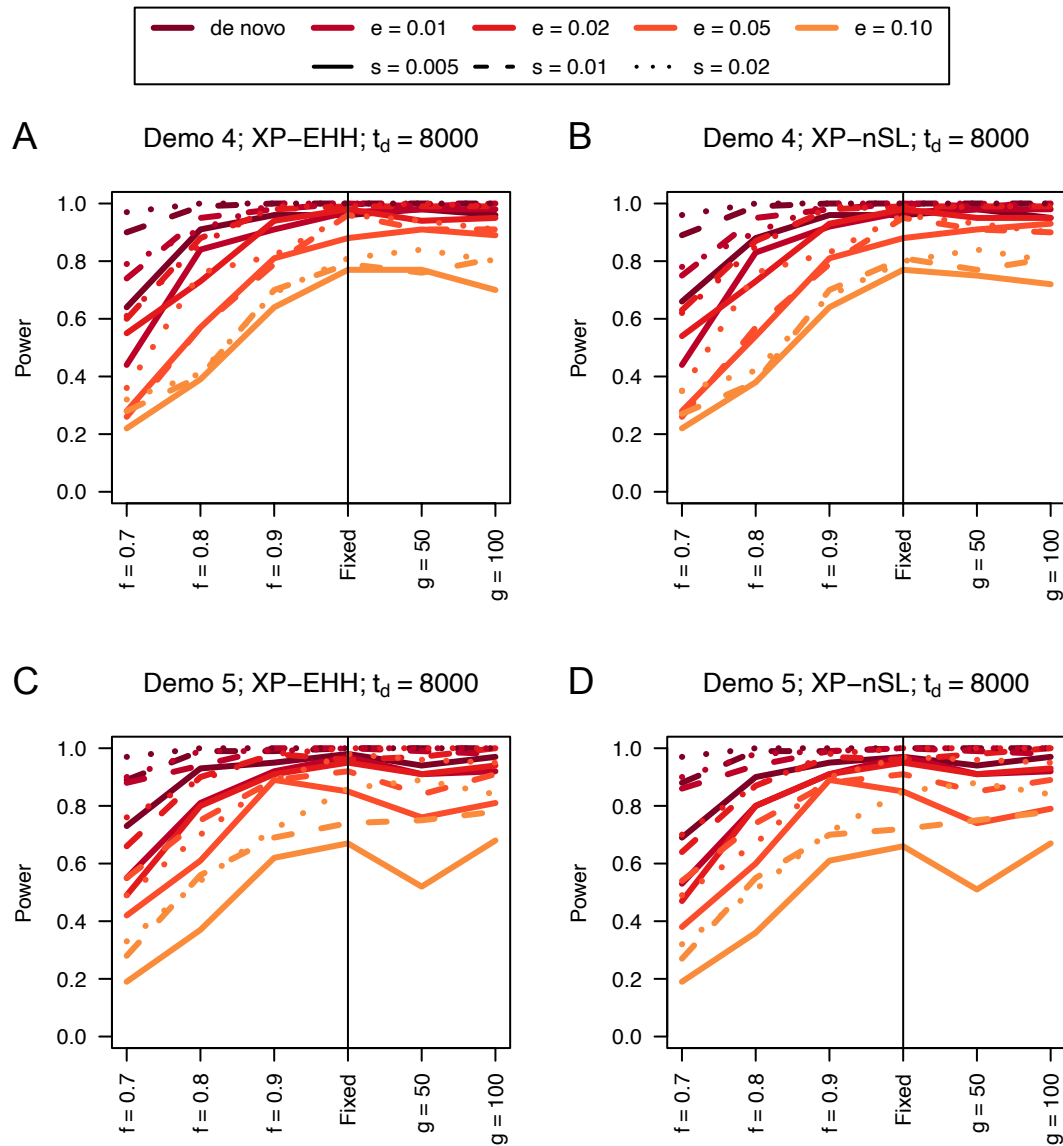
327

328

329

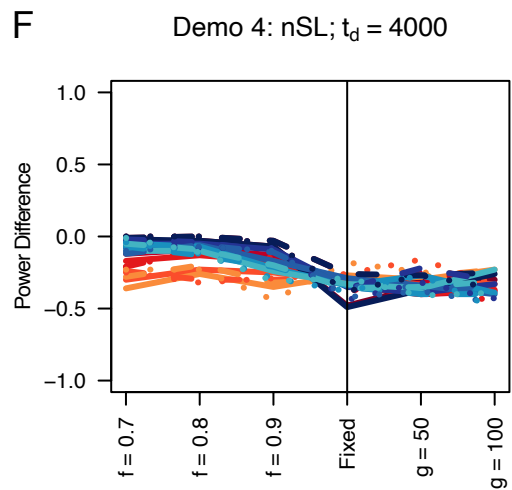
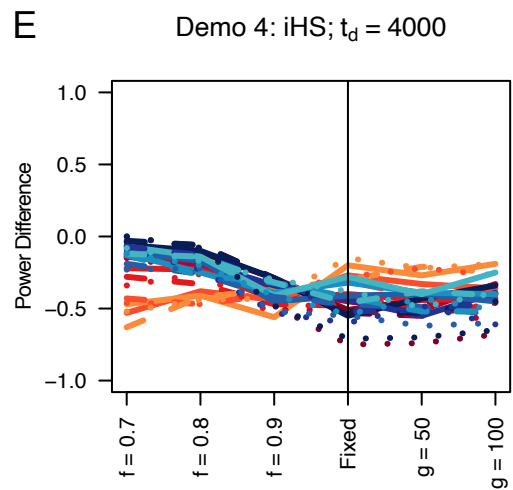
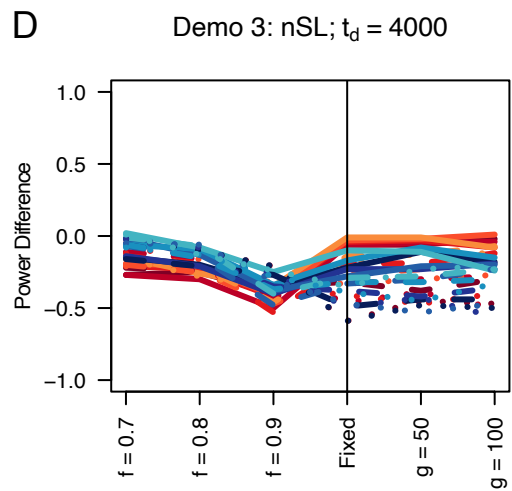
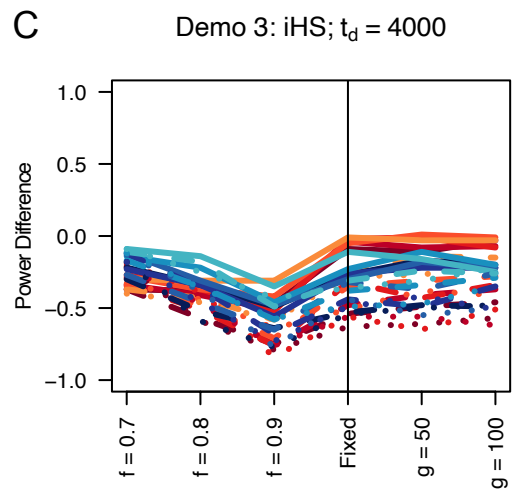
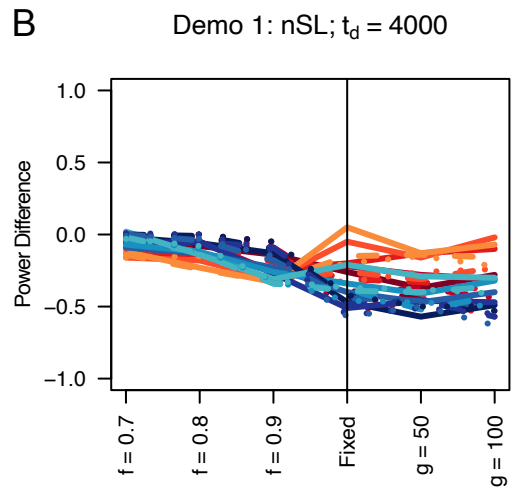
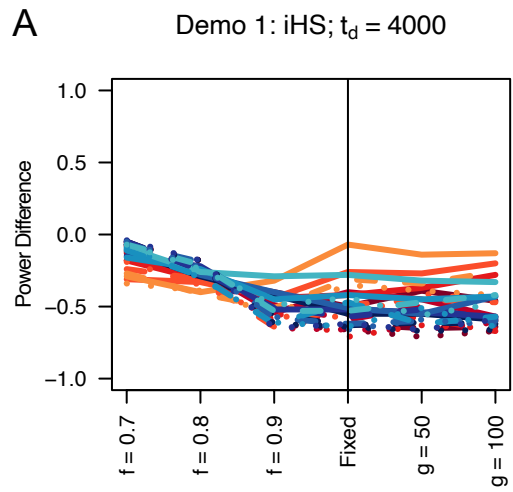
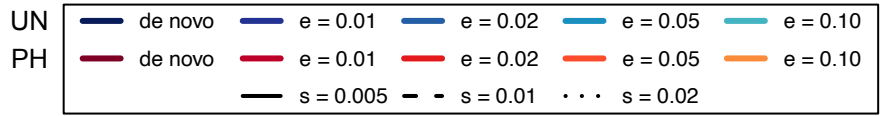
330

**Figure S5.** Power curves for unphased implementations of XP-EHH (A and C) and XP-nSL (B and D) under demographic histories Demo 4 (A and B), and Demo 5 (C and D).  $s$  is the selection coefficient,  $f$  is the frequency of the adaptive allele at time of sampling,  $g$  is the number of generations at time of sampling since fixation,  $e$  is the frequency at which selection began, and  $t_d = 4000$  is the time in generations since the two populations diverged.

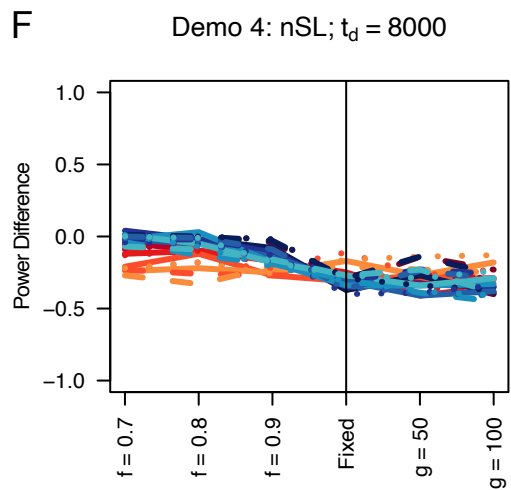
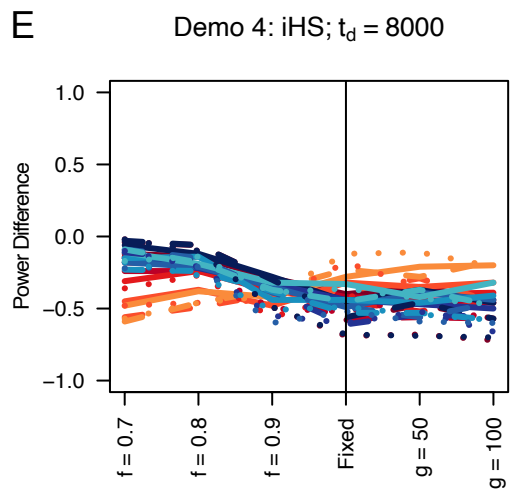
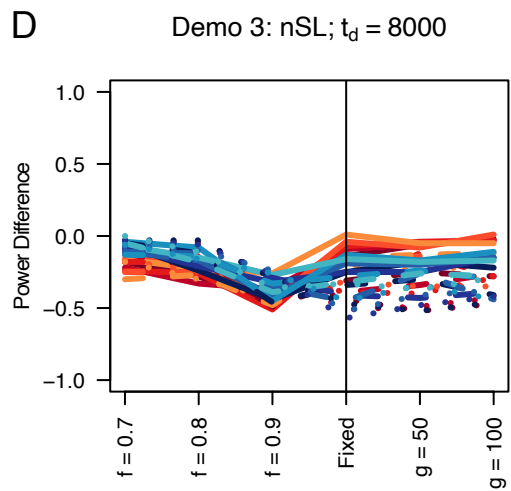
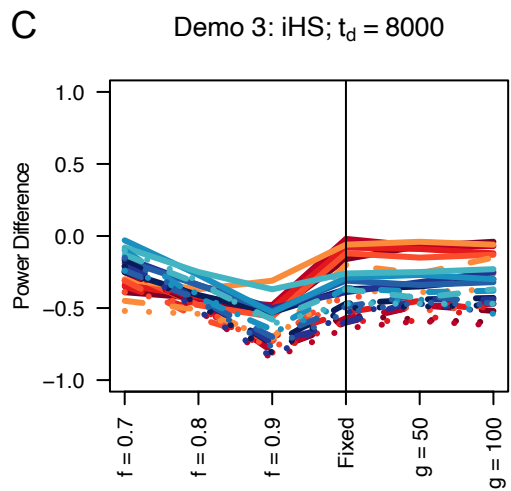
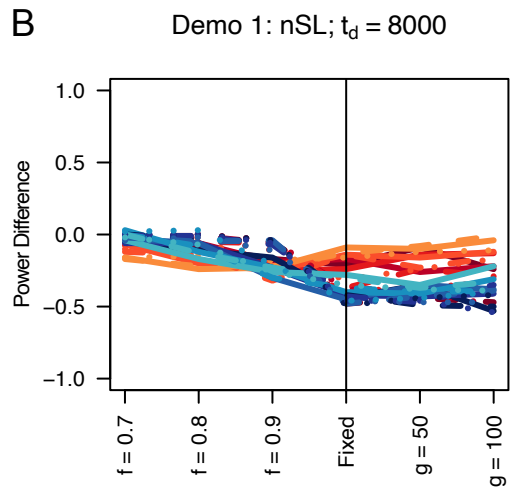
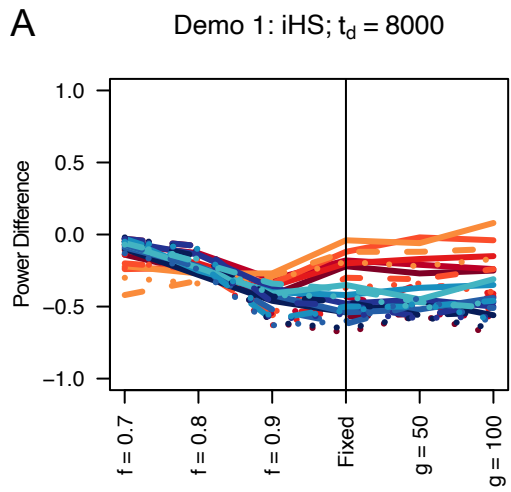
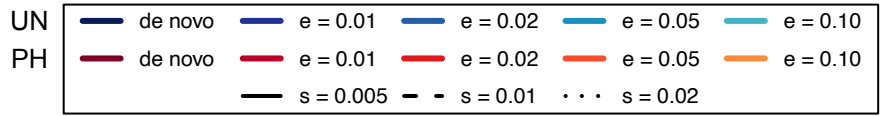


331

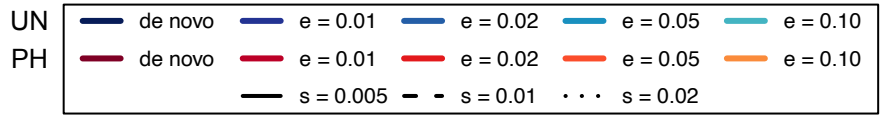
332 **Figure S6.** Power curves for unphased implementations of XP-EHH (A and C) and XP-nSL (B  
 333 and D) under demographic histories Demo 4 (A and B), and Demo 5 (C and D).  $s$  is the  
 334 selection coefficient,  $f$  is the frequency of the adaptive allele at time of sampling,  $g$  is the  
 335 number of generations at time of sampling since fixation,  $e$  is the frequency at which selection  
 336 began, and  $t_d = 8000$  is the time in generations since the two populations diverged.



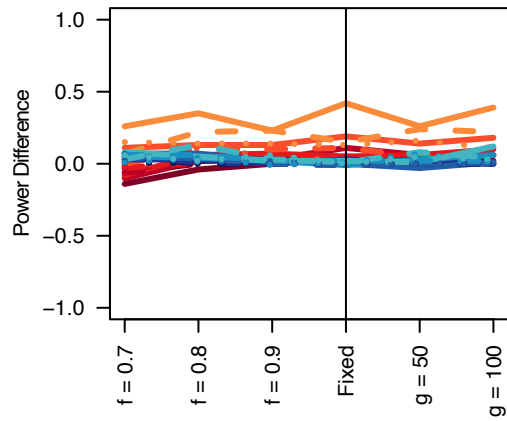
338 **Figure S7.** Power difference between unphased implementations of iHS (A, C, and E) and nSL  
339 (B, D, and F) and phased implementations. Blue curves represent the power difference between  
340 the unphased and phased statistics when applied to unphased data (UN). Red curves represent  
341 the power difference between the unphased and phased statistics when applied to perfectly  
342 phased data (PH). Values greater than 0 indicate the unphased statistic had higher power.  
343 Applied to demographic histories Demo 1 (A and B), Demo 3 (C and D), and Demo 4 (E and F).  
344  $s$  is the selection coefficient,  $f$  is the frequency of the adaptive allele at time of sampling,  $g$  is  
345 the number of generations at time of sampling since fixation,  $e$  is the frequency at which  
346 selection began, and  $t_d = 4000$  is the time in generations since the two populations diverged.



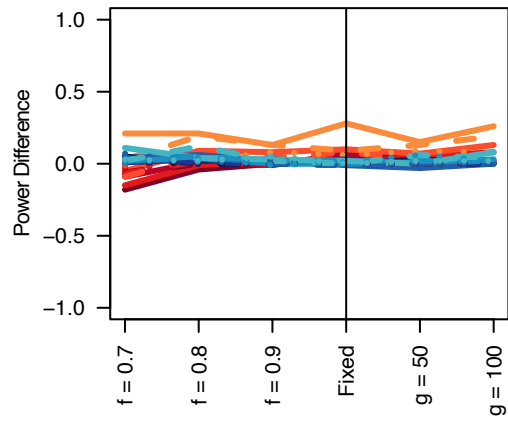
348 **Figure S8.** Power difference between unphased implementations of iHS (A, C, and E) and nSL  
349 (B, D, and F) and phased implementations. Blue curves represent the power difference between  
350 the unphased and phased statistics when applied to unphased data (UN). Red curves represent  
351 the power difference between the unphased and phased statistics when applied to perfectly  
352 phased data (PH). Values greater than 0 indicate the unphased statistic had higher power.  
353 Applied to demographic histories Demo 1 (A and B), Demo 3 (C and D), and Demo 4 (E and F).  
354  $s$  is the selection coefficient,  $f$  is the frequency of the adaptive allele at time of sampling,  $g$  is  
355 the number of generations at time of sampling since fixation,  $e$  is the frequency at which  
356 selection began, and  $t_d = 8000$  is the time in generations since the two populations diverged.



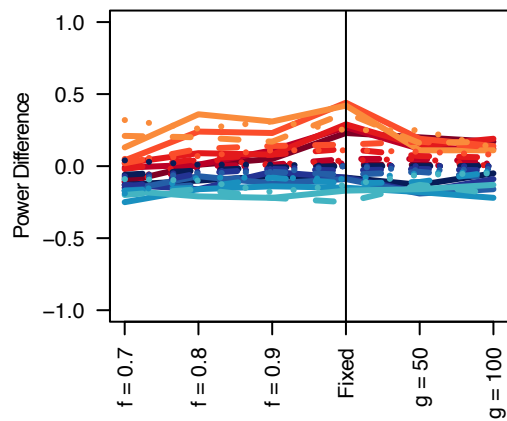
**A** Demo 1: XP-EHH;  $t_d = 4000$



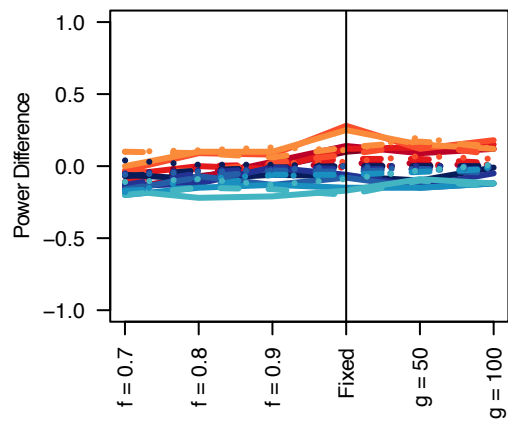
**B** Demo 1: XP-nSL;  $t_d = 4000$



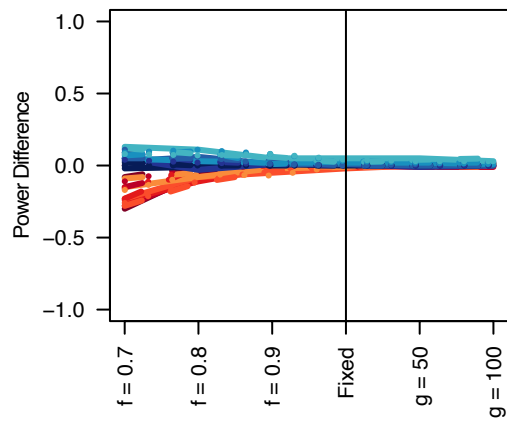
**C** Demo 2: XP-EHH;  $t_d = 4000$



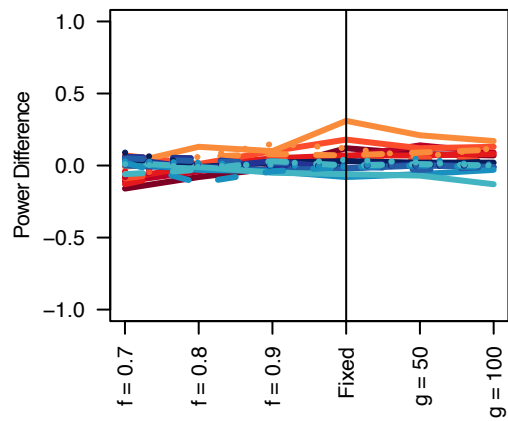
**D** Demo 2: XP-nSL;  $t_d = 4000$



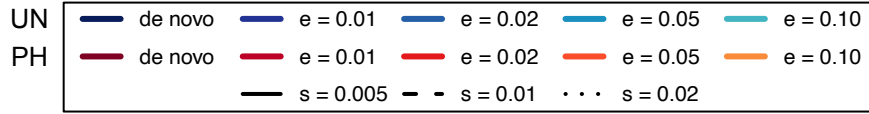
**E** Demo 3: XP-EHH;  $t_d = 4000$



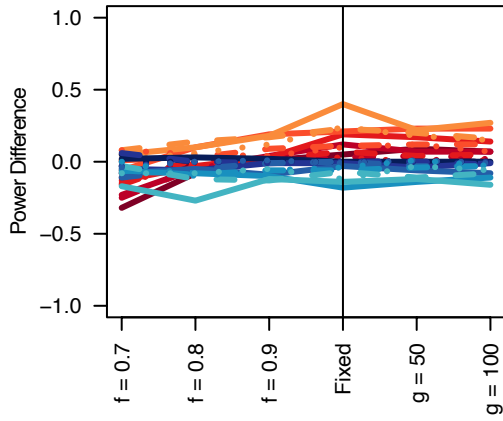
**F** Demo 3: XP-nSL;  $t_d = 4000$



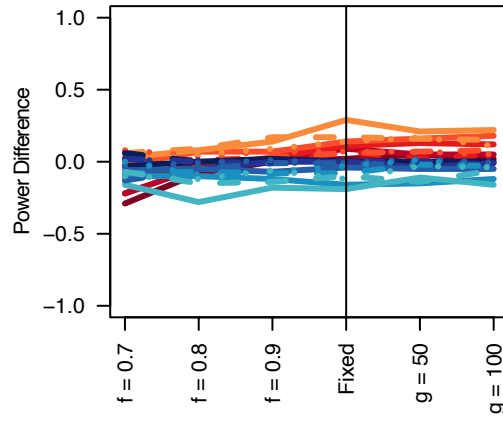
358 **Figure S9.** Power difference between unphased implementations of XP-EHH (A, C, and E) and  
359 XP-nSL (B, D, and F) and phased implementations. Blue curves represent the power difference  
360 between the unphased and phased statistics when applied to unphased data (UN). Red curves  
361 represent the power difference between the unphased and phased statistics when applied to  
362 perfectly phased data (PH). Values greater than 0 indicate the unphased statistic had higher  
363 power. Applied to demographic histories Demo 1 (A and B), Demo 2 (C and D), and Demo 3 (E  
364 and F).  $s$  is the selection coefficient,  $f$  is the frequency of the adaptive allele at time of sampling,  
365  $g$  is the number of generations at time of sampling since fixation,  $e$  is the frequency at which  
366 selection began, and  $t_d = 4000$  is the time in generations since the two populations diverged.



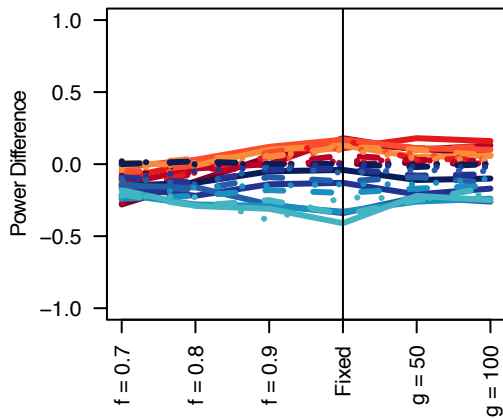
**A** Demo 1: XP-EHH;  $t_d = 8000$



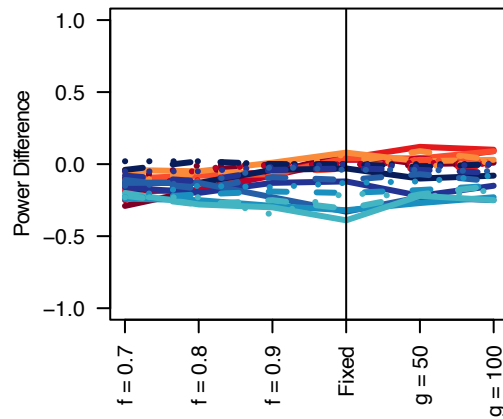
**B** Demo 1: XP-nSL;  $t_d = 8000$



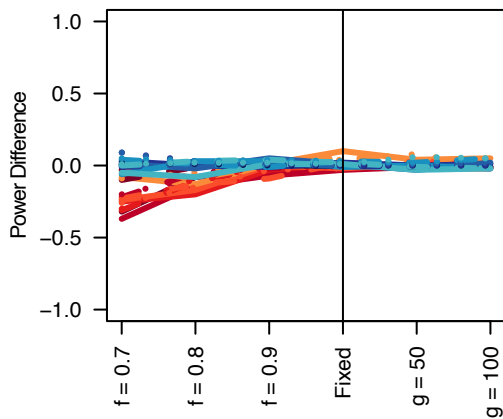
**C** Demo 2: XP-EHH;  $t_d = 8000$



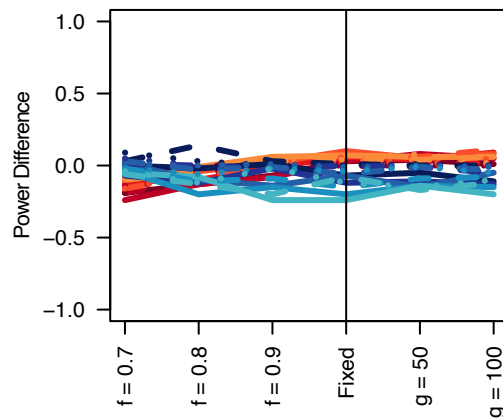
**D** Demo 2: XP-nSL;  $t_d = 8000$



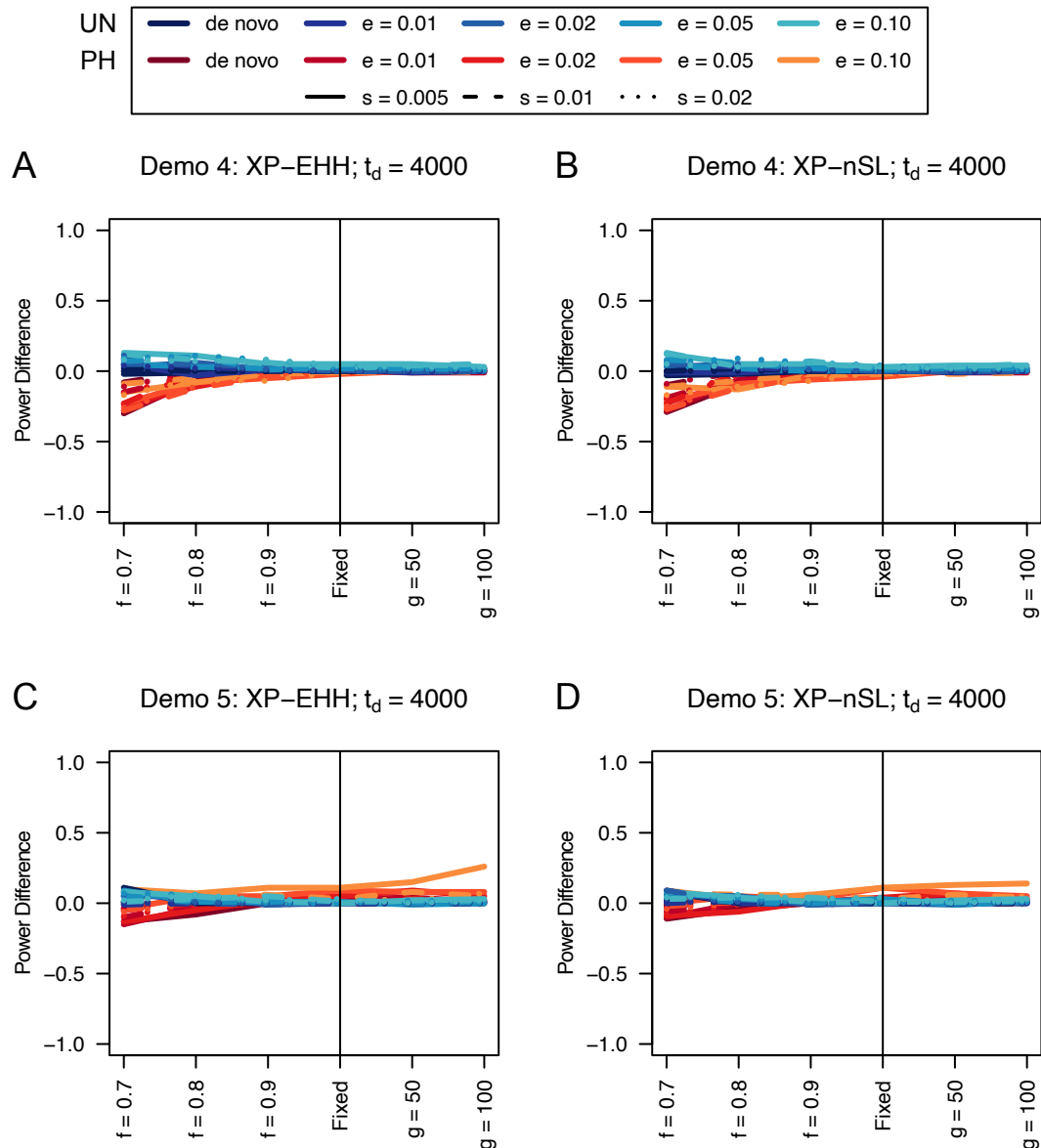
**E** Demo 3: XP-EHH;  $t_d = 8000$



**F** Demo 3: XP-nSL;  $t_d = 8000$



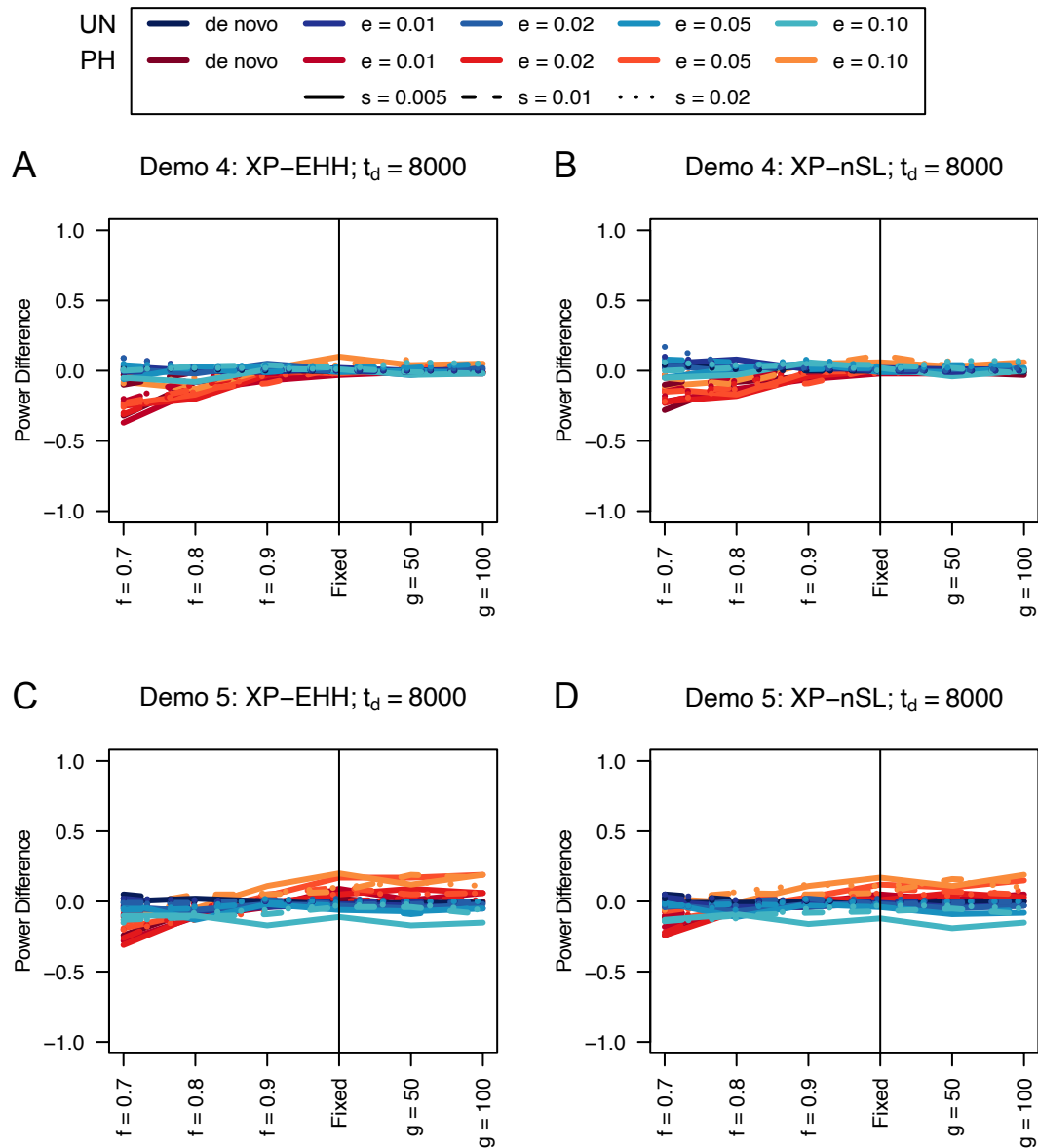
368 **Figure S10.** Power difference between unphased implementations of XP-EHH (A, C, and E)  
369 and XP-nSL (B, D, and F) and phased implementations. Blue curves represent the power  
370 difference between the unphased and phased statistics when applied to unphased data (UN).  
371 Red curves represent the power difference between the unphased and phased statistics when  
372 applied to perfectly phased data (PH). Values greater than 0 indicate the unphased statistic had  
373 higher power. Applied to demographic histories Demo 1 (A and B), Demo 2 (C and D), and  
374 Demo 3 (E and F).  $s$  is the selection coefficient,  $f$  is the frequency of the adaptive allele at time  
375 of sampling,  $g$  is the number of generations at time of sampling since fixation,  $e$  is the frequency  
376 at which selection began, and  $t_d = 8000$  is the time in generations since the two populations  
377 diverged.



378

379 **Figure S11.** Power difference between unphased implementations of XP-EHH (A and C) and  
 380 XP-nSL (B and D) and phased implementations. Blue curves represent the power difference  
 381 between the unphased and phased statistics when applied to unphased data (UN). Red curves  
 382 represent the power difference between the unphased and phased statistics when applied to  
 383 perfectly phased data (PH). Values greater than 0 indicate the unphased statistic had higher  
 384 power. Applied to demographic histories Demo 4 (A and B), and Demo 5 (C and D).  $s$  is the  
 385 selection coefficient,  $f$  is the frequency of the adaptive allele at time of sampling,  $g$  is the

386 number of generations at time of sampling since fixation,  $e$  is the frequency at which selection  
 387 began, and  $t_d = 4000$  is the time in generations since the two populations diverged.



388  
 389 **Figure S12.** Power difference between unphased implementations of XP-EHH (A and C) and  
 390 XP-nSL (B and D) and phased implementations. Blue curves represent the power difference  
 391 between the unphased and phased statistics when applied to unphased data (UN). Red curves  
 392 represent the power difference between the unphased and phased statistics when applied to  
 393 perfectly phased data (PH). Values greater than 0 indicate the unphased statistic had higher  
 394 power. Applied to demographic histories Demo 4 (A and B), and Demo 5 (C and D).  $s$  is the

395 selection coefficient,  $f$  is the frequency of the adaptive allele at time of sampling,  $g$  is the  
 396 number of generations at time of sampling since fixation,  $e$  is the frequency at which selection  
 397 began, and  $t_d = 8000$  is the time in generations since the two populations diverged.

398  
 399 **Table 1.** Demographic history parameters for simulations.  $N_A$  represents the ancestral effective  
 400 population size.  $N_0$  represents the effective population size of the population experiencing the  
 401 sweep.  $N_0$  represents the effective population size of the non-sweep population.  $t_d$  represents  
 402 the split time between the two populations.

	$N_A$	$N_0$ at split	$N_0$ at present	$N_1$ at split	$N_1$ at present	$t_d$
Demo 1	10,000	10,000	10,000	10,000	10,000	2,000/4,000/8,000
Demo 2	10,000	10,000	10,000	5,000	5,000	2,000/4,000/8,000
Demo 3	10,000	5,000	5,000	10,000	10,000	2,000/4,000/8,000
Demo 4	10,000	10,000	50,000 <sup>†</sup>	10,000	10,000	2,000/4,000/8,000
Demo 5	10,000	10,000	10,000	10,000	50,000 <sup>†</sup>	2,000/4,000/8,000

403 <sup>†</sup>The reached via exponential growth starting 2,000 generations ago.

404  
 405 **Table S1.** False positive rate computed from neutral simulations for varying  $t_d$  and demographic  
 406 history.

		$t_d = 2000$	$t_d = 4000$	$t_d = 8000$
iHS	Demo 1	0.013	0.1	0.009
	Demo 3	0.007	0.013	0.007
	Demo 4	0.015	0.018	0.008
nSL	Demo 1	0.01	0.015	0.008
	Demo 3	0.008	0.011	0.007
	Demo 4	0.014	0.021	0.014
XP-EHH	Demo 1	0.013	0.013	0.016
	Demo 2	0.017	0.009	0.015
	Demo 3	0.01	0.011	0.012
	Demo 4	0.012	0.014	0.014
	Demo 5	0.011	0.012	0.013
XP-nSL	Demo 1	0.014	0.011	0.013
	Demo 2	0.019	0.011	0.012
	Demo 3	0.011	0.011	0.012
	Demo 4	0.012	0.012	0.014
	Demo 5	0.011	0.012	0.014

407  
 408  
 409  
 410  
 411  
 412

## 413 References

414  
 415 Browning BL, Tian X, Zhou Y, Browning SR. 2021. Fast two-stage phasing of large-scale  
 416 sequence data. Am J Hum Genet 108:1880-1890.

417 Campagna L, Toews DPL. 2022. The genomics of adaptation in birds. *Curr Biol* 32:R1173-  
418 R1186.

419 Colonna V, Ayub Q, Chen Y, Pagani L, Luisi P, Pybus M, Garrison E, Xue Y, Tyler-Smith C,  
420 Genomes Project C, et al. 2014. Human genomic regions with exceptionally high levels  
421 of population differentiation identified from 911 whole-genome sequences. *Genome*  
422 *Biol* 15:R88.

423 Crawford NG, Kelly DE, Hansen MEB, Beltrame MH, Fan S, Bowman SL, Jewett E, Ranciaro A,  
424 Thompson S, Lo Y, et al. 2017. Loci associated with skin pigmentation identified in  
425 African populations. *Science* 358.

426 DeGiorgio M, Szpiech ZA. 2022. A spatially aware likelihood test to detect sweeps from  
427 haplotype distributions. *PLoS Genet* 18:e1010134.

428 Delaneau O, Zagury JF, Marchini J. 2013. Improved whole-chromosome phasing for disease  
429 and population genetic studies. *Nat Methods* 10:5-6.

430 Ellegren H. 2014. Genome sequencing and population genomics in non-model organisms.  
431 *Trends Ecol Evol* 29:51-63.

432 Fagny M, Patin E, Enard D, Barreiro LB, Quintana-Murci L, Laval G. 2014. Exploring the  
433 occurrence of classic selective sweeps in humans using whole-genome sequencing data  
434 sets. *Mol Biol Evol* 31:1850-1868.

435 Ferrer-Admetlla A, Liang M, Korneliusen T, Nielsen R. 2014. On detecting incomplete soft  
436 or hard selective sweeps using haplotype structure. *Mol Biol Evol* 31:1275-1291.

437 Harris AM, DeGiorgio M. 2020. A likelihood approach for uncovering selective sweep  
438 signatures from haplotype data. *Mol Biol Evol*.

439 Harris AM, Garud NR, DeGiorgio M. 2018. Detection and Classification of Hard and Soft  
440 Sweeps from Unphased Genotypes by Multilocus Genotype Identity. *Genetics* 210:1429-  
441 1452.

442 Kern AD, Schrider DR. 2016. Discoal: flexible coalescent simulations with selection.  
443 *Bioinformatics* 32:3839-3841.

444 Lu K, Wei L, Li X, Wang Y, Wu J, Liu M, Zhang C, Chen Z, Xiao Z, Jian H, et al. 2019. Whole-  
445 genome resequencing reveals *Brassica napus* origin and genetic loci involved in its  
446 improvement. *Nat Commun* 10:1154.

447 Meier JI, Marques DA, Wagner CE, Excoffier L, Seehausen O. 2018. Genomics of Parallel  
448 Ecological Speciation in Lake Victoria Cichlids. *Mol Biol Evol* 35:1489-1506.

449 Nedelec Y, Sanz J, Baharian G, Szpiech ZA, Pacis A, Dumaine A, Grenier JC, Freiman A, Sams  
450 AJ, Hebert S, et al. 2016. Genetic Ancestry and Natural Selection Drive Population  
451 Differences in Immune Responses to Pathogens. *Cell* 167:657-669 e621.

452 Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, Cotsapas C, Xie X, Byrne EH, McCarroll  
453 SA, Gaudet R, et al. 2007. Genome-wide detection and characterization of positive  
454 selection in human populations. *Nature* 449:913-918.

455 Salmon P, Jacobs A, Ahren D, Biard C, Dingemanse NJ, Dominoni DM, Helm B, Lundberg M,  
456 Senar JC, Sprau P, et al. 2021. Continent-wide genomic signatures of adaptation to  
457 urbanisation in a songbird across Europe. *Nat Commun* 12:2983.

458 Schriber DR. 2020. Background Selection Does Not Mimic the Patterns of Genetic Diversity  
459 Produced by Selective Sweeps. *Genetics* 216:499-519.

460 Szpiech ZA, Hernandez RD. 2014. selscan: an efficient multithreaded program to perform  
461 EHH-based scans for positive selection. *Mol Biol Evol* 31:2824-2827.

462 Szpiech ZA, Novak TE, Bailey NP, Stevison LS. 2021. Application of a novel haplotype-based  
463 scan for local adaptation to study high-altitude adaptation in rhesus macaques. *Evol*  
464 *Lett* 5:408-421.

465 Voight BF, Kudravalli S, Wen X, Pritchard JK. 2006. A map of recent positive selection in  
466 the human genome. *Plos Biology* 4:e72.

467 Zhang SJ, Wang GD, Ma P, Zhang LL, Yin TT, Liu YH, Otecko NO, Wang M, Ma YP, Wang L, et  
468 al. 2020. Genomic regions under selection in the feralization of the dingoes. *Nat*  
469 *Commun* 11:671.

470 Zoledziwska M, Sidore C, Chiang CWK, Sanna S, Mulas A, Steri M, Busonero F, Marcus JH,  
471 Marongiu M, Maschio A, et al. 2015. Height-reducing variants and selection for short  
472 stature in Sardinia. *Nat Genet* 47:1352-1356.

473