

## Genetics and population analysis

# selscan 2.0: scanning for sweeps in unphased data

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

<sup>1</sup>Department of Biology, Penn State University, University Park, PA 16802, United States

<sup>2</sup>Institute for Computational and Data Sciences, Penn State University, University Park, PA 16802, United States

\*Corresponding author. Department of Biology, Penn State University, 514A Wartik Laboratory, University Park, PA 16802, United States.

E-mail: szpiech@psu.edu

Associate Editor: Russell Schwartz

### Abstract

**Summary:** Several popular haplotype-based statistics for identifying recent or ongoing positive selection in genomes require knowledge of haplotype phase. Here, we provide an update to selscan which implements a re-definition of these statistics for use in unphased data.

**Availability and implementation:** Source code and binaries are freely available at <https://github.com/szpiech/selscan>, implemented in C/C++, and supported on Linux, Windows, and MacOS.

### 1 Introduction

Haplotype-based summary statistics—such as iHS (Voight *et al.* 2006), nSL (Ferrer-Admetlla *et al.* 2014), XP-EHH (Sabeti *et al.* 2007), and XP-nSL (Szpiech *et al.* 2021)—have become commonplace in evolutionary genomics studies to identify recent and ongoing positive selection in populations (e.g. Colonna *et al.* 2014, Zoledziewska *et al.* 2015, Nédélec *et al.* 2016, Crawford *et al.* 2017, Meier *et al.* 2018, Lu *et al.* 2019, Zhang *et al.* 2020, Salmón *et al.* 2021). When an adaptive allele sweeps through a population, it leaves a characteristic pattern of long high-frequency haplotypes and low genetic diversity in the vicinity of the allele. These statistics aim to capture these signals by summarizing the decay of haplotype homozygosity as a function of distance from a putatively selected region, either within a single population (iHS and nSL) or between two populations (XP-EHH and XP-nSL).

These haplotype-based statistics are powerful for detecting recent positive selection (Colonna *et al.* 2014, Zoledziewska *et al.* 2015, Nédélec *et al.* 2016, Crawford *et al.* 2017, Meier *et al.* 2018, Lu *et al.* 2019, Zhang *et al.* 2020, Salmón *et al.* 2021), and the two-population versions can even out-perform pairwise *F<sub>st</sub>* scans on a large swath of the parameter space (Szpiech *et al.* 2021). Furthermore, haplotype-based methods have also been shown to be robust to background selection (Fagny *et al.* 2014, Schrider 2020). However, each of these statistics presumes that haplotype phase is known or well-estimated.

As the generation of genomic sequencing data for non-model organisms is becoming routine (Ellegren 2014), there are many great opportunities for studying recent adaptation across the tree of life (e.g. Campagna and Toews 2022). However, often these organisms/populations do not have a well-characterized demographic history or recombination rate

map, two pieces of information which are important inputs for statistical phasing methods (Delaneau *et al.* 2013, Browning *et al.* 2021).

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

### 2 Materials and methods

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

#### 2.1 Extended haplotype homozygosity

In a sample of *n* diploid individuals, let *C* denote the set of all possible genotypes at locus *x*<sub>0</sub>. For multi-locus genotypes, *C* := {0, 1, 2}, representing the total counts of a derived allele. Let *C*(*x*<sub>*i*</sub>) be the set of all unique haplotypes extending from site *x*<sub>0</sub> to site *x*<sub>*i*</sub> either upstream or downstream

Received: 18 August 2023; Revised: 26 December 2023; Editorial Decision: 29 December 2023; Accepted: 3 January 2024

© The Author(s) 2024. Published by Oxford University Press.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

of  $x_0$ . If  $x_1$  is a site immediately adjacent to  $x_0$ , then  $\mathcal{C}(x_1) := \{00, 01, 02, 10, 11, 12, 20, 21, 22\}$ , representing all possible two-site multi-locus genotypes. We can then compute the extended haplotype homozygosity (EHH) of a set of multi-locus genotypes as

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

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

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

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

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

$$\text{cEHH}_c(x_i) = \sum_{h \in \mathcal{C}(x_i) \setminus \mathcal{H}_c(x_i)} \binom{n_h}{2} / \binom{n_{c'}}{2},$$

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

## 2.2 iHS and nSL

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

$$\begin{aligned} \text{iHH}_c &= \sum_{i=1}^{|\mathcal{D}|} \frac{1}{2} (\text{EHH}_c(x_{i-1}) + \text{EHH}_c(x_i)) d(x_{i-1}, x_i) \\ &+ \sum_{i=1}^{|\mathcal{U}|} \frac{1}{2} (\text{EHH}_c(x_{i-1}) + \text{EHH}_c(x_i)) d(x_{i-1}, x_i), \end{aligned}$$

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

$$\begin{aligned} \text{ciHH}_c &= \sum_{i=1}^{|\mathcal{D}|} \frac{1}{2} (\text{cEHH}_c(x_{i-1}) + \text{cEHH}_c(x_i)) d(x_{i-1}, x_i) \\ &+ \sum_{i=1}^{|\mathcal{U}|} \frac{1}{2} (\text{cEHH}_c(x_{i-1}) + \text{cEHH}_c(x_i)) d(x_{i-1}, x_i). \end{aligned}$$

The (unstandardized) unphased iHS is then calculated as

$$\text{iHS} = \begin{cases} \text{iHS}_2, & \text{if } \text{iHS}_2 > \text{iHS}_0 \\ -\text{iHS}_0, & \text{otherwise} \end{cases},$$

where  $\text{iHS}_2 = \log_{10}(\text{iHH}_2/\text{ciHH}_2)$  and  $\text{iHS}_0 = \log_{10}(\text{iHH}_0/\text{ciHH}_0)$ . Conceptually, this is nearly identical to the phased version of iHS, where the log ratio of the integrated haplotype homozygosity is computed between all haplotypes carrying the ancestral allele at the core locus versus all haplotypes carrying the derived allele at the core locus. In this case, however, we compare the iHH of the haplotypes containing homozygous genotypes of one allele at the core locus to the iHH of the haplotypes containing all other genotypes at the core locus. Doing this for both homozygous derived and homozygous ancestral haplotypes separately, we then choose the most extreme value. We assign a positive sign for long low-diversity haplotypes containing the derived homozygous genotype at the core locus, and we assign a negative sign for long low-diversity haplotypes containing the ancestral homozygous genotype at the core locus. Unstandardized iHS scores are then normalized in frequency bins, as previously described (Voight *et al.* 2006, Ferrer-Admetlla *et al.* 2014). Unstandardized unphased nSL is computed similarly with the appropriate distance measure [see Ferrer-Admetlla *et al.* (2014) where they show that nSL can be reformulated as iHS with a different distance measure]. Large positive scores indicate long high-frequency haplotypes with a homozygous derived core genotype, and large negative scores indicate long high-frequency haplotypes with a homozygous ancestral core genotype. Clusters of extreme scores in both directions indicate evidence for a sweep.

## 2.3 XP-EHH and XP-nSL

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

$$\begin{aligned} \text{iHH}_P &= \sum_{i=1}^{|\mathcal{D}|} \frac{1}{2} (\text{EHH}(x_{i-1}) + \text{EHH}(x_i)) d(x_{i-1}, x_i) \\ &+ \sum_{i=1}^{|\mathcal{U}|} \frac{1}{2} (\text{EHH}(x_{i-1}) + \text{EHH}(x_i)) d(x_{i-1}, x_i), \end{aligned}$$

where the distance measure is given as centimorgans or basepairs for XP-EHH (Sabeti *et al.* 2007) and number of sites observed for XP-nSL (Szpiech *et al.* 2021). The XP statistics between population  $A$  and  $B$  are then computed as  $\text{XP} = \log_{10}(\text{iHH}_A/\text{iHH}_B)$  and are normalized genome wide. Large positive scores indicate long high-frequency haplotypes in population  $A$ , and large negative scores indicate long high-frequency haplotypes in population  $B$ . Clusters of extreme scores in one direction indicate evidence for a sweep in that population.

## 2.4 Simulations

We evaluate the performance of the phased and unphased versions of iHS, nSL, XP-EHH, and XP-nSL under a generic two-population divergence model using the coalescent simulation program discoal (Kern and Schrider 2016). We explore five versions of this generic model and name them Demo 1 through Demo 5 (Supplementary Table S1). Let  $N_0$  and  $N_1$  be the effective population sizes of Population 0 and Population 1 after the split from their ancestral population (of size  $N_A$ ). For Demo 1, we keep a constant population size post-split

and let  $N_0 = N_1 = 10\,000$ . For Demo 2, we keep a constant population size post-split and let  $N_0 = 2N_1 = 10\,000$ . For Demo 3, we keep a constant population size post-split and let  $2N_0 = N_1 = 10\,000$ . For Demo 4, we initially set  $N_0 = N_1 = 10\,000$  and let  $N_0$  grow stepwise exponentially every 50 generations starting at 2000 generations ago until  $N_0 = 5N_1 = 50\,000$ . For Demo 5, we initially set  $N_0 = N_1 = 10\,000$  and let  $N_1$  grow stepwise exponentially every 50 generations starting at 2000 generations ago until  $5N_0 = N_1 = 50\,000$ .

For each demographic history we vary the population divergence time  $t_d \in \{2000, 4000, 8000\}$  generations ago. For non-neutral simulations, we simulate a sweep in Population 0 in the middle of the simulated region across a range of selection coefficients  $s \in \{0.005, 0.01, 0.02\}$ . We vary the frequency at which the adaptive allele starts sweeping as  $e \in \{0, 0.01, 0.02, 0.05, 0.10\}$ , where  $e = 0$  indicates a hard sweep and  $e > 0$  indicates a soft sweep, and we also vary the frequency of the selected allele at time of sampling  $f \in \{0.7, 0.8, 0.9, 1.0\}$  as well as  $g \in \{50, 100\}$  representing fixation of the sweeping allele  $g$  generations ago. For all simulations we set the genome length to be  $L = 500\,000$  basepairs, the ancestral effective population size to be  $N_A = 10\,000$ , the per site per generation mutation rate at  $\mu = 2.35 \times 10^{-8}$ , and the per site per generation recombination rate at  $r = 1.2 \times 10^{-8}$ . For neutral simulations, we simulate 1000 replicates for each parameter set, and for non-neutral simulations we simulate 100 replicates for each parameter set. We sample  $2n \in \{200, 100, 40, 20\}$  haplotypes, randomly paired together to form  $n \in \{100, 50, 20, 10\}$  diploid individuals, from each population for analysis. These datasets represent the case where phase is known perfectly. We also create a set of “unphased” datasets from these phased datasets by swapping the alleles of each heterozygote to the opposing haplotype with probability 0.5.

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

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

## 2.5 Power and false positive rate

Here we evaluate the power and false positive rate for the unphased version of iHS, nSL, XP-EHH, and XP-nSL. For comparison, we also compute the power for the original phased versions of these statistics in two different ways. We compute the phased statistics for a set of simulated datasets where perfect phase is known, and we compute them again for a set of simulated datasets where we destroy phase information (see Section 2.4). As the unphased statistics collapse genotypes into derived allele counts, there is no functional difference between these two datasets for these statistics. We compute power in the same way for each statistic regardless of underlying dataset analyzed as described below.

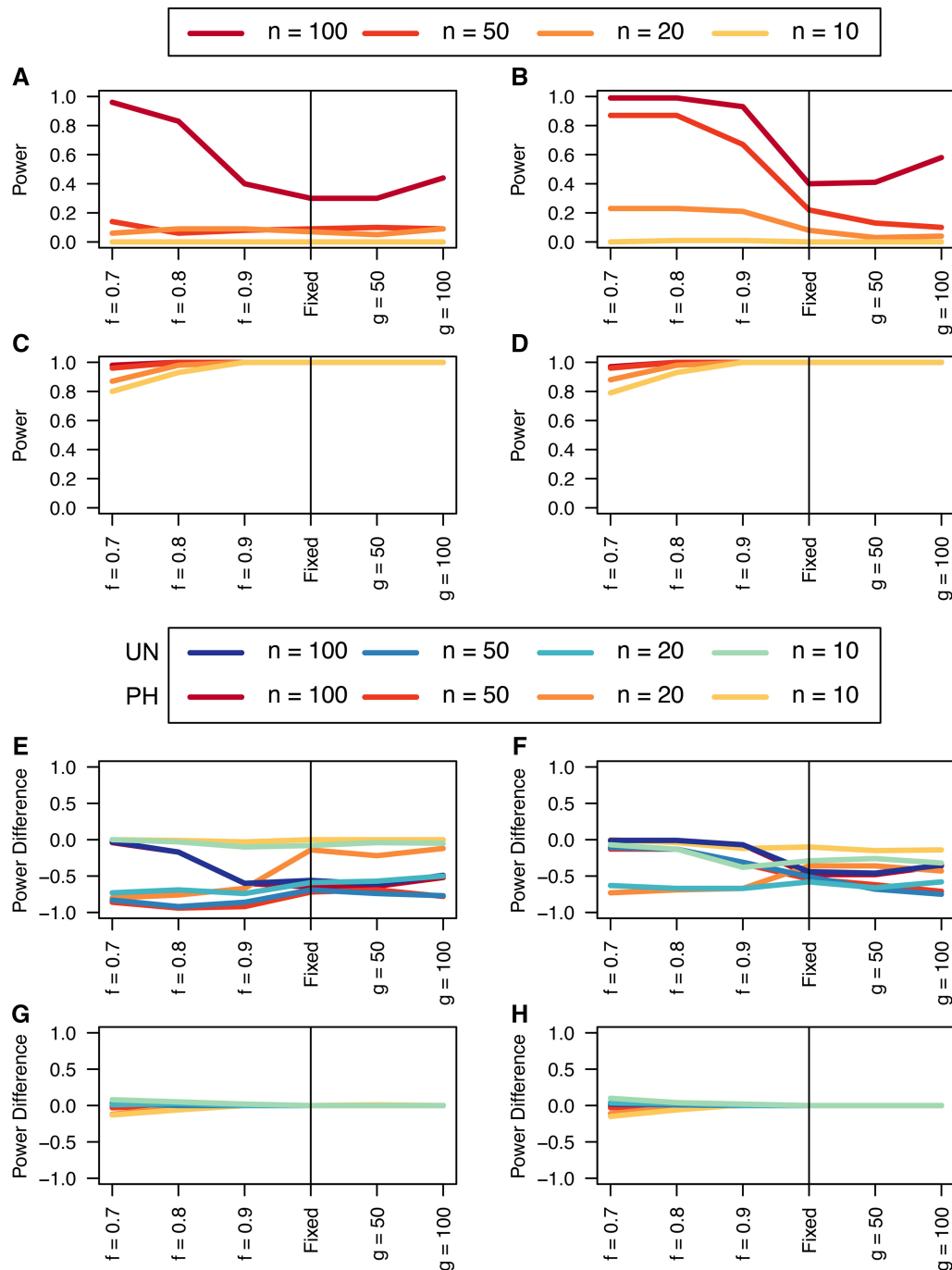
To compute power for iHS and nSL, we follow the approach of Voight *et al.* (2006). For these statistics, each non-neutral replicate is individually normalized jointly with all neutral replicates with matching demographic history in 1% allele frequency bins. Because extreme values of the statistic are likely to be clustered along the genome (Voight *et al.* 2006), we then compute the proportion of extreme scores ( $|iHS| > 2$  or  $|nSL| > 2$ ) within 100kbp non-overlapping windows. We then bin these windows into 10 quantile bins based on the number of scores observed in each window and call the top 1% of these windows as putatively under selection. We calculate the proportion of non-neutral replicates that fall in this top 1% as the power. To compute the false positive rate, we compute the proportion of neutral simulations that fall within the top 1%.

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

## 3 Results

We find that the unphased versions of iHS and nSL generally have good power at large sample sizes (Fig. 1A and B, Supplementary Figs S1, S7, and S8) to detect selection prior to fixation of the allele, with nSL generally outperforming iHS. In smaller populations (Supplementary Fig. S1C and D), power does suffer relative to larger populations (Supplementary Fig. S1A, B, E, and F). We note that these statistics struggle to identify soft sweeps when the population is undergoing exponential growth (Supplementary Fig. S1E and F). Each of these statistics also have low false positive rates hovering around 1% (Supplementary Tables S2–S5). These single-population statistics only perform well for relatively large samples (Fig. 1A and B, and Supplementary Figs S19, S25, S26, S31, S32, S37, S43, S44, S55, S61, and S62).

Similarly, we find that the unphased versions of XP-EHH and XP-nSL have good power as well even for relatively low sample sizes (Fig. 1C, D, G, and H and Supplementary Figs S2, S3, S9–S12, S20, S21, S27–30, S38, S39, S45–48, S56, S57, S63–S66). When the sweep takes place in the smaller of the two populations (Supplementary Figs S2C, S2D, S20C, S20D, S38C, S38D, S56C, and S56D), we see a similar decrease in power, likely related to the lower efficiency of selection in small populations. When one population is undergoing exponential growth (Supplementary Figs S3, S21, S39, and S57) performance is generally quite good, likely the result of a larger effective selection coefficient in large populations. These two-population statistics generally outperform their single-population counterparts, especially at small diploid sample sizes and for sweeps that have reached fixation recently. Each of these statistics also have low false positive rates hovering around 1% (Supplementary Tables S2–S5).



**Figure 1.** Unphased power. Power curves for unphased implementations of iHS (A), nSL (B), XP-EHH (C), and XP-nSL (D), and power difference between unphased implementations of iHS (E), nSL (F), XP-EHH (G), and XP-nSL (H) and phased implementations. Blue curves represent the power difference between the unphased and phased statistics when applied to unphased data (UN). Red curves represent the power difference between the unphased and phased statistics when applied to perfectly phased data (PH). Values greater than 0 indicate the unphased statistic had higher power. All panels represent analyses with demographic history Demo 1 and  $n = 100, 50, 20$ , or 10 diploid samples. For these plots the selection coefficient is set at  $s = 0.01$ , the frequency at which selection began is set at  $e = 0$  (i.e. a hard sweep), and the divergence time in generations is set at  $t_d = 2000$ .  $f$  is the frequency of the adaptive allele at time of sampling,  $g$  is the number of generations at time of sampling since fixation.

Next, we turn to comparing the performance of these unphased statistics to their phased counterparts when they are used to analyze either phased data or unphased data. In Fig. 1E–H and Supplementary Figs S4–S6, S13–S18, S22–S24, S31–S36, S40–S42, S49–S54, S58–S60, and S67–S72, we plot the difference in power between the unphased statistics and the phased counterpart applied to data with phase

known (red lines) or phase scrambled (blue lines). Where these lines are greater than or equal to 0 indicates that the unphased statistic performed as well as or better than the phased counterpart.

We find that iHS tends to underperform the traditional phased implementations, but nSL tends to perform as well as the phased versions (Fig. 1E and F and Supplementary Figs



S4, S13, S14, S22, S31, S32, S40, S49, S50, S58, S67, and S68). Although we note noticeable drops in unphased nSL power for softer sweeps in exponential growth scenarios (Supplementary Figs S4F, S13F, S14F, S22F, S31F, S32F, S40F, S49F, S50F, S58F, S67F, and S68F) and for sweeps near completion in small population sizes (Supplementary Figs S4E, S13E, S14E, S22E, S31E, S32E, S40E, S49E, S50E, S58E, S67E, and S68E).

When comparing the unphased versions of XP-EHH and XP-nSL, we find that they consistently perform as well or better than their phased counterparts (Fig. 1G and H and Supplementary Figs S5, S6, S17, S18, S23, S24, S35, S36, S41, S42, S53, S54, S59, S60, S71, and S72), except in limited circumstances where phase is known, and the sweep is fairly young (sweeping allele at 0.7 frequency) or the divergence time is further in the past.

## 4 Discussion

We introduce multi-locus genotype versions of four popular haplotype-based selection statistics—iHS (Voight *et al.* 2006), nSL (Ferrer-Admetlla *et al.* 2014), XP-EHH (Sabeti *et al.* 2007), and XP-nSL (Szpiech *et al.* 2021)—that can be used when the phase of genotypes is unknown. Although phase would seem to be a critically important component of any haplotype-based method for detecting selection, here we show that, by collapsing haplotypes into derived allele counts (thus erasing phase information), we can achieve similar power to using this information. We observed that single-population statistics such as iHS and nSL require relatively large diploid sample sizes ( $n \geq 100$  for iHS,  $n \geq 50$  for nSL), but the two-population statistics XP-EHH and XP-nSL perform well even for diploid sample sizes down to  $n = 10$  per population. This follows other work that has shown similar patterns with other haplotype-based statistics for detecting selection (Harris *et al.* 2018, Harris and DeGiorgio 2020, DeGiorgio and Szpiech 2022, Klassmann and Gautier 2022). Importantly, this approach now opens up the application of several popular haplotype-based selection statistics (based on extended haplotype homozygosity) to more species where phase information is challenging to know or infer.

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

## Acknowledgements

Computations for this research were performed using the Pennsylvania State University's Institute for Computational Data Sciences' Roar supercomputer.

## Supplementary data

Supplementary data are available at *Bioinformatics* online.

## Conflict of interest

None declared.

## Funding

This work was supported by the National Institute of General Medical Sciences of the National Institutes of Health [award number R35GM146926]; and start-up funds from the Pennsylvania State University's Department of Biology.

## Data availability

The data underlying this article are available in the article and in its online [supplementary material](#).

## References

- Browning BL, Tian X, Zhou Y *et al.* Fast two-stage phasing of large-scale sequence data. *Am J Hum Genet* 2021;108:1880–90.
- Campagna L, Toews DPL. The genomics of adaptation in birds. *Curr Biol* 2022;32:R1173–86.
- Colonna V, Ayub Q, Chen Y *et al.*; 1000 Genomes Project Consortium. Human genomic regions with exceptionally high levels of population differentiation identified from 911 whole-genome sequences. *Genome Biol* 2014;15:R88.
- Crawford NG, Kelly DE, Hansen MTEB *et al.*; NISC Comparative Sequencing Program. Loci associated with skin pigmentation identified in African populations. *Science* 2017;358:eaan8433.
- DeGiorgio M, Szpiech ZA. A spatially aware likelihood test to detect sweeps from haplotype distributions. *PLoS Genet* 2022;18:e1010134.
- Delaneau O, Zagury JF, Marchini J. Improved whole-chromosome phasing for disease and population genetic studies. *Nat Methods* 2013;10:5–6.
- Ellegren H. Genome sequencing and population genomics in non-model organisms. *Trends Ecol Evol* 2014;29:51–63.
- Fagny M, Patin E, Enard D *et al.* Exploring the occurrence of classic selective sweeps in humans using whole-genome sequencing data sets. *Mol Biol Evol* 2014;31:1850–68.
- Ferrer-Admetlla A, Liang M, Korneliussen T *et al.* On detecting incomplete soft or hard selective sweeps using haplotype structure. *Mol Biol Evol* 2014;31:1275–91.
- Harris AM, DeGiorgio M. A likelihood approach for uncovering selective sweep signatures from haplotype data. *Mol Biol Evol* 2020;37:3023–46.
- Harris AM, Garud NR, DeGiorgio M. Detection and classification of hard and soft sweeps from unphased genotypes by multilocus genotype identity. *Genetics* 2018;210:1429–52.
- Kern AD, Schrider DR. Discoal: flexible coalescent simulations with selection. *Bioinformatics* 2016;32:3839–41.
- Klassmann A, Gautier M. Detecting selection using extended haplotype homozygosity (EHH)-based statistics in unphased or unpolarized data. *PLoS One* 2022;17:e0262024.
- Lu K, Wei L, Li X *et al.* Whole-genome resequencing reveals *Brassica napus* origin and genetic loci involved in its improvement. *Nat Commun* 2019;10:1154.
- Meier JI, Marques DA, Wagner CE *et al.* Genomics of parallel ecological speciation in Lake Victoria Cichlids. *Mol Biol Evol* 2018;35:1489–506.
- Nédélec Y, Sanz J, Baharian G *et al.* Genetic ancestry and natural selection drive population differences in immune responses to pathogens. *Cell* 2016;167:657–69.e621.
- Sabeti PC, Varilly P, Fry B *et al.*; International HapMap Consortium. Genome-wide detection and characterization of positive selection in human populations. *Nature* 2007;449:913–8.
- Salmón P, Jacobs A, Ahrén D *et al.* Continent-wide genomic signatures of adaptation to urbanisation in a songbird across Europe. *Nat Commun* 2021;12:2983.
- Schrider DR. Background selection does not mimic the patterns of genetic diversity produced by selective sweeps. *Genetics* 2020;216:499–519.

- Szpiech ZA, Hernandez RD. selscan: an efficient multithreaded program to perform EHH-based scans for positive selection. *Mol Biol Evol* 2014;**31**:2824–7.
- Szpiech ZA, Novak TE, Bailey NP *et al.* Application of a novel haplotype-based scan for local adaptation to study high-altitude adaptation in rhesus macaques. *Evol Lett* 2021;**5**: 408–21.
- Voight BF, Kudaravalli S, Wen X *et al.* A map of recent positive selection in the human genome. *PLoS Biol* 2006;**4**:e72.
- Zhang S-J, Wang G-D, Ma P *et al.* Genomic regions under selection in the feralization of the dingoes. *Nat Commun* 2020;**11**:671.
- Zoledziewska M, Sidore C, Chiang CWK *et al.*; Understanding Society Scientific Group. Height-reducing variants and selection for short stature in Sardinia. *Nat Genet* 2015;**47**:1352–6.