

Detection and classification of hard and soft sweeps from unphased genotypes by multilocus genotype identity

Alexandre M. Harris^{1,2}, Nandita R. Garud³, Michael DeGiorgio^{1,4,5,*}

August 14, 2018

¹ *Department of Biology, Pennsylvania State University, University Park, PA 16802, USA*

² *Program in Molecular, Cellular, and Integrative Biosciences at the Huck Institutes of the Life Sciences,*

³ *Pennsylvania State University, University Park, PA 16802, USA*

⁴ *Gladstone Institute, University of California, San Francisco, CA, 94158, USA*

⁵ *Department of Statistics, Pennsylvania State University, University Park, PA 16802, USA*

⁶ *Institute for CyberScience, Pennsylvania State University, University Park, PA 16802, USA*

⁷ *Corresponding author: mxd60@psu.edu

To be appreciated in the review :

- inclusion of admixture scenarios
- inference of k

To be criticised :

- connection to PGT unfortunate, esp. in case of soft sweep; perhaps not very illuminating / conclusive

⁸ Keywords: Expected haplotype homozygosity, multilocus genotype, positive selection, hard sweep, soft sweep

¹⁰ Running title: Detecting hard and soft sweeps

Abstract

Positive natural selection can lead to a decrease in genomic diversity at the selected site and at linked sites, producing a characteristic signature of elevated expected haplotype homozygosity. These selective sweeps can be hard or soft. In the case of a hard selective sweep, a single adaptive haplotype rises to high population frequency, whereas **multiple adaptive haplotypes sweep through the population simultaneously in a soft sweep**, producing distinct patterns of genetic variation in the vicinity of the selected site. Measures of expected haplotype homozygosity have previously been used to detect sweeps in multiple study systems. However, these methods are formulated for phased haplotype data, **typically** unavailable for nonmodel organisms, and may have reduced power to detect soft sweeps due to their increased genetic diversity relative to hard sweeps. To address these limitations, we applied the H12 and H2/H1 statistics of Garud et al. [2015] to unphased multilocus genotypes, denoting them as G12 and G2/G1. G12 (**and** the more direct expected homozygosity analogue to H12, denoted G123) has comparable power to H12 for detecting both hard and soft sweeps. G2/G1 can be used to classify hard and soft sweeps analogously to H2/H1, conditional on a genomic region having high G12 or G123 values. The reason for this power is that under random mating, the most frequent haplotypes will yield the most frequent multilocus genotypes. Simulations based on **parameters compatible with our recent understanding of human demographic history** suggest that **expected homozygosity** methods are best suited for detecting recent sweeps, and increase in power under recent population expansions. Finally, we find candidates for selective sweeps within the 1000 Genomes CEU, YRI, GIH, and CHB populations, which corroborate and complement existing studies.

1 Introduction

2 Positive natural selection is the process by which an advantageous genetic variant rises to high frequency
3 in a population, thereby reducing site diversity and creating a tract of elevated expected homozygosity and
4 linkage disequilibrium (LD) surrounding that variant [Sabeti et al., 2002]. As beneficial alleles increase to
5 high frequency in a population, the signature of a selective sweep emerges, which we can characterize from the
6 number of haplotypes involved in the sweep [Maynard Smith and Haigh, 1974, Schweinsberg and Durrett,
7 2005, Hermisson and Pennings, 2017]. A hard sweep is an event in which a single haplotype harboring
8 a selectively advantageous allele rises in frequency, while in a soft sweep, multiple haplotypes harboring
9 advantageous mutations can rise in frequency simultaneously. Thus, selective sweeps represent a broad and
10 non-homogenous spectrum of genomic signatures. A selective event that persists until the beneficial allele
11 reaches fixation is a *complete* sweep, while a *partial* sweep is one in which the selected allele does not reach
12 fixation. Consequently, expected haplotype homozygosity surrounding the selected site is greatest once the
13 selected allele has fixed and before recombination and mutation break up local LD [Przeworski, 2002].

14 Two modes of soft sweeps have been proposed across three seminal papers, consisting of sweeps from
15 standing genetic variation that becomes beneficial in a changing environment, or new recurrent *de novo* adap-
16 tive mutations [Hermisson and Pennings, 2005, Pennings and Hermisson, 2006a,b], and these can be complete
17 and partial as well. Unlike hard sweeps, where haplotypic diversity is decreased, in a soft sweep, haplotypic
18 diversity remains, since multiple haplotypes carrying the adaptive allele rise to high frequency. [Przeworski
19 et al., 2005, Berg and Coop, 2015]. Patterns of diversity surrounding the selected site begin to resemble those
20 expected under neutrality as the number of unique haplotypic backgrounds carrying the beneficial allele (the
21 softness of the sweep) increases, potentially obscuring the presence of the sweep. Accordingly, the effect of
22 a soft sweep may be unnoticeable, even if the selected allele has reached fixation.

23 Popular modern methods for identifying recent selective sweeps from haplotype data identify distortions
24 in the haplotype structure following a sweep, making use of either the signature of elevated LD or reduced
25 haplotypic diversity surrounding the site of selection. Methods in the former category [Kelly, 1997, Kim
26 and Nielsen, 2004, Pavlidis et al., 2010] can detect both hard and soft sweeps, as neighboring neutral vari-
27 ants hitchhike to high frequency under either scenario. Indeed, LD-based methods may have an increased
28 sensitivity to softer sweeps [Pennings and Hermisson, 2006b], especially relative to methods that do not use
29 haplotype data, such as composite likelihood approaches [Kim and Stephan, 2002, Nielsen et al., 2005, Chen
30 et al., 2010, Vy and Kim, 2015, Racimo, 2016]. Haplotype homozygosity-based methods include iHS [Voight
31 et al., 2006], its extension, nS_L [Ferrer-Admetlla et al., 2014], and H-scan [Schlamp et al., 2016]. These
32 approaches identify a site under selection from the presence of a high-frequency haplotype. Additionally,

1 Chen et al. [2015] developed a hidden Markov model-based approach that similarly identifies sites under
2 selection from the surrounding long, high-frequency haplotype.

3 While the aforementioned methods are powerful tools for identifying **selective sweeps** in the genome,
4 they lack the ability to distinguish between hard and soft sweeps. It is this concern that Garud et al. [2015]
5 address with the statistics H12 and H2/H1. H12, a haplotype homozygosity-based method, identifies selective
6 sweeps from elevated expected haplotype homozygosity surrounding the selected site. It is computed as **the**
7 expected haplotype homozygosity, with the frequencies of the two most frequent haplotypes pooled into a
8 single frequency:

$$H12 = (p_1 + p_2)^2 + \sum_{i=3}^I p_i^2, \quad (1)$$

9 where there are I distinct haplotypes in the population, and p_i is the frequency of the i th most frequent
10 haplotype, with $p_1 \geq p_2 \geq \dots \geq p_I$. Pooling **the two largest haplotype frequencies** provides little additional
11 power to detect hard sweeps relative to H1, the standard measure of expected haplotype homozygosity,
12 where $H1 = \sum_{i=1}^I p_i^2$ (Figure 1A, left panel). However, pooling provides more power to detect soft sweeps,
13 in which at least two haplotypes rise to high frequency, and the distortion of their joint frequency produces
14 an elevated expected haplotype homozygosity consistent with a sweep (Figure 1A, right panel).

15 In conjunction with an elevated value of H12, the ratio H2/H1 serves as a measure of sweep softness,
16 and is not meaningful on its own. H2 is **the** expected haplotype homozygosity, omitting the most frequent
17 haplotype, computed as $H2 = H1 - p_1^2$, and is larger for softer sweeps. In the case of a soft sweep, the
18 frequencies of the first- and second-most frequent haplotypes are both large, and omitting the most frequent
19 haplotype still yields a frequency distribution in which one haplotype predominates. Under a hard sweep,
20 the second through I th haplotypes are likely to be at low frequency and closer in value, such that their
21 expected homozygosity is small. Thus, while $H2 < H1$ in all cases, the value of H2 is closer to that of H1
22 under a soft sweep.

23 To leverage the power of H12 and H2/H1 to detect sweeps in nonmodel organisms, for which phased
24 haplotype data are often unavailable, we extend the application of these statistics to unphased multilocus
25 genotype (MLG) data as G12 and G2/G1. MLGs are **single strings** representing a diploid individual's allelic
26 state at each site as homozygous for the reference allele, homozygous for the alternate allele, or heterozygous.
27 Similarly to H12, we define G12 as

$$G12 = (q_1 + q_2)^2 + \sum_{j=3}^J q_j^2, \quad (2)$$

28 where there are J distinct unphased MLGs in the population, and q_j is the frequency of the j th most
29 frequent MLG, with $q_1 \geq q_2 \geq \dots \geq q_J$. As with **haplotypes**, pooling the most frequent MLGs only provides
30 marginally more resolution to detect hard sweeps, as only a single predominant unphased MLG is expected

under random mating (Figure 1B, left panel). However, because the input data for G12 and G2/G1 are unphased MLGs, we define another statistic that is uniquely meaningful in this context. The presence of multiple unique frequent haplotypes under a soft sweep implies not only that the frequency of individuals homozygous for these haplotypes will be elevated, but also that the frequencies of their heterozygotes will be elevated. When haplotypes X and Y both exist at high frequency, diploid individuals of type XX , YY , and XY will also exist at high frequency, assuming individuals mate randomly within the population (Figure 1B, right panel). Therefore, we can define a statistic truly analogous to H12 for unphased MLG data, G123.

This statistic is calculated as

$$G123 = (q_1 + q_2 + q_3)^2 + \sum_{j=4}^J q_j^2. \quad (3)$$

We note, however, that with this approach we do not explicitly enforce a constraint on the presence of particular high-frequency MLGs in the sample. That is, we only assume that the presence of one or more high-frequency MLGs implies a sweep, even if any one of the XX , YY , or XY MLGs is absent.

We show through simulation and empirical application that the statistics G12 and G123, in conjunction with the ratio G2/G1, both maintain the power of H12 to detect and classify sweeps, without requiring phased haplotype input data. Furthermore, as a closer analogue to H12, the use of G123 with G2/G1 more closely maintains the classification ability of H12 with H2/H1 than does G12. Generally, we find that the selective events visible with H12 in phased haplotype data are visible to G12 and G123 in unphased MLG data, with trends in power and genomic signature of the applications remaining consistent with one another. Accordingly, we recover well-documented sweep signatures at *LCT* and *SLC24A5* in individuals with European ancestry [Bersaglieri et al., 2004, Sabeti et al., 2007, Gerbault et al., 2009], with the latter also detected in South Asian individuals [Coop et al., 2009, Mallick et al., 2013], as well as the region linked to *EDAR* in East Asian populations [Fujimoto et al., 2007, Bryk et al., 2008, Pickrell et al., 2009], and *SYT1* in African individuals [Voight et al., 2006]. In addition, we identify novel candidates *RGS18* in African individuals, *P4HA1* in South Asian individuals, and *FMNL3* in East Asian individuals.

Results

To detect selective sweeps, we must have power to identify loci with elevated haplotype homozygosity relative to expectations under neutral demographic scenarios. We compared the power of the MLG-based methods G12 and G123 to that of the haplotype-based methods H12 and H123 [Garud et al., 2015], at the 1% false positive rate (FPR) obtained from simulations under neutral demographic models (see *Materials and methods*). We performed simulations under population-genetic parameters inferred for human data [Takahata et al., 1995, Nachman and Crowell, 2000, Payseur and Nachman, 2000] with the forward-time simulator

1 SLiM 2 [Haller and Messer, 2017]. Because SLiM outputs paired phased haplotypes for each diploid indi-
2 vidual, we manually merged each individual's haplotypes to apply the MLG-based methods. Our simulated
3 replicates included scenarios of selective neutrality, hard sweeps, and soft sweeps. We evaluated methods
4 across simulations of constant demographic history, as well as realistic human models of bottleneck and
5 expansion [Lohmueller et al., 2009] (Figure 2). We then use an approximate Bayesian computation (ABC)
6 approach to evaluate the ability of the MLG-based methods with G2/G1, and the haplotype-based methods
7 with H2/H1, to differentiate between hard and soft sweeps. Finally, we evaluated empirical data from the
8 1000 Genomes Project [Auton et al., 2015], manually merging each study individual's phased haplotypes into
9 MLGs to observe the effect of phasing on our ability to detect selective events. See *Materials and methods*
10 for a detailed explanation of experiments.

11 Using G12 and G123 to detect sweeps

12 We demonstrate the range of sensitivity of G12 and G123 relative to H12 and H123 for selective sweeps
13 occurring at time points between 400 and 4,000 generations before the time of sampling. We evaluated G123
14 to determine whether it is a more direct analogue of H12 as we expected, while our application of H123
15 follows from the work of Garud et al. [2015], which suggested that H123 **yields little difference in power to**
16 **detect sweeps relative to H12** for given sample and window size parameters. In the following experiments,
17 we simulated 100 kilobase (kb) chromosomes carrying a selected allele at their center (sweep simulations), or
18 carrying no selected allele for neutrality, performing 10^3 replicates for each scenario with sample size $n = 100$
19 diploid individuals.

20 For each series of simulations, we **detected sweeps** using a sliding window of size 40 kb shifting by 4 kb
21 increments across the chromosome. We selected this window size to ensure that the effect of short-range
22 LD would not inflate the values of our statistics (Figure S1). This additionally matched the window size
23 we selected for analysis of empirical data in non-African populations (see *Analysis of empirical data for*
24 *signatures of sweeps*). **According** to theoretical expectations [Gillespie, 2004, Garud et al., 2015, Hermisson
25 and Pennings, 2017], a window of size 40 kb under our simulated parameters is sensitive to sweeps with
26 selection strength $s \geq 0.004$ (see *Materials and methods*). **Additionally**, although we used a nucleotide-
27 delimited window in our analysis, one can also fix the number of single-nucleotide polymorphisms (SNPs)
28 included in each window (SNP-delimited window), though this somewhat changes the properties of the
29 methods (see *Discussion*). A SNP-delimited window corresponding to approximately 40 kb for our simulated
30 data **contains** on average 235 SNPs under neutrality. To supplement experiments measuring **the power of**
31 **each method**, we also assessed the **genomic** distribution of G12 and G123 values to characterize their patterns
32 under sweep scenarios.

¹ **Tests for detection of hard sweeps**

² Methods that detect selective sweeps typically focus on the signature of hard sweeps, though many can
³ detect soft sweeps as well. Accordingly, we began by measuring the ability of G12, G123, H12, and H123
⁴ to detect both partial and complete hard sweeps, under scenarios in which a single haplotype acquires a
⁵ selected mutation and rises in frequency. We examined selection start times (*t*) of 400, 1,000, 2,000, and
⁶ 4,000 generations before the time of sampling. These values of *t* span the time periods of various sweeps in
⁷ human history [Przeworski, 2002, Sabeti et al., 2007, Beleza et al., 2012, Jones et al., 2013, Clemente et al.,
⁸ 2014, Fagny et al., 2014]. For each *t*, we simulated hard sweeps under the aforementioned parameters to
⁹ sweep frequencies (*f*) between 0.1 and 1 for the selected allele (Figures 3 and S2). Sweeps to smaller *f* have
¹⁰ a smaller effect on the surrounding expected haplotype homozygosity and are more difficult to detect. We
¹¹ performed hard sweep simulations for a large selection coefficient of *s* = 0.1 and a more moderate selection
¹² coefficient of *s* = 0.01.

¹³ The values of *t* and *f* both impact the ability of methods to identify hard sweeps (Figure 3). At the 1%
¹⁴ FPR, all methods are suited to the detection of more recent sweeps for simulated data, losing considerable
¹⁵ power to resolve hard sweep events occurring prior to 2,000 generations before sampling, and losing power
¹⁶ entirely for hard sweeps occurring prior to 4,000 generations before sampling. For selection within 2,000
¹⁷ generations of sampling, trends in the power of the MLG-based methods resemble those of the haplotype-
¹⁸ based methods, with the power of the MLG-based methods either matching or approaching that of the
¹⁹ haplotype-based methods for *s* = 0.1 (Figures 3A and S2A), and following similar trends in power for
²⁰ *s* = 0.01 (though with slightly reduced power overall; Figures 3B and S2B), indicating that the two highest-
²¹ frequency MLGs and the two highest-frequency haplotypes have a similar ability to convey the signature of
²² a sweep.

²³ For data simulated under strong selection, *s* = 0.1 (Figure 3A), G12 and H12 achieve their maximum
²⁴ power for recent selective sweeps originating within the past 1,000 generations (with little to no power lost
²⁵ over this interval for sweeps to large *f*). This result is expected because sweeps with such a high selection
²⁶ coefficient quickly reach fixation, at which point mutation and recombination break down tracts of elevated
²⁷ expected homozygosity until the signal fully decays, obscuring more ancient events. For a given value of *s*,
²⁸ selective sweeps to larger values of *f* for the selected allele additionally produce a stronger signal because
²⁹ more diversity is ablated the longer a sweep lasts. Thus, G12 and H12 are best able to detect sweeps over
³⁰ recent time intervals, especially as the sweep goes to larger values of *f*. Strong hard sweeps additionally
³¹ create a peak in signal surrounding the site of selection that increases in magnitude with increasing duration
³² of a sweep. This signal is broad and extends across the one Mb interval that we modeled in Figure 3C. These
³³ patterns repeat for G123 and H123 (Figure S2A), yielding little difference in power between H12 and H123,

1 and no difference in power between G123 and G12 (along with a nearly-identical spatial signature **along the**
2 **chromosome**; Figure S2C).

3 At a smaller selection coefficient of $s = 0.01$ (Figure 3B), G12 and H12 have a **distinct range of sweep**
4 **detection** from $s = 0.1$. The reduced strength of selection **here** leads beneficial mutations to rise more slowly
5 in frequency than for stronger selection. Consequently, after 400 generations of selection, the distribution of
6 haplotype (and therefore MLG) frequencies has scarcely changed from neutrality, and G12 and H12 cannot
7 reliably detect the signal of a sweep. However, the powers of G12 and H12, as well as G123 and H123
8 (Figure S2B), are greatest for **a moderate sweep to $f \geq 0.9$ starting 2,000 generations prior to sampling. As**
9 **with stronger selection, pooling the three largest frequencies had little effect on power relative to pooling the**
10 **two largest frequencies (Figure S2)**. We **could not** detect adaptive mutations appearing more anciently than
11 2,000 generations before sampling, **indicating** that all methods lose power to detect sweeps for smaller values
12 of s , and that haplotype methods may outperform MLG methods for smaller values of s as well. Furthermore,
13 the range of time over which methods detect a sweep narrows and shifts to more ancient time periods with
14 decreasing s . Weaker selection nonetheless produces a signal peak distinct from the neutral background and
15 proportional in magnitude to **the value of f** (Figures 3D and S2D), though expected haplotype homozygosity,
16 and therefore expected MLG homozygosity, is reduced for moderate selection (compare vertical axes of
17 Figures 3C and D and of Figures S2C and D).

18 Tests for detection of sweeps on standing variation

19 We characterized the properties of G12, G123, H12, and H123 for simulated soft sweeps from selection on
20 standing genetic variation (SSV). We generated results analogous to those for hard sweeps: measures of
21 power **for each method**, and the chromosome-wide spatial distribution of the G12 and G123 signals. Across
22 identical times of selection (t) and selection coefficients (s) as for hard sweep simulations, we simulated SSV
23 scenarios by introducing the selected mutation on multiple haplotypes simultaneously. We evaluated method
24 ability to correctly **distinguish** sweeps on $k = 2, 4, 8, 16$, and 32 initially-selected **different** haplotypes from
25 **neutrality. One copy of the selected allele is guaranteed to remain in the population for the entire simulation,**
26 **but we do not condition on the number of sweeping haplotypes at the time of sampling.** Indeed, we do not
27 expect that for larger values of k , all haplotypes carrying the selected allele will remain at high frequency,
28 or remain at all by the time of sampling (Figure S4). For our scaled (see *Materials and methods*) simulated
29 population size of 500 diploids (unscaled 10^4 diploids), this corresponds to having the beneficial allele present
30 on 0.2 to 3.2% of haplotypes at the **onset of selection**. Our results for these tests mirror those for hard sweeps,
31 with stronger selection on fewer distinct haplotypes yielding the most readily detectable genomic signatures
32 (Figures 4 and S3).

1 SSV once again produces a signal of elevated MLG homozygosity for $s = 0.1$ that all methods most
2 readily detect if it is recent, and rapidly lose power to detect as t increases. G12 and H12 reliably detect
3 signals of SSV in simulated 100 kb chromosomes, retaining power for SSV on as many as $k \leq 16$ haplotypes
4 within the first 400 generations after the start of selection (Figure 4A). However, the relatively smaller
5 expected homozygosity under SSV leads the power of each method to decay more rapidly than under a hard
6 sweep. The levels of expected homozygosity produced under SSV are consequently smaller in magnitude than
7 those generated under hard sweeps, but unambiguously distinct from neutrality for at least one combination
8 of each tested k and t , with $k = 2$ most closely resembling a hard sweep throughout (Figure 4C). As with
9 the hard sweep scenario, G123 and H123 yield little change in resolution for detecting strong soft sweeps
10 from SSV, suggesting that the third-most frequent haplotype may have little importance in detecting sweeps
11 (Figures S3A and C). Once again, H123 maintains slightly greater power than does G123.

12 G12 and H12 perform comparably well for moderate ($s = 0.01$) sweeps from SSV (Figure 4B). Similarly
13 to hard sweep scenarios for $s = 0.01$, G12 and H12 detected soft sweeps from SSV occurring between 1,000
14 and 2,000 generations before sampling. Once again, the power of H12 was greater than that of G12, with
15 trends in power for G12 following those of H12. For both MLG and haplotype data, the inclusion of additional
16 selected haplotypes at the start of selection up to $k = 8$ only slightly reduced the maximum power of G12
17 and H12 to detect sweeps, but with time at which maximum power is reached shifting from 2,000 generations
18 before sampling for $k \leq 8$ to 1,000 generations before sampling for $k \geq 16$. Additionally, the spatial signal
19 for moderate sweeps was comparable between SSV and hard sweep scenarios (Figure 4D). This result may be
20 because at lower selection strengths, haplotypes harboring adaptive alleles are more likely to be lost by drift,
21 leaving fewer distinct selected haplotypes rising to appreciable frequency. These trends persist for G123 and
22 H123, which display similar powers to G12 and H12 across all scenarios (Figures S3B and D).

23 Effect of population size changes on detection capabilities of G12 and G123

24 Changes in population size that occur simultaneously with or after the time of selection may impact the ability
25 of methods to detect sweeps because haplotypic diversity may decrease under a population bottleneck, or
26 increase under a population expansion [Campbell and Tishkoff, 2008]. To test the robustness of the expected
27 homozygosity statistics to these potentially confounding scenarios, we modeled hard sweeps following the
28 human population bottleneck and expansion parameters inferred by Lohmueller et al. [2009] (Figure 2). We
29 measured the powers of the MLG- and haplotype-based methods across our previously-tested parameters,
30 using simulated 100 kb chromosomes and sliding windows, and approaching these scenarios in two ways.

31 First, we applied a 40 kb window as previously to evaluate the effect of population size change on the
32 power of expected homozygosity methods. Under a bottleneck, a 40 kb window is expected to carry fewer

1 SNPs than under a constant-size demographic history, whereas an expansion results in greater diversity
2 per window. Second, we examined whether we could increase the robustness of the expected homozygosity
3 methods to population size changes by adjusting the window size for each scenario to match the expected
4 number of segregating sites for a 40 kb window under constant demographic history. To do this, we followed
5 the approach outlined in DeGiorgio et al. [2014], increasing window size for bottleneck simulations and
6 decreasing window size for expansion simulations. We employed windows of size 56,060 nucleotides for
7 bottleneck, and of size 35,048 nucleotides for expansion scenarios [see DeGiorgio et al., 2014].

8 A recent population bottleneck reduces the powers of all methods to detect sweeps, whereas a recent
9 population expansion enhances power (Figures S5 and S6). This results from the genome-wide reduction
10 in haplotypic diversity under a bottleneck relative to the constant-size demographic history. Thus, the
11 maximum values of the expected homozygosity statistics in the absence of a sweep are inflated, resulting in
12 a distribution of maximum values under neutrality that has increased overlap with the distribution under
13 selective sweeps. In contrast, haplotypic diversity is greater under the population expansion than what
14 is expected for the constant-size demographic history, rendering easier the detection of elevated expected
15 homozygosity due to a sweep.

16 For strong selection ($s = 0.1$) under a population bottleneck, all methods using unadjusted windows have
17 reliable power to detect only recent hard sweeps to large f occurring within 1,000 generations of sampling
18 (Figures S5A and S6A). Adjusting window size has little effect on this trend, with powers for sweeps beginning
19 400 generations before sampling increasing only slightly (Figures S5C and S6C). This result indicates that
20 we can apply the expected homozygosity methods to populations that have experienced a severe bottleneck
21 and make accurate inferences about their selective histories. Similarly, adjusting window size had little
22 effect on the power of methods to detect a sweep under a population expansion, wherein power is already
23 elevated. As with the bottleneck scenario, reducing the size of a 40 kb window (Figure S5B and S6B) to
24 35,048 bases (Figure S5D and S6D) provided a minor increase in power to detect selective events occurring
25 within 2,000 generations of sampling, with high power for larger values of f extending to 2,000 generations
26 prior to sampling.

27 Distinguishing hard and soft sweeps with G2/G1

28 Having identified selective sweeps with the statistics G12 or G123, our goal is to make an inference about the
29 number of sweeping haplotypes. To distinguish between hard and soft sweeps, Garud et al. [2015] defined
30 the ratio H2/H1, which is larger under a soft sweep and smaller under a hard sweep. The H2/H1 ratio
31 leverages the observation that haplotypic diversity following a soft sweep is greater than that under a hard
32 sweep. Garud and Rosenberg [2015] showed that the value of H2/H1 is inversely correlated with that of

1 H12, and that identical values of H2/H1 have different interpretations depending on their associated H12
2 value. Therefore, H2/H1 **should** only be applied in conjunction with H12 when H12 is large enough to be
3 distinguished from neutrality.

4 Here, we extend the application of H2/H1 to MLGs. As with the haplotype approach, G2/G1 is larger
5 under a soft sweep and smaller under a hard sweep, because MLG diversity following a soft sweep is greater
6 than under a hard sweep. G2/G1 **should** therefore distinguish between hard and soft sweeps similarly to
7 H2/H1, conditional on a high G12 or G123 value. To demonstrate the classification ability of the MLG-
8 based methods with respect to the haplotype-based methods, we **began by generating** 10^6 simulated replicates
9 of 40 kb chromosomes with sample size $n = 100$ diploids for hard sweep and SSV scenarios, treating each
10 chromosome as a single window and recording its G12, G123, and G2/G1 values (see *Materials and methods*).

11 We evaluated the ability of G2/G1 with G12 or G123 to distinguish between hard sweeps and soft
12 sweeps from SSV **specifically** from $k = 3$ and $k = 5$ **drawn** haplotypes, both within the range of method
13 detection (**Figures 4 and S3**), with all sweeps allowed but not guaranteed to go to fixation. We examined
14 two values of k , distinct from one another and from hard sweeps, to illustrate the effect of model choice on
15 sweep classification. Each experiment evaluated the likelihood that a soft sweep scenario would produce a
16 particular paired (G12, G2/G1) or (G123, G2/G1) value relative to a hard sweep scenario. We measured
17 this relative likelihood by plotting the Bayes factors (BFs) for paired (G12, G2/G1) and (G123, G2/G1) test
18 points generated from an approximate Bayesian computation (**ABC**) approach (see *Materials and methods*).
19 A $\text{BF} > 1$ indicates a greater likelihood of a soft sweep generating the paired values of a test point, and
20 a $\text{BF} < 1$ indicates that a hard sweep is more likely to have generated **such values**. In practice, however,
21 we only assign $\text{BF} \leq 1/3$ as hard and $\text{BF} \geq 3$ as soft to avoid making inferences about borderline cases
22 (**Figure 5**). For each replicate, time of selection (t) and selection strength (s) were drawn uniformly at
23 random on a log-scale from $t \in [40, 2000]$ generations before sampling and $s \in [0.005, 0.5]$.

24 The comparison of hard sweep and SSV scenarios provides a distribution of **BFs** broadly in agreement
25 with expectations for the **haplotype-based approaches** (Garud et al. [2015], Garud and Rosenberg [2015];
26 **Figure 5**). In **Figure 5**, colored in blue **are the values** most likely to be generated under **SSV**, and colored
27 in red **are the values** most likely to be generated under hard **sweeps**. In all scenarios tested, hard sweeps
28 produce relatively smaller G2/G1 values than do soft sweeps. Intermediate **G12 and G123** paired with large
29 values of G2/G1 are more likely to result from soft sweeps than from hard sweeps. SSV cannot generate
30 large values of G12 or G123 because these sweeps are too soft to elevate homozygosity levels to the extent
31 observed under hard sweeps. This is particularly so when soft sweeps are simulated with $k = 5$. Therefore,
32 the majority of test points with extreme values of G12 and G123, regardless of G2/G1, have $\text{BF} \leq 1/3$
33 (**meaning** only one SSV observation within a Euclidean distance of 0.1 **for every three or more hard sweep**

1 observations), and this is in line with the results from the constant-size demographic model of Garud et al.
2 [2015] for comparisons between hard sweeps and the softest soft sweeps. Additionally, we cannot classify
3 sweeps if the values of G12 and G123 are too low, as these values are unlikely to be distinct from neutrality.
4 Thus, our ability to distinguish between hard and soft sweeps is greatest for intermediate values of G12 and
5 G123. In practice, our empirical top sweep candidates all converge over this range of the (G12, G2/G1) and
6 (G123, G2/G1) values (Figure 6), meaning that we can confidently classify sweeps from outlying values of
7 G12 and G123 in our data as hard or soft.

8 In Figure S7, we repeat our ABC procedure for the phased haplotype data corresponding to our preceding
9 analyses. We find that a small proportion of (G12, G2/G1) and (G123, G2/G1) values for which we lack the
10 ability to distinguish hard and soft sweeps (gray points), corresponds to (H12, H2/H1) values that do classify
11 sweeps as soft. Additionally, the (H123, H2/H1) values yielded a still larger proportion of SSV-classified
12 (blue) values. This result may indicate that the haplotype approaches maintain a somewhat greater ability
13 to classify sweeps than do the MLG approaches. Accordingly, the skew toward larger BFs among the (G123,
14 G2/G1) values relative to (G12, G2/G1) may indicate that classification with the former may more closely
15 resemble classification using (H12, H2/H1) values.

16 To further characterize the classification properties of both the MLG- and haplotype-based approaches,
17 we next employed an alternative ABC approach in which we determined the posterior distribution of k for a
18 range of (G12, G2/G1), (G123, G2/G1), (H12, H2/H1), and (H123, H2/H1) value combinations. For these
19 experiments, we generated 5×10^6 replicates of sweep scenarios with $k \in \{1, 2, \dots, 16\}$ drawn uniformly at
20 random for each replicate, maintaining all other relevant parameters identical to the BF experiments (see
21 *Materials and methods*). From the posterior distribution of k values, we assigned the most probable k for
22 a wide range of points using both MLG and haplotype data (Figure S8), and generated probability density
23 functions across H12, H2/H1, G123, and G2/G1 for each value of k (Figure S9). G12, G123, H12, and H123
24 values were larger for sweeps with smaller k , and G2/G1 values were smaller for these sweeps, as expected.
25 We achieved a finer resolution from haplotypes than from MLGs, as in the BF experiments (Figures 5 and S7),
26 and found our inference of the most probable values of k across test points to be concordant with BF-based
27 results. As previously, hard sweeps ($k = 1$) occupied larger values of G12, G123, H12, and H123 and smaller
28 values of G2/G1 and H2/H1, with inferred k (similarly to inferred BF) increasing with increasing G2/G1
29 and H2/H1, regardless of G12, G123, H12, and H123 value. Thus, our alternative ABC approach can assign
30 a most probable k from the entire tested range of $k \in \{1, 2, \dots, 16\}$, allowing for sweep classification without
31 the ambiguity of BFs.

¹ Analysis of empirical data for signatures of sweeps

² We applied G12, G123, and H12 to whole-genome variant calls on human autosomes from the 1000 Genomes
³ Project [Auton et al., 2015] to compare the detective properties for each method on empirical data (Fig-
⁴ ures 7 and S11-S18; Tables S3-S14). This approach allowed us to understand method performance in the
⁵ absence of confounding factors such as missing data and small sample size. The choice of human data
⁶ additionally allowed us to validate our results from the wealth of identified candidates for selective sweeps
⁷ within human populations worldwide that has emerged from more than a decade of research [e.g., Sabeti
⁸ et al., 2002, Bersaglieri et al., 2004, Voight et al., 2006, Bhatia et al., 2011, Chen et al., 2015, Schrider and
⁹ Kern, 2016, Cheng et al., 2017]. To apply our MLG-based methods to the empirical dataset, consisting of
¹⁰ haplotype data, we manually merged the haplotypes for each study individual to generate MLGs. Thus, all
¹¹ comparisons of G12 and G123 with H12 were for the same data, as in our simulation experiments.

¹² For our analysis of human data, we focused on individuals from European (CEU), African (YRI), South
¹³ Asian (GIH), and East Asian (CHB) descent. Across all populations, we assigned *p*-values and BFs, as well
¹⁴ as maximum posterior estimates and Bayesian credible intervals on *k*, for the top 40 selection candidates (see
¹⁵ *Materials and methods*). Our Bonferroni-corrected significance threshold [Neyman and Pearson, 1928] was
¹⁶ 2.10659×10^{-6} , with critical values for each statistic in each population displayed in Table S1. We defined
¹⁷ soft sweeps as those with $\text{BF} \geq 3$ or inferred $k \geq 2$, and hard sweeps as those with $\text{BF} \leq 1/3$ or inferred
¹⁸ $k = 1$. Following each genome-wide scan, we filtered our raw results using a mappability and alignability
¹⁹ measure (see *Materials and methods*), following the approach of Huber et al. [2016]. We additionally omitted
²⁰ genomic windows from our analysis with fewer than 40 SNPs, the expected number of SNPs in our genomic
²¹ windows [Watterson, 1975] under the assumption that a strong recent sweep has affected all but one of the
²² sampled haplotypes. This is thus a conservative approach. We display the filtered top 40 outlying sweep
²³ candidates for G12, G123, and H12, including *p*-values, BFs, and inferred *k* (with credible interval), in
²⁴ Tables S3-S14. We also overlay the top 40 selection candidates for each population onto (G123, G2/G1)
²⁵ test points (Figures 6 and S10). For all populations, we see that top candidates, regardless of assignment as
²⁶ hard or soft, generate broadly similar G123 values within a narrow band of paired (G123, G2/G1) values.
²⁷ Finally, we indicate the top 10 selection candidates in chromosome-wide Manhattan plots for both G12 and
²⁸ G123 (Figures S11-S18). Expectedly, G12 and G123 plots are nearly identical in their profiles.

²⁹ We recovered significant signals from the well-documented region of CEU chromosome 2 harboring the
³⁰ *LCT* gene, which confers lactase persistence beyond childhood [Bersaglieri et al., 2004]. Although filtering
³¹ removed *SLC24A5*, another expected top candidate controlling skin pigmentation, the adjacent *SLC12A1*
³² gene remained. Assigned BFs and inferred values of *k* suggest that hard sweeps in each of these regions yield
³³ the observed signals (Tables S3 and S4). In YRI (Tables S6 and S7), we most notably found the previously-

1 identified *SYT1*, *HEMGN*, and *NNT* [Voight et al., 2006, Pickrell et al., 2009, Fagny et al., 2014, Pierron
2 et al., 2014]. *SYT1* and *HEMGN* were significant for G12, G123, and H12 analyses, with *SYT1* yielding
3 the strongest signal by a large margin, while *NNT* was not significant. Of these, we could only confidently
4 classify *HEMGN*, which we uniformly identified as hard. Though we were more likely to confidently classify
5 candidate sweeps in YRI as hard from their MLG-based BFs, the proportion of top candidates assigned as
6 hard from the posterior distribution of k remained comparable across data types, and generally greater than
7 the levels we observed in other populations (see *Discussion* for further analysis). The most outlying target
8 of selection in GIH (Tables S9 and S10) for all methods was at *SLC12A1*, a significant signal corresponding
9 to a sweep shared among Indo-European populations [Mallick et al., 2013], which we also recovered as a top
10 candidate in CEU. We could classify this signal as hard from haplotype data, but we assigned $k = 2$ from
11 MLGs, despite a $\text{BF} < 1$. Finally, our analysis of CHB returned *EDAR*-adjacent genes among the top sweep
12 candidates, including *LIMS1*, *CCDC138*, and *RANBP2* (each below the significance threshold), though not
13 *EDAR* itself (Tables S12 and S13), and additionally *MIR548AE2* and *LONP2*, adjacent to the site of a
14 proposed sweep on earwax texture within *ABCC11* [Ohashi et al., 2010], which we recovered as another top
15 candidate.

16 In Figure 7, we highlight for each population one example of a sweep candidate, including its G12
17 signal profile, with the genomic window of maximum value highlighted, and a visual representation of the
18 MLG diversity within that region. For the CEU population, we present *LCT* ($p < 10^{-6}$), and additionally
19 highlight the nearby outlying candidates, each of which was within the top 10 outlying G12 signals in the
20 population (Figure 7A, left panel). The distribution of MLGs surrounding *LCT* in the sample showed a single
21 predominant MLG comprising approximately half of individuals, consistent with a hard sweep (Figure 7A,
22 right panel). Accordingly, *LCT* yielded a $\text{BF} \approx 0.1$, indicating that a hard sweep is tenfold more likely to
23 yield this signal than a soft sweep (from $k = 5$), and an inferred $k = 1$ supports this result. For the YRI
24 population, the top selection signal for all analyses was *SYT1* ($p = 10^{-6}$), previously identified by Voight et al.
25 [2006] (Figure 7B, left panel). Here, one high-frequency and one intermediate-frequency MLG predominated
26 in the population (Figure 7B, right panel), but we could not confidently assign the signal as hard or soft, with
27 haplotypes suggesting $k = 1$ and MLGs suggesting $k = 2$. This is because one high-frequency haplotype
28 exists in the population, carried by approximately half of individuals, while another haplotype exists in
29 approximately one quarter of individuals. In GIH, we found *P4HA1* as a selection candidate exceeding the
30 significance threshold for haplotype data ($p = 10^{-6}$), but not for MLG data. Although we were unable
31 to confidently assign the putative sweep on *P4HA1* as hard or soft from BFs, we note that two MLGs, as
32 well as two haplotypes, exist at elevated frequency here, and that all methods yielded $\text{BF} > 1$ and $k > 1$,
33 suggesting that *P4HA1* is likely the site of a soft sweep, but on fewer than $k = 5$ haplotypes (Figure 7C,

1 right panel). Finally, our scan in CHB returned the undocumented *FMNL3* gene as a top candidate from
2 the G12 analysis ($p = 5 \times 10^{-6}$; Figure 7D, left panel). A single high-frequency MLG predominated at this
3 site, and this yielded a BF from MLG data of 0.147, and inferred $k = 1$ from all data, indicating a hard
4 sweep (Figure 7D, right panel).

5 Through the application of G123 and G2/G1 we have identified and classified a number of interesting
6 sweep candidates. We further explored the existence of a more general relationship between top sweep
7 candidates and the prevalence and length of runs of homozygosity. Previous research has indicated that
8 short-to-intermediate runs of homozygosity spanning tens to hundreds of kilobases are characteristic of
9 recent sweeps [Pemberton et al., 2012, Blant et al., 2017], and we sought to examine whether there was a
10 correlation of G123 or sweep softness (using $\log_{10}(\text{BF})$ as proxy) with the proportion of individuals falling
11 in a run of homozygosity of specific length. To this end, we intersected our top candidates lists with the
12 inferred coordinates of short to intermediate runs of homozygosity from Blant et al. [2017]. We found that the
13 proportion of individuals with runs of homozygosity of intermediate length (class 4) is positively correlated
14 (correlation coefficient = 0.32, p -value = 3.66×10^{-5}) with G123 (Table S2), likely due to stronger and more
15 recent sweeps generating larger G123. Moreover, the proportion of individuals with runs of homozygosity
16 of intermediate length is negatively correlated (correlation coefficient = -0.26, p -value = 1.02×10^{-3}) with
17 $\log_{10}(\text{BF})$ (Table S2), likely due to the narrower genomic signature left behind by soft sweeps relative to
18 hard sweeps. In contrast, we observe no significant correlation for smaller runs of homozygosity (classes 2
19 and 3), which have also been proposed to potentially be affected by selective sweeps [Pemberton et al., 2012,
20 Blant et al., 2017].

7

potentially
mislea-
ding /
not very
helpful

21 Discussion

22 Selective sweeps represent an important outcome of adaptation in natural populations, and detecting these
23 signatures is key to understanding the history of adaptation in a population. We have extended the existing
24 statistics H12 and H2/H1 [Garud et al., 2015] from phased haplotypes to unphased MLGs as G12, G123, and
25 G2/G1, and demonstrated that the ability to detect and classify selective sweeps as hard or soft remains.
26 Across simulated selective sweep scenarios covering multiple selection start times and strengths, as well as
27 sweep types and demographic models, we found that both G12 and G123 maintain comparable power to
28 H12. The most immediate implication of these results is that signatures of selective sweeps can be identified
29 and classified in organisms for which genotype data are available, without the need to generate phased
30 haplotypes. Because phasing may be difficult or impossible given the resources available to a study system,
31 while also not being error-free [Browning and Browning, 2011, O'Connell et al., 2014, Laver et al., 2016,
32 Castel et al., 2016, Zhang et al., 2017], the importance of our MLG-based approach is apparent. Although

1 phased haplotypes tend to be preferable for use with expected homozygosity statistics based on our findings,
2 we nonetheless observe a high degree of congruence in practice between the lists of selection candidates for
3 human empirical data emerging from analyses on haplotypes and MLGs (Tables S3-S14).

4 Performance of G12 and G123 for simulated data

5 G12 and G123, similarly to H12 and H123, are best suited to the detection of recent **and strong** selective
6 sweeps in which the beneficial allele has risen to appreciable frequency. This is as expected because haplotype (red)
7 (and therefore MLG) homozygosity increases under sweeps, which results in a distinct signature from which to
8 infer the sweep. This extended tract of sequence identity within the population erodes over time and returns
9 to neutral levels due to the effects of recombination and mutation. The strength of selection and range of
10 time over which the expected homozygosity-based methods can detect selection are inversely correlated. Our
11 approach detects weaker selective events only if they started far enough back in time, and has a narrower
12 time interval of detection than do stronger events (compare panels A and B across Figures 3, 4, S2, and S3).
13 This is because alleles under weaker selection increase in frequency toward fixation more slowly than those
14 under stronger selection, and so more time is required to generate a detectable signal. In the process, the
15 size of the genomic tract that hitchhikes with the beneficial allele decreases due to recombination and is
16 smaller than under a hard sweep. Panels C and D from Figures 3, 4, S2, and S3 motivate this point. Across
17 all simulation scenarios, stronger selection produces on average a wider and larger signature surrounding the
18 site of selection, while weaker sweeps are more difficult to detect and classify. For empirical analyses, this
19 means we are more likely to detect stronger sweeps, as reductions in diversity from strong selection persist
20 for hundreds of generations and can leave footprints on order of hundreds of kilobases [Gillespie, 2004, Garud
21 et al., 2015, Hermisson and Pennings, 2017].

22 Expectedly, the signatures of sweeps, and the power of the expected homozygosity methods to detect
23 them, vary across selective sweep scenarios, with nearly identical trends in haplotype and MLG data. Strong
24 ($s = 0.1$) hard sweeps to high sweep frequency f are easiest to detect, as the single, large tract of sequence
25 identity generated under a strong hard sweep remains distinct from neutrality for the longest time interval
26 relative to other scenarios (Figures 3A and C and Figures S2A and C). Nonetheless, power to distinguish
27 soft sweeps is large for the most recent simulated sweeps. Indeed, a soft sweep yields a smaller tract of
28 sequence identity that requires a shorter time to break apart, but for strong selection on up to $k = 16$
29 different haplotypic backgrounds (1.6% of the total population), both the MLG and haplotype methods
30 have perfect or nearly-perfect power (Figures 4A and S3A). While this power rapidly fades for selection
31 within 1,000 generations of sampling for $k > 4$, our strong sweep results illustrate that selection coefficient s , more than partial sweep frequency f or number of initially-selected haplotypes k , influences the power

1 of our pooled expected homozygosity methods, and that pooling can allow for similar detection of hard
2 and soft sweeps. Our moderate selection ($s = 0.01$) results further highlight this. Once again, we see a
3 distinct concordance in power trends between hard (Figures 3B and D and Figures S2B and D) and soft
4 (Figures 4B and D and Figures S3B and D) sweeps that depends primarily on the value of s and secondarily
5 on f or k .

6 Because genomic scans using G12, G123, H12 and H123 are window-based, the choice of window size is
7 an **important determinant of** the methods' sensitivity. As do Garud et al. [2015], we recommend a choice
8 of window size that minimizes the influence of background LD on window diversity, while maximizing the
9 proportion of sites in the window affected by the sweep. Windows that are too small may contain extended
10 homozygous tracts not resulting from a sweep, while windows that are too large will contain an excess
11 of neutral diversity leading to a weaker signal, **while overlooking weaker selective events** [Gillespie, 2004,
12 Garud et al., 2015, Hermisson and Pennings, 2017]. Accordingly, our choice of a 40 kb sliding window
13 to analyze simulation results derives from our observation that the value of LD between pairs of SNPs
14 separated by 40 kb in these simulations is less than one-third of the LD between pairs separated by **one** kb,
15 as measured from the squared correlation, r^2 (Figure S1). We also found that for recent selection within 400
16 generations of sampling, power under bottleneck or **expansion** does not change **for a 40 kb analysis window**
17 (Figures S5 and S6). This is especially important in the context of a population bottleneck, in which levels of
18 short-range LD are elevated beyond their expected value under a constant-size demographic history [Slatkin,
19 2008, DeGiorgio et al., 2009]. Thus, our population size change experiments indicated that for sufficiently
20 large analysis windows, further adjusting window size does not improve power. The trends in power that
21 we observed for samples of $n = 100$ diploids and 40 kb genomic windows also persisted for experiments with
22 a smaller sample size of $n = 25$ (Figure S19). The expected homozygosity methods are therefore suitable
23 for detecting sweeps from a wide range of sample sizes, though samples need to be large enough to capture
24 the difference in variation between selected and neutral regions of the genome, as smaller samples result in
25 fewer sampled haplotypes [Pennings and Hermisson, 2006a]. Accordingly, the classification of sweeps requires
26 substantially larger sample sizes, as differentiating between hard and soft sweeps requires the detection of a
27 more subtle signal than does distinguishing selection from neutrality.

28 Although we exclusively used a nucleotide-delimited window in our present analyses, it is possible to
29 search for signals of selection using a SNP-delimited window, and this was the approach of Garud et al.
30 [2015]. Similarly to our present approach, the number of SNPs to include in a window could be determined
31 based on the decay in pairwise LD between two sites separated by a SNP-delimited interval. Under the
32 SNP-delimitation approach, each analyzed genomic window includes a specified number of SNPs. Thus, the
33 range of physical window sizes may be broad. In principle, the use of a SNP-delimited window prevents the

1 inclusion of SNP-poor windows. Accordingly, SNP delimitation may be inherently robust to the effect of
2 bottlenecks, or to the misidentification of heterochromatic regions as sweep targets. In practice, however, we
3 can filter out nucleotide-delimited genomic windows carrying too few SNPs to overcome confounding signals.
4 More importantly, allowing for a variable number of SNPs in a window allows the genomic scan to identify
5 sweeps not only from distortions in the haplotype frequency spectrum, but also from reductions in the total
6 number of distinct haplotypes, which are more constrained in their range of values when conditioned on a
7 specific number of SNPs. Because both of these signatures can indicate a sweep, it may be useful to consider
8 each. Even so, the use of a SNP-delimited window may be preferable for SNP chip data. That is, SNP
9 density can be low relative to whole-genome data, resulting in an excess of regions spuriously appearing
10 to be under selection within a nucleotide-delimited window. Indeed, Schlamp et al. [2016] employ a SNP-
11 delimited window approach for their canine SNP array dataset.

shorten

-1

|

May
be
to a
restrictive

underline
why

12 During a genomic scan, it may also be helpful to account for sources of uncertainty in the data. Foremost
13 among these is uncertainty in genotype calls [Marchini and Howie, 2010, Nielsen et al., 2011]. Modern geno-
14 type calling methods provide a posterior probability for each genotype [He et al., 2014, Korneliussen et al.,
15 2014, Fumagalli et al., 2014], and so it may be possible to assign to each analysis window a weighted mean
16 G12 or G123 score from this posterior to produce a more accurate representation of sweep events throughout
17 the study population's genome. It is also possible that windows of elevated G12 and G123 value may arise in
18 the absence of random mating. That is, although our approach assumes elevated MLG homozygosity derives
19 from elevated haplotype homozygosity as a result of random mating, we do not specifically evaluate whether
20 observed patterns of MLG diversity are compatible with the random mating assumption. Such an approach
21 could condition on the presence of one high-frequency MLG with only homozygous sites in the case of a hard
22 sweep, or at least two high-frequency homozygous MLGs in the case of a soft sweep. To further consider
23 this point, we rescanned the 1000 Genomes dataset, but randomly paired haplotypes into diploid MLGs to
24 simulate random mating. Our lists of outlying sweep candidates for G123 across each study population after
25 random reshuffling were highly concordant with the lists for the true set of diploid individuals (Tables S5,
26 S8, S11, and S14).

|

27 While power to detect hard and soft sweeps is comparable, the possible values of G12 and G2/G1 that
28 can be generated under hard versus soft sweeps for a variety of k values are distinct. Thus, we can properly
29 classify sweeps from MLG data (Figure 5, 6, S8, and S10). This result matched our theoretical expectations
30 (Figure 1), and corresponded to the results from haplotype data as well (Figure S7). However, we note that
31 with the BF-based ABC approach there is substantial ambiguity in classification over which $1/3 \leq BF \leq 3$
32 (where BF is computed as Probability(soft)/Probability(hard)), meaning that distinguishing between hard
33 and soft sweeps for these paired values remains difficult or not meaningful. In addition, we find that MLGs

1 (Figure 5) provide a greater proportion of $BF \leq 1/3$ than do haplotypes (Figure S7), which yield a greater
2 proportion of $BF \geq 3$. This observation may indicate that a hard sweep with a small associated BF for
3 MLGs will also have a small haplotype-based BF, while a hard sweep with an associated BF closer to 1, may
4 be called as ambiguous or soft from haplotypes. We were able to address the issue of classification ambiguity
5 with our alternative ABC approach, which assigned each test point a most probable underlying k . Although
6 haplotypes provided better ability over MLGs to assign a posterior value of k , our results here were as
7 expected, showing a clear increase in assigned k as G2/G1 or H2/H1 increased (Figure S8). For application
8 to empirical data, however, most top sweep candidates are likely to be classifiable as hard or soft from
9 BFs (Tables S3-S14). Pooling frequencies beyond the greatest two also increased the occupancy associated
10 with larger BFs, and this effect was greater for haplotype data. Ultimately, the use of G123 with G2/G1 to
11 classify sweeps and assign k from MLGs may be preferable because (G123, G2/G1) classification more closely
12 resembles (H12, H2/H1) than does (G12, G2/G1). The true value of pooling additional frequencies may
13 thus lie in sweep classification rather than detection, as G123 and H123 are not appreciably more powerful
14 than G12 and H12 (Figures S2 and S3).

15 Application of G12 and G123 to empirical data

16 Our analysis of human empirical data from the 1000 Genomes Project [Auton et al., 2015] recovered multiple
17 positive controls from each study population, as well as novel candidates. Across many of these candidates,
18 a single high-frequency MLG predominated (Figure 7). Additionally, more top candidates in CEU appear
19 as hard sweeps than in other populations (Tables S3 and S4), though all populations had more hard sweeps
20 than soft. The top outlying genes we detected in CEU following the application of a filter to remove
21 heterochromatic regions with low mappability and alignability consisted of *LCT* and the adjacent loci of
22 chromosome 2 (Figure 7A), as well as *SLC12A1* of chromosome 15 (Table S3). All of these sites are well-
23 represented in the literature as targets of sweeps [Bersaglieri et al., 2004, Sabeti et al., 2007, Liu et al., 2013,
24 Chen et al., 2015]. Diet-mediated selection on *LCT* likely drives the former signal cluster, as dairy farming
25 has been a feature of European civilizations since antiquity [Itan et al., 2009, Edwards et al., 2011, Ermini
26 et al., 2015]. Accordingly, we see that most individuals in the sample carry the most frequent MLG, and we
27 assign this signal to be a hard sweep from its BF and from the posterior distribution of k generated under our
28 demographic model for CEU (see *Materials and Methods*; Tables S3 and S4). Meanwhile, the latter signal
29 peak is associated with the known target of selection *SLC24A5*, a melanosome solute transporter responsible
30 for skin pigmentation [Lamason et al., 2005], also a hard sweep.

31 The assignment of sweeps as hard or soft in CEU, as well as their assigned k , were highly concordant
32 between haplotype and MLG approaches, with the sole exception of *PRKDC*, a protein kinase involved in

1 DNA repair [Fushan et al., 2015]. Our haplotype results indicate the presence of $k = 3$ high-frequency
2 haplotypes at *PRKDC*, but MLG results suggest a hard sweep. This is because the window of maximum
3 signal differs between both data types. The maximal haplotype-based window features multiple haplotypes
4 and MLGs at high frequency, while the maximal MLG window approximately 35 kb upstream more closely
5 resembles a hard sweep for both data types. We found such classification discrepancies to be rare across
6 our top candidates, and typically inverted, with the MLG signal more often appearing softer (see *SYT1* and
7 *RGS18*; Figure 7). Furthermore, we emphasize that classification discrepancies do not appear to impact
8 the power of MLG-based methods to detect sweeps, as we generated highly concordant lists of outlying
9 candidates for both haplotype and MLG data.

10 Large tracts of MLG homozygosity surround the *SYT1*, *RGS18*, *HEMGN*, *KIAA0825*, and *NNT* genes in
11 YRI. Unlike for CEU, we found that assigning BFs to top signals was difficult, both for haplotype and MLG
12 data (Tables S6 and S7). We also note a greater proportion of soft sweeps among top signals in YRI relative
13 to other populations (Tables S6 and S7). This is likely due to the greater ease of detecting soft sweeps in
14 more genetically diverse populations rather than any non-adaptive confounding factor (see next subsection),
15 and we indeed see a larger occupancy of soft BFs among (G123, G2/G1) values (Figure 6). In addition, BFs
16 for the two top candidates, *SYT1* and *RGS18*, yielded values close to 1/3 (hard) for haplotype data, but
17 closer to 3 (soft, $k = 2$) for MLG data, indicating disproportionately large MLG diversity resulting from
18 low haplotypic diversity, as the presence of a high-frequency haplotype alongside one or more intermediate-
19 frequency haplotypes may generate comparatively more diversity among MLGs than haplotypes. Voight
20 et al. [2006] previously identified our strongest selection target, *SYT1*, as a target of selection in the YRI
21 population, and The International HapMap Consortium [2007] corroborated this, but neither speculated as
22 to the implications of selection at this site. *SYT1* (Figure 7B) is a cell surface receptor by which the type
23 B botulinum neurotoxin enters human neurons [Connan et al., 2017]. Selection here may be a response
24 to pervasive foodborne bacterial contamination by *Clostridium botulinum*, similar to what exists in modern
25 times [Chukwu et al., 2016]. Pierron et al. [2014] named *HEMGN* (which Pickrell et al. [2009] also identified),
26 involved in erythrocyte differentiation, as a selection signal common to Malagasy populations derived from
27 common ancestry with YRI. Racimo [2016] also identified *KIAA0825* as a target of selection, but in the
28 ancestor to **African** and Eurasian populations. Our identification of *NNT* in YRI matches the result of
29 Fagny et al. [2014], who identified this gene using a combination of iHS [Voight et al., 2006] and their derived
30 intraallelic nucleotide diversity (DIND) method. Fagny et al. [2014] point out that *NNT* is involved in the
31 glucocorticoid response, which is variable among global populations. Our **most** noteworthy candidate of
32 selection in YRI, *RGS18*, has not been previously characterized as the location of a sweep. However, Chang
33 et al. [2007] point to *RGS18* as a contributor to familial hypertrophic cardiomyopathy (HCM) pathogenesis.

could
be
here!

1 HCM is the primary cause of sudden cardiac death in American athletes [Barsheshet et al., 2011], and
2 particularly affects African-American athletes [Maron et al., 2003].

3 Our scan for selection in the GIH population once again revealed the *SLC12A1* site as the strongest sweep
4 signal (Tables S9 and S10). Because this signal is common to Indo-European populations [Liu et al., 2013,
5 Ali et al., 2014], this was expected. However, we found that we could not confidently classify this sweep from
6 MLG data (with inferred $k = 2$), though haplotype data suggests that this is a hard sweep. We additionally
7 find *P4HA1* (Figure 7C) as a novel sweep candidate in GIH that exceeds the significance threshold for
8 haplotype data, and appears as a near-soft sweep for MLGs ($BF > 2.5$) with inferred $k \geq 2$ for both
9 haplotype and MLG data. Two high-frequency MLGs predominate at the location of this candidate sweep,
10 and their pooled frequency yields a prominent signal peak. *P4HA1* is involved in collagen biosynthesis,
11 with functions including wound repair [Baxter et al., 2013], and the population-variable hypoxia-induced
12 remodeling of the extracellular matrix [Petousi et al., 2013, Chakravarthi et al., 2014]. Because selection on
13 *P4HA1* has been documented among both the tropical forest-dwelling African pygmy population [Mendizabal
14 et al., 2012, Amorim et al., 2015] and now in individuals of Gujarati descent, and is known to present a
15 differing expression profile among low- and high-altitude populations [Petousi et al., 2013], this gene may be
16 involved in a number of adaptations to harsh climatic conditions, potentially in wound repair, which is more
17 difficult in tropical climates.

18 Of the sweep candidates we identified in the CHB population (Tables S12 and S13), we found that the
19 inference of significance from G123 was considerably more concordant with H12 than was G12. We recovered
20 as top candidates *EXOC6B*, which produces a protein component of the exocyst [Evers et al., 2014] and
21 *LONP2*, both previously documented [Baye et al., 2009, Ohashi et al., 2010, Durbin and Consortium, 2011,
22 Pybus et al., 2014]. *EXOC6B* is a characteristic signal in East Asian populations alongside *EDAR*, which
23 we did not specifically recover in our scan (but nearby candidates *LIMS1*, *CCDC138*, and *RANBP2* did
24 appear), while *LONP2* is adjacent to *ABCC11*, which controls earwax texture. *FMNL3* yielded elevated
25 values of G12 and G123 in CHB, but was only significant from its H12 value. A single MLG predominates at
26 *FMNL3* in the sample (Figure 7D), and all approaches assign this sweep as hard. The function of *FMNL3* is
27 related to actin polymerization [Hetheridge et al., 2012, Gauvin et al., 2014], and has a role in shaping the
28 cytoskeleton, which it shares with *EXOC6B*. Moreover, the signal at *FMNL3* may be additionally associated
29 with the outlier *RANBP10*, which also interacts with the cytoskeleton, but with microtubules [Schulze et al.,
30 2008]. Though it is unclear why we identify an enrichment in cytoskeleton-associated genes, future studies
31 may shed light on why variants in such genes could be phenotypically-relevant specifically in individuals of
32 East Asian descent. Finally, we found *SPATA31D3* as a hard sweep within the top H12 signals in CHB,

1 as well as in GIH, and while it did not exceed our significance threshold, this is in line with the results of
2 Schrider and Kern [2017].

3 Addressing confounding scenarios

4 A variety of processes, both adaptive and non-adaptive, may produce elevated values of expected homozy-
5 gosity in the absence of selective sweeps in a sampled population, or small values of expected homozygosity
6 despite a sweep, thereby misleading expected homozygosity methods. To understand the impacts of poten-
7 tially confounding processes on the power of the expected homozygosity methods, we evaluated the effects of
8 long-term background selection, long-term population substructure, and pulse admixture on G12, G123, H12,
9 and H123. We additionally consider the confounding effect of missing data, as the manner in which missing
10 sites is addressed during computations can change analyzed patterns of MLG and haplotype diversity.

11 We first addressed long-term background selection as a potentially common confounding factor with
12 a brief experiment to determine the susceptibility of all methods to the misidentification of background
13 selection as a sweep. Signatures of background selection are ubiquitous in a number of systems [McVicker
14 et al., 2009, Comeron, 2014], and the effect of background selection is a reduction in nucleotide diversity
15 and a distortion of the site frequency spectrum, which to many methods may spuriously resemble a sweep
16 [Charlesworth et al., 1993, 1995, Seger et al., 2010, Charlesworth, 2012, Nicolaisen and Desai, 2013, Cutter
17 and Payseur, 2013, Huber et al., 2016]. Here, we simulated chromosomes containing a centrally-located genic
18 region of length 11 kb in which deleterious alleles arise throughout the course of the simulation. Our model
19 involved a gene with exons, introns, and untranslated regions (UTRs) with properties based on human-
20 inspired parameters (see *Materials and methods*). In agreement with the result of Enard et al. [2014], we
21 found that background selection did not distort the haplotype (and therefore MLG) frequency spectrum to
22 resemble that of a sweep, such that G12 and G123 were thoroughly robust to background selection. We
23 demonstrate this by displaying the concordance in the distributions of maximum G12, G123, H12, and H123
24 scores for background selection and neutral evolution scenarios (Figure S20). Thus, we do not expect that
25 outlying G12, G123, H12, or H123 values can result from background selection.

26 Methods to detect recent sweeps may be confounded by the effect of long-term population substructure, as
27 well as from admixture. Structured populations contain a greater proportion of homozygous genotypes than
28 would be expected under an equally-sized, randomly-mating population [Sinnock, 1975], thereby increasing
29 the chance that an elevated level of expected homozygosity will arise in the absence of a sweep. We examined
30 the possibility that a symmetric island migration model with six demes (Figure S21A), and migration rates
31 (m) between demes of $m \in \{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$ per generation (a proportion m of the haplotypes
32 in a deme derives from each of the other five demes for a total proportion of $5m$ haplotypes) could yield

1 elevated values of H12 and G123 under neutrality. We found that compared to a model with no substructure,
2 H12 and G123 values were moderately impacted for a model with population substructure. These values
3 were substantially lower than expected H12 and G123 values under a recent strong hard sweep. However,
4 these values are more comparable to an ancient sweep, and so caution is warranted in the study of structured
5 populations for all but the most outlying signals.

refer
to
fig.
showing
results

(Fig.
S22)

May
be
due
to
RJ.?

6 The expected homozygosity methods are similarly robust to the effect of admixture under most scenarios.
7 Specifically, we evaluated whether any admixture scenario can falsely generate a signature of a sweep.
8 We simulated a model in which a single ancestral population diverges into two descendant populations
9 (Figure S21B; see also *Materials and methods*). We maintained the size of one descendant population (the
10 target) at $N = 10^4$ diploid individuals, and varied the size of the unsampled (donor) population ($N = 10^3$,
11 10^4 , or 10^5 diploids), admixing at rate $m \in \{0.05, 0.10, 0.15, 0.20, 0.25, 0.30, 0.35, 0.40\}$ as a single pulse 200
12 generations before sampling. We find that admixture with donor sizes $N = 10^5$ or $N = 10^4$ produces only
13 small values of H12 (Figure S23, left and center) and G123 (Figure S24, left and center) in the sampled
14 population in the absence of a sweep. However, admixture with a donor population of small size ($N = 10^3$)
15 can produce elevated values of H12 and H2/H1, as well as G123 and G2/G1 when migration is sufficiently
16 large ($m \geq 0.15$), thus spuriously resembling the pattern of a soft sweep in the absence of selection (Fig-
17 ures S23 and S24, right). In this scenario, with a large enough admixture fraction, there will be a high
18 probability that many sampled lineages from the target population will derive from the donor population,
19 which will coalesce rapidly due to the small effective size, which will in turn lead to elevated homozygosity.
20 Small donor population sizes with large migration rates therefore represent the only admixture scenario that
21 we considered under which the expected homozygosity methods are susceptible to misclassifying neutrality
22 as selection, specifically as a soft sweep. Otherwise, our methodology remains robust under a wide range of
23 other admixture scenarios. We note therefore that the elevated number of soft sweeps we detected within the
24 YRI population (Tables S6 and S7) is unlikely to be due to the effect of the admixture described in Busby
25 et al. [2016], as this would produce a genome-wide pattern, which we do not observe (Figures S13 and S14).

26 Finally, we note that accounting for missing data is a practical consideration that must be undertaken
27 when searching for signals of selection, and the manner in which missing data are removed affects our ability
28 to identify sweeps. We explored the effects of two corrective strategies to account for missing data. Our
29 strategies were to remove sites with missing data or to define MLGs and haplotypes with missing data
30 as new distinct MLGs and haplotypes. Relative to the ideal of no missing data (Figure 3A), removing
31 sites resulted in a slight inflation of power observed in the absence of missing data. This was true for
32 G12 and H12 (Figure S25A), as well as G123 and H123 (Figure S25C). After removing sites, the overall
33 polymorphism in the sample decreases, but windows containing the site of selection are still likely to be the

1 least polymorphic, and therefore identifiable. Even so, weaker sweeps are likely to be obscured by the lower
2 background diversity after removing sites. Conservatively defining MLGs and haplotypes with missing data
3 as new distinct MLGs and haplotypes inflates the total observed diversity and results in a more rapid decay
4 of power compared to complete data (Figures S25B and D). This result is because individuals affected by
5 the sweep may have different patterns in their missing data, and therefore different assigned sequences after
6 accounting for missingness. Overall, the choice of strategy will likely depend on the level of missing data in
7 the sample. Removing too many sites is likely to generate false positive signals, while removing no sites may
8 lead to false negatives.

9 Concluding remarks

10 Our results emphasize that detecting selective sweeps does not require phased haplotype data, as distortions
11 in the frequency spectrum of MLGs capture the reduction in diversity under a sweep similarly well to
12 phased haplotypes. Accordingly, the advent of rapid and cost-effective genotyping-by-sequencing technologies
13 [Elshire et al., 2011] across diverse taxa including bovine, marine-dwelling, and avian populations means that
14 the adaptive histories of myriad organisms may now be inferred from genome-wide data [Daetwyler et al.,
15 2014, Drury et al., 2011, Zhu et al., 2016]. Furthermore, we have shown that the inferences emerging
16 from MLG-based scans align with those of phased haplotype-based scans, with empirical analyses of human
17 populations yielding concordant top outlying candidates for selection, both documented and novel. We
18 demonstrate as well that paired (G12, G2/G1) and (G123, G2/G1) values properly distinguish hard sweeps
19 from soft sweeps. In addition to identifying sweeps from single large values of G12 and G123, we find that
20 the genomic signature of these MLG-based statistics surrounding the site of selection provides a means of
21 distinguishing a sweep from other types of selection (e.g., balancing selection). This additional layer of
22 differentiation motivates the use of MLG identity statistics as a signature in a statistical learning framework,
23 as such approaches have increasing in prominence for genome analysis [Grossman et al., 2010, Lin et al.,
24 2011, Pavlidis et al., 2010, Ronen et al., 2013, Pybus et al., 2015, Ronen et al., 2015, Sheehan and Song,
25 2016, Schrider and Kern, 2016, Akbari et al., 2017, Kern and Schrider, 2018, Mughal and DeGiorgio, 2018].
26 We expect that the MLG-based approaches G12 and G123, in conjunction with G2/G1, will be invaluable
27 in localizing and classifying adaptive targets in both model and non-model study systems.

¹ Materials and methods

² Simulation parameters

³ To compare the powers of G12 and G123 to detect sweeps relative to H12 and H123 [Garud et al., 2015],
⁴ we performed simulations for neutral and selection scenarios using SLiM 2 ([version 2.6](#)) [Haller and Messer,
⁵ 2017]. SLiM is a general-purpose forward-time simulator that models a population according to Wright-
⁶ Fisher dynamics [Fisher, 1930, Wright, 1931, Hartl and Clark, 2007] and can simulate complex population
⁷ structure, selection events, [recombination](#), and demographic histories. For our present work, we used SLiM
⁸ 2 to model scenarios of recent selective sweeps, long-term background selection, and neutrality, [additionally](#)
⁹ [including models of population substructure and pulse admixture](#). Our models of sweeps comprised complete
¹⁰ and partial hard sweeps, as well as soft sweeps from selection on standing variation (SSV). For background
¹¹ selection, we simulated a gene with introns, exons, and untranslated regions in which deleterious mutations
¹² arose randomly. We additionally tested the effect of demographic history on power by examining constant
¹³ population size, population expansion, and population bottleneck models for hard sweep scenarios.

¹⁴ General approach

¹⁵ We first simulated data according to human-[specific](#) parameters for a constant population size model. For
¹⁶ simulated sequences (Figures 2A and D), we chose a mutation rate of $\mu = 2.5 \times 10^{-8}$ per site per generation,
¹⁷ a recombination rate of $r = 10^{-8}$ per site per generation, and a diploid population size of $N = 10^4$ [Takahata
¹⁸ et al., 1995, Nachman and Crowell, 2000, Payseur and Nachman, 2000]. All simulations ran for a duration
¹⁹ of $12N$ generations, where N is the starting population size for a simulation, equal to the diploid effective
²⁰ population size. The duration of simulations is the sum of a $10N$ generation burn-in period of neutral
²¹ evolution to generate equilibrium levels of variation across simulated individuals [Messer, 2013], and the
²² expected time to coalescence for two lineages of $2N$ generations. Simulation parameters were scaled, as is
²³ common practice, to reduce runtime while maintaining expected levels of population-genetic variation, such
²⁴ that mutation and recombination rates were multiplied by a factor λ , while population size and simulation
²⁵ duration were divided by λ . For simulations of constant population size, we used $\lambda = 20$.

²⁶ Scenarios involving population expansion and bottleneck were modeled on the demographic histories
²⁷ inferred by Lohmueller et al. [2009]. For population expansion (Figures 2B and D), we used $\lambda = 20$, and
²⁸ implemented the expansion at 1,920 unscaled generations before the simulation end time. After expansion,
²⁹ the size of the simulated population doubled from 10^4 to 2×10^4 diploid individuals. This growth in size
³⁰ corresponds to the increase in effective size of African populations that occurred approximately 48,000 years
³¹ ago [Lohmueller et al., 2009], assuming a generation time of 25 years. Population bottleneck simulations

1 (Figures 2C and D) were scaled by $\lambda = 10$, began at 1,200 generations before the simulation end time,
2 and ended at 880 generations before the simulation end time. During the bottleneck, population size fell to
3 550 diploid individuals. This drop represents the approximately 8,000-year bottleneck that the population
4 ancestral to non-African humans experienced as it migrated out of Africa [Lohmueller et al., 2009], assuming
5 a generation time of 25 years.

6 Simulating selection

7 Our simulated selection scenarios encompassed a variety of selection modes and parameters. Though we
8 primarily focused on selective sweeps, we additionally modeled a history of long-term background selection
9 to test the specificity of methods for sweeps. Background selection may decrease genetic diversity relative to
10 neutrality. For sweep experiments specifically, we tested the power of methods to detect selection occurring
11 between 40 and 4,000 generations prior to the simulation end time (thus, within $2N$ generations prior to
12 sampling). We set the site of selection to be at the center of the simulated chromosome, and performed
13 two categories of simulations, allowing us to answer two distinct types of questions about the power of
14 our approach: whether G12 and G123 properly identify the signature of a selective sweep (the detection
15 experiments), and whether G12 or G123 in conjunction with G2/G1 can distinguish between hard and soft
16 sweeps and ultimately infer the number of selected haplotypes (k ; the classification experiments), and hence
17 “softness” of the sweep.

18 For the detection experiments (see *Detecting sweeps*), we simulated chromosomes of length 100 kb under
19 neutrality and for each set of selection parameters, performed 10^3 replicates of sample size $n = 100$ diploids
20 (and $n = 25$ for hard sweep experiments in Figure S19). Here, we fixed the times (t) at which selected alleles
21 arise to be 400, 1,000, 2,000, or 4,000 generations prior to sampling (Figure 2), and selection coefficients (s)
22 to be either 0.1 or 0.01, respectively representing strong and moderate selection. The parameters t and s were
23 common to all selection simulations of the first type, with additional scenario-specific parameters which we
24 subsequently define. For the classification experiments (see *Differentiating between hard and soft sweeps*), we
25 performed two types of simulations. First, we simulated 10^6 replicates of $n = 100$ diploids for each scenario,
26 with $s \in [0.005, 0.5]$, drawn uniformly at random from a natural log-scale, and $t \in [40, 2000]$ (also drawn
27 uniformly at random from a natural log-scale), across chromosomes of length 40 kb. With these simulations,
28 we assessed the occupancy of specific hard and soft sweeps among (G12, G2/G1), (G123, G2/G1), (H12,
29 H2/H1), and (H123, H2/H1) test points. Second, we simulated 5×10^6 replicates with $s \in [0.05, 0.5]$ and
30 $t \in [200, 2000]$ and all other parameters as previously. Here, we assigned the most probable k to each test
31 point from the posterior distribution of k among nearby test points, drawing $k \in \{1, 2, \dots, 16\}$ uniformly at
32 random. We scaled selection simulations as previously described.

1 We first examined hard sweeps, in which the beneficial mutation was added to one randomly-drawn
2 haplotype from the population at time t , remaining selectively advantageous until reaching a simulation-
3 specified sweep frequency (f) between 0.1 and 1.0 at intervals of 0.1, where $f < 1.0$ represents a partial
4 sweep and $f = 1.0$ is a complete sweep (to fixation of the selected allele). Although we conditioned on
5 the selected allele not being lost during the simulation, we did not require the selected allele to reach f .
6 We additionally modeled soft sweeps from selection on standing genetic variation (SSV). For this scenario,
7 we introduced the selected mutation to multiple different, but not necessarily distinct, randomly-drawn
8 haplotypes (k) such that $k = 2, 4, 8, 16$, or 32 haplotypes out of $2N = 10^3$ (scaled haploid population size)
9 acquired the mutation at the time of selection. We did not condition on the number of remaining selected
10 haplotypes at the time of sampling as long as the selected mutation was not lost.

11 For hard sweeps only, we additionally examined the effects of three common scenarios—population
12 substructure, pulse admixture, and missing data—on performance. The population substructure model
13 consisted of six demes in a symmetric island migration model in which migration between each deme is
14 constant at rate m per generation for the duration of the simulation (Figure S21A). We simulated $m \in$
15 $\{10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}\}$. All demes were identical in size at $N = 1,660$ (unscaled) diploid individuals,
16 and samples consisted of $n = 100$ diploid individuals, with 50 individuals sampled from each of two demes.
17 Thus, as m increases, the structured model converges to the unstructured model of $N = 10^4$ (unscaled)
18 diploid individuals. Our admixture scenarios examined a single pulse of gene flow from an unsampled
19 donor population into the sampled target at rate $m \in \{0.05, 0.10, 0.15, 0.20, 0.25, 0.30, 0.35, 0.40\}$, occurring
20 200 generations prior to sampling (Figure S21B). We performed experiments in which the donor had a
21 (unscaled) diploid size of $N = 10^3, 10^4$, or 10^5 , keeping the size of the target fixed at $N = 10^4$. For
22 admixture simulations, a single population of size 10^4 diploids evolves neutrally until it splits into two
23 subpopulations at 4,000 generations before sampling. We selected the divergence and admixture times to
24 approximately match the timing of these events in sub-Saharan African populations [Veeramah et al., 2011,
25 Busby et al., 2016]. Sample sizes were of $n = 100$ diploids, matching the standard hard sweep experiments.

26 To simulate missing data in the sampled population, we followed a random approach. Using data
27 generated for the previous simple hard sweep experiment, we removed data from a random number of SNPs
28 in each replicate sample, between 25 and 50, drawing these sites from locations throughout the simulated
29 sequence uniformly at random. At each missing site, we assigned a number of the sampled individuals,
30 between 1 and 5, uniformly at random, to have their genotypes missing at the site. We then accounted for
31 missing data in one of two ways. First, we omitted any SNP with missing data in each analysis window. This
32 reduced the number of SNPs included in each computation. Second, we assigned any haplotype or MLG

1 with missing data as an entirely new string. Thus, the number of distinct haplotypes and MLGs increases
2 when sites are missing, providing a more conservative approach than the first.

3 Finally, our single scenario of background selection was intended to quantify the extent to which the
4 long-term removal of deleterious alleles in a population, which reduces nearby neutral genetic diversity,
5 would mislead each method to make false inferences of selective sweeps. We generated a 100 kb chromosome
6 containing an 11 kb gene at its center and allowed it to evolve over $12N$ generations under a constant-size
7 demographic model. The gene was composed of 10 exons of length 100 bases with 1 kb introns separating
8 each adjacent exon pair. The first and last exons were flanked by untranslated regions (UTRs) of length 200
9 bases at the 5' end and 800 bases at the 3' end. Strongly deleterious mutations ($s = -0.1$) arose at a rate of
10 50% in the UTRs, 75% in exons, and 10% in introns, while mutations occurring outside of the genic region
11 were neutral. To measure the confounding effect of background selection, we observed the overlap between
12 the distributions of maximum G12, G123, H12, and H123 values of 10^3 simulated replicates under neutrality
13 and background selection. Our model here is identical to that of Cheng et al. [2017], with the sizes of genetic
14 elements based on human mean values [Mignone et al., 2002, Sakharkar et al., 2004].

15 Detecting sweeps

16 We performed scans across simulated 100 kb and one Mb chromosomes with all methods using sliding genomic
17 windows of length 40 kb, advancing by four kb increments. We chose this window size primarily because the
18 mean value of LD between pairs of loci across the chromosome decays below one-third of its maximum value
19 over this interval (Figure S1), and because this was the window size with which we analyzed all non-African
20 populations from the 1000 Genomes dataset. Window size also affects sensitivity to sweeps by constraining
21 the minimum strength of selective sweeps we can detect. That is, with our chosen window size, we are likely
22 to detect sweeps with $s > 0.004$, because such sweeps will generate genomic footprints on the order of 40
23 kb for our simulated population size of $N = 10^4$. We computed this value as $F = s/(2r \ln(4Ns))$, where
24 F is the size of the footprint in nucleotides, s is the per-generation selection coefficient, r is the per-base,
25 per-generation recombination rate, and N is the effective population size [Gillespie, 2004, Garud et al., 2015,
26 Hermissen and Pennings, 2017].

27 For experiments measuring power at defined time points, we recorded the chromosomal maximum value
28 of G12, G123, H12, or H123 across all windows as the score for each of 10^3 replicates of 100 kb chromosomes.
29 Selection simulation scores provided us with a distribution of values that we compared with the distribution
30 of scores generated under neutral parameters. We define a method's power for each of our specified time
31 intervals at the 1% false positive rate (FPR). This measures the proportion of our 1,000 replicates generated
32 under selection parameters with a score greater than the top 1% of scores from the neutral replicates. The

1 method performs ideally if the distribution of its scores under a sweep does not overlap the distribution of
2 scores for neutral simulations; *i.e.*, if neutrality can never produce scores as large as a sweep.

3 In addition to power, we also tracked the mean scores of G12 and G123 across simulated one Mb
4 chromosomes at each 40 kb window for all selection scenarios at the time point for which power was greatest.

5 In situations where G12 or G123 had the same power at more than one time point (this occurred for strong
6 selection within 1,000 generations of sampling), we selected the most recent time point in order to represent
7 the maximum signal, since mutation and recombination erode expected haplotype homozygosity over time.

8 This analysis allowed us to observe the interval over which elevated scores are expected, and additionally
9 define the shape of the sweep signal.

10 Differentiating between selection scenarios

11 Experiments to test the ability of G2/G1 to correctly differentiate between soft and hard sweeps, as H2/H1
12 can (conditioning on a G12 or G123 value for G2/G1, or an H12 or H123 value for H2/H1), required a
13 different simulation approach than did the simple detection of selective sweeps. Whereas multiple methods
14 exist to identify sweeps from extended tracts of expected haplotype homozygosity, the method of Garud
15 et al. [2015] classifies this signal further to identify it as deriving from a soft or hard sweep. As did Garud
16 et al. [2015], we undertook an approximate Bayesian computation (ABC) approach to test the ability of
17 our method to distinguish soft and hard sweeps. To demonstrate the ability of G2/G1 conditional on G12
18 and G123 to differentiate between sweep scenarios and establish the basic properties of the (G12, G2/G1)
19 and (G123, G2/G1) distributions, we simulated sequences of length 40 kb under a constant population
20 size demographic history (Figure 2A) with a centrally-located site of selection. Here, we treated the whole
21 simulated sequence as a single window.

22 For ABC experiments to classify test points as hard or soft from a fixed number of different selected
23 haplotypes k , we performed 10^6 simulations for each selection scenario, drawing selection coefficients s and
24 selection times t uniformly at random from a log-scale as previously described. Soft sweeps from SSV were
25 generated for $k = 5$ and $k = 3$ starting haplotypes (out of a scaled $2N = 10^3$ haploids). Soft sweeps generated
26 under random t and s were compared with hard sweeps generated under random t and s , with completion of
27 the sweep possible but not guaranteed. From the resulting distribution of scores for each simulation type, we
28 computed Bayes factors (BFs) for direct comparisons between a hard sweep scenario and either soft sweep
29 scenario.

30 For two selection scenarios A and B and a (G12, G2/G1) or (G123, G2/G1) test point (or haplotype
31 statistic test point), we compute BFs as the number of simulations of type A yielding results within a
32 Euclidean distance of 0.1 from the test point, divided by the number of simulations of type B within that

1 distance. Here, test values of $(G_{12}, G_2/G_1)$ and $(G_{123}, G_2/G_1)$ are each plotted as a 100×100 grid, with
2 both dimensions bounded by $[0.005, 0.995]$ at increments of 0.01. In the work of Garud et al. [2015], soft
3 sweeps were of type *A* and hard sweeps were of type *B*, and we retain this orientation in our present work.
4 Following these definitions, a BF less than one at a test coordinate indicates that a hard sweep is more likely
5 to generate such a $(G_{12}, G_2/G_1)$ or $(G_{123}, G_2/G_1)$ pair, whereas a BF larger than one indicates greater
6 support for a recent soft sweep generating that value pair. As do Lee and Wagenmakers [2013], we define
7 $BF \geq 3$ as representing evidence for selection scenario *A* producing a similar paired $(G_{12}, G_2/G_1)$ or $(G_{123},$
8 $G_2/G_1)$ value as the test point, and $BF \geq 10$ to represent strong evidence. Similarly, $BF \leq 1/3$ is evidence in
9 favor of scenario *B*, and $BF \leq 1/10$ is strong evidence. We performed analyses for both MLG and haplotype
10 data to demonstrate the effect of data type on sweep type inference.

11 We followed a similar approach for ABC experiments to assign a most probable k to test points within
12 the aforementioned 100×100 grids. Here, we generated 5×10^6 replicates, drawing t and s uniformly at
13 random on a log scale as previously, and $k \in \{1, 2, \dots, 16\}$ uniformly at random. For each $(G_{12}, G_2/G_1)$,
14 $(G_{123}, G_2/G_1)$, $(H_{12}, H_2/H_1)$, or $(H_{123}, H_2/H_1)$ test point, we retained the value of k for each replicate
15 within a Euclidean distance of 0.1, and assigned the most frequently-occurring k as the most probable value
16 for the test point. Thus, unlike for BF experiments, no test point yielded an ambiguous result, and all test
17 points were assigned a most probable k .

18 Analysis of empirical data

19 We evaluated the ability of G_{12} , G_{123} , and H_{12} to corroborate and complement the results of existing
20 analyses on human data. Because G_{12} and G_{123} take unphased diploid MLGs as input, we manually
21 merged pairs of haplotype strings for this dataset (1000 Genomes Project, Phase 3 [Auton et al., 2015])
22 into MLGs, merging haplotype pairs that belonged to the same individual. We also complemented the
23 individual-centered approach by randomly merging pairs of haplotypes to produce a sample of individuals
24 that could arise under random mating. Our approaches therefore allowed us to determine the effect of using
25 different data types to infer selection. Unlike biallelic haplotypes, MLGs are triallelic, with an indicator
26 for each homozygous state and the heterozygous state. Thus, there are at least as many possible MLGs as
27 haplotypes, such that a sample with I distinct haplotypes can produce up to $I(I + 1)/2$ distinct MLGs.

28 We scanned all autosomes using nucleotide-delimited genomic windows, proportional to the effective size
29 of the study population, and the interval over which the rate of decay in pairwise LD plateaus empirically [see
30 Jakobsson et al., 2008]. For the 1000 Genomes YRI population, we employed a window of length 20 kb sliding
31 by increments of two kb, whereas for non-African populations (effective population size approximately half
32 of YRI) we used a window of 40 kb sliding by increments of five kb (see *Results*). This means that we were

1 sensitive to sweeps from approximately $s \geq 0.002$ for YRI, and approximately $s \geq 0.004$ for the others. We
2 recorded G12, G123, and H12 scores for all genomic windows, and subsequently filtered windows for which
3 the observed number of SNPs was less than a certain threshold value in order to avoid biasing our results
4 with heterochromatic regions for which sequence diversity is low in the absence of a sweep. Specifically, we
5 removed windows containing fewer SNPs than would be expected [Watterson, 1975] when two lineages are
6 sampled, which is the extreme case in which the selected allele has swept across all haplotypes except for one.
7 For our chosen genomic windows and all populations, this value is $4N_e\mu \times (\text{window size in nucleotides}) \approx 40$
8 SNPs, where N_e is the diploid effective population size and μ is the per-site per-generation mutation rate.
9 As in Huber et al. [2016], we additionally divided each chromosome into non-overlapping 100 kb bins and
10 removed sites within bins whose mean CRG100 score [Derrien et al., 2012], a measure of site mappability
11 and alignability, was less than 0.9. Filtering thereby removed additional sites for which variant calls were
12 unreliable, making no distinction between genic and non-genic regions.

13 Following a scan, we intersected selection signal peaks with the coordinates for protein- and RNA-coding
14 genes and generated a ranked list of all genomic hits discovered in the scan for each population. We used
15 the coordinates for human genome build hg19 for our data, to which Phase 3 of the 1000 Genomes Project
16 is mapped. The top 40 candidates for each study population were recorded and assigned *p*-values and
17 BFs. Specifically, we simulated sequences following the estimates of population size generated by Terhorst
18 et al. [2017] from smc++ using *ms* [Hudson, 2002] to assign *p*-values and SLiM 2 to assign BFs, with per-
19 generation, per-site mutation and recombination rates of 1.25×10^{-8} and 3.125×10^{-9} [Terhorst et al., 2017,
20 Narasimhan et al., 2017], and sample sizes for each population matching those of the 1000 Genomes Project.
21 For *p*-value simulations, we selected a sequence length uniformly at random from the set of all hg19 gene
22 lengths, appended the window size used for scanning that population's empirical data to this sequence, and
23 used a sliding window approach, retaining information from the window of maximum G12, G123, or H12
24 value. For BF simulations, we used simulated sequence lengths of either 20 kb for YRI or 40 kb for others,
25 to match the strategy of empirical scans. That is, once we have identified an elevated sweep signal within a
26 window, we then seek to classify it as hard or soft.

27 We assigned *p*-values by generating 10^6 replicates of neutrally-evolving sequences, where the *p*-value for
28 a gene is the proportion of maximum G12 (or G123 or H12) scores generated under neutrality that is greater
29 than the score assigned to that gene. After Bonferroni correction for multiple testing [Neyman and Pearson,
30 1928], a significant *p*-value was $p < 0.05/23,735 \approx 2.10659 \times 10^{-6}$, where 23,735 is the number of protein-
31 and RNA-coding genes for which we assigned a G12 (or G123 or H12) score. To assign BFs, we simulated
32 10^6 replicates of hard sweep and SSV ($k = 5$) scenarios for each study population (thus, 2×10^6 replicates
33 for each population), wherein the site of selection was at the center of the sequence. We drew $t \in [40, 2000]$

1 and $s \in [0.005, 0.5]$ uniformly at random from a log-scale, and defined BFs as previously. Additionally, we
2 assigned the most probable values of k from the posterior distribution for each top 40 sweep candidate for
3 each population, following the previous protocol. Values of t were chosen to reflect selective events within the
4 range of detection of G12, G123, and H12, while also being after the out-of-Africa event, whereas values of s
5 represent a range of selection strengths from weak to strong. We once again conditioned on the selected allele
6 remaining in the population throughout the simulation, though not on its frequency beyond this constraint.

7 We affirm that all data necessary for confirming the conclusions of the article are present within the
8 article, figures, and tables. Any other materials and resources are available upon request.

9 Acknowledgments

10 This work was supported by National Institutes of Health grant R35GM128590, by the Alfred P. Sloan
11 Foundation, and by Pennsylvania State University startup funds. We also thank Jonathan Terhorst for
12 providing demographic information on our study populations, estimated from his method `smc++`, as well
13 as Dmitri Petrov, Pleuni Pennings, and Arbel Harpak for helpful conversations. Finally, we thank three
14 anonymous reviewers for evaluating the merit of this work and providing comments that improved its over-
15 all quality. Portions of this research were conducted with Advanced CyberInfrastructure computational
16 resources provided by the Institute for CyberScience at Pennsylvania State University.

17 References

- 18 A Akbari, A Iranmehr, M Bakhtiari, S Mirarab, and V Bafna. Fine-mapping the Favored Mutation in a
19 Positive Selective Sweep. *bioRxiv*, pages 1–33, 2017.
- 20 M Ali, X Liu, E N Pillai, P Chen, C Khor, R T Ong, and Y Teo. Characterizing the genetic differences
21 between two distinct migrant groups from Indo-European and Dravidian speaking populations in India.
22 *BMC Genet.*, 15:86, 2014.
- 23 C E G Amorim, J T Daub, F M Salzano, M Foll, and L Excoffier. Detection of Convergent Genome-Wide
24 Signals of Adaptation to Tropical Forests in Humans. *PLoS ONE*, 10:e0121557, 2015.
- 25 A Auton, G R Abecasis, and The 1000 Genomes Project Consortium. A global reference for human genetic
26 variation. *Nature*, 526:68–74, 2015.
- 27 A Barsheshet, A Brenyo, A J Moss, and I Goldenberg. Genetics of Sudden Cardiac Death. *Curr. Cardiol.*
28 *Rep.*, 13:364–376, 2011.

- 1 R M Baxter, T Dai, J Kimball, E Wang, M R Hamblin, W P Wiesmann, S J McCarthy, and S M Baker.
2 Chitosan dressing promotes healing in third degree burns in mice: Gene expression analysis shows biphasic
3 effects for rapid tissue regeneration and decreased fibrotic signaling. *J. Biomed. Mater. Res. A*, 101:340–
4 348, 2013.
- 5 T M Baye, R A Wilke, and M Olivier. Genomic and geographic distribution of private SNPs and pathways
6 in human populations. *Pers. Med.*, 6:623–641, 2009.
- 7 S Beleza, A M Santos, B McEvoy, I Alves, C Martinho, E Cameron, M D Shriver, E J Parra, and J Rocha.
8 The Timing of Pigmentation Lightening in Europeans. *Mol. Biol. Evol.*, 30:24–35, 2012.
- 9 J J Berg and G Coop. A Coalescent Model for a Sweep of a Unique Standing Variant. *Genetics*, 201:707–725,
10 2015.
- 11 T Bersaglieri, P C Sabeti, N Patterson, T Vanderploeg, S F Schaffner, J A Drake, M Rhodes, D E Reich,
12 and J N Hirschhorn. Genetic Signatures of Strong Recent Positive Selection at the Lactase Gene. *Am. J.
13 Hum. Genet.*, 74:1111–1120, 2004.
- 14 G Bhatia, N Patterson, B Pasaniuc, N Zaitlen, G Genovese, S Pollack, S Mallick, S Myers, A Tandon,
15 C Spencer, C D Palmer, A A Adeyemo, E L Akylbekova, L A Cupples, J Divers, M Fornage, W H L
16 Kao, L Lange, M Li, S Musani, J C Mychaleckyj, A Ogunniyi, G Papanicolaou, C N Rotimi, J I Rotter,
17 I Ruczinski, B Salako, D S Siscovick, B O Tayo, Q Yang, S McCarroll, P Sabeti, G Lettre, P De Jager,
18 J Hirschhorn, X Zhu, R Cooper, D Reich, J G Wilson, and A L Price. Genome-wide Comparison of
19 African-Ancestry Populations from CARE and Other Cohorts Reveals Signals of Natural Selection. *Am.
20 J. Hum. Genet.*, 89:368–381, 2011.
- 21 A Blant, M Kwong, Z A Szpiech, and T J Pemberton. Weighted likelihood inference of genomic autozygosity
22 patterns in dense genotype data. *BMC Genomics*, 18:928, 2017.
- 23 S R Browning and B L Browning. Haplotype phasing: existing methods and new developments. *Nat. Rev.
24 Genet.*, 12:703–714, 2011.
- 25 J Bryk, E Hardouin, I Pugach, D Hughes, R Strotmann, M Stoneking, and S Myles. Positive Selection in
26 East Asians for an *EDAR* Allele that Enhances NF- κ B Activation. *PLoS ONE*, 3:e2209, 2008.
- 27 G B J Busby, G Band, Q S Le, M Jallow, E Bougama, V D Mangano, L N Amenga-Etego, A Enimil,
28 T Apinjoh, C M Ndila, A Manjurano, V Nyirongo, O Doumba, K A Rockett, D P Kwiatkowski, C C A
29 Spencer, and Malaria Genomic Epidemiology Network. Admixture into and within sub-Saharan Africa.
30 *eLife*, 5:e15266, 2016.

- 1 M C Campbell and S A Tishkoff. African Genetic Diversity: Implications for Human Demographic History,
2 Modern Human Origins, and Complex Disease Mapping. *Annu. Rev. Genom. Hum. G.*, 9:403–433, 2008.
- 3 S E Castel, P Mohammadi, W K Chung, Y Shen, and T Lappalainen. Rare variant phasing and haplotypic
4 expression from RNA sequencing with phASER. *Nat. Commun.*, 7:12817, 2016.
- 5 B V S K Chakravarthi, S S Pathi, M T Goswami, M Cieślik, H Zheng, S Nallasivam, S R Arekapudi, X Jing,
6 J Siddiqui, J Athanikar, S L Carskadon, R J Lonigro, L P Kunju, A M Chinnayan, N Palanisamy, and
7 S Varamballi. The miR-124-Prolyl Hydroxylase P4HA1-MMP1 axis plays a critical role in prostate cancer
8 progression. *Oncotarget*, 5:6654–6669, 2014.
- 9 Y C Chang, X Liu, J D O Kim, M A Ikeda, M R Layton, A B Weder, R S Cooper, S L R Kardia, D C
10 Rao, S C Hunt, A Luke, E Boerwinkle, and A Chakravarti. Multiple Genes for Essential-Hypertension
11 Susceptibility on Chromosome 1q. *Am. J. Hum. Genet.*, 80:253–264, 2007.
- 12 B Charlesworth. The Effects of Deleterious Mutations on Evolution at Linked Sites. *Genetics*, 190:5–22,
13 2012.
- 14 B Charlesworth, M T Morgan, and D Charlesworth. The Effect of Deleterious Mutations on Neutral Molecular
15 Variation. *Genetics*, 134:1289–1303, 1993.
- 16 B Charlesworth, D Charlesworth, and M T Morgan. The Pattern of Neutral Molecular Variation Under the
17 Background Selection Model. *Genetics*, 141:1619–1632, 1995.
- 18 H Chen, N J Patterson, and D E Reich. Population differentiation as a test for selective sweeps. *Genome Res.*, 20:393–402, 2010.
- 19 H Chen, J Hey, and M Slatkin. A hidden Markov model for investigating recent positive selection through
20 haplotype structure. *Theor. Popul. Biol.*, 99:18–30, 2015.
- 21 X Cheng, C Xu, and M DeGiorgio. Fast and robust detection of ancestral selective sweeps. *Mol. Ecol.*, 2017.
22 doi: 10.1111/mec.14416.
- 23 E E Chukwu, F O Nwaokorie, A O Coker, M J Avila-Campos, R L Solis, L A Llanco, and F T Ogunsola.
24 Detection of toxigenic *Clostridium perfringens* and *Clostridium botulinum* from food sold in Lagos, Nigeria.
25 *Anaerobe*, 42:176–181, 2016.
- 26 F J Clemente, A Cardona, C E Inchley, B M Peter, G Jacobs, L Pagani, D J Lawson, T Antão, M Vicente,
27 M Mitt, M DeGiorgio, Z Faltyskova, Y Xue, Q Ayub, M Szpak, R Mägi, A Eriksson, A Manica, M Raghavan,
28 M Rasmussen, S Rasmussen, E Willerslev, A Vidal-Puig, C Tyler-Smith, R Villemans, R Nielsen,

- 1 M Metspalu, B Malyarchuk, M Derenko, and T Kivisild. A Selective Sweep on a Deleterious Mutation in
2 *CPT1A* in Arctic Populations. *Am. J. Hum. Genet.*, 95:584–589, 2014.
- 3 J M Comeron. Background Selection as Baseline for Nucleotide Variation across the *Drosophila* Genome.
4 *PLoS Genet.*, 10:e1004434, 2014.
- 5 C Connan, M Voillequin, C V Chavez, C Mazuet, C Levesque, S Vitry, A Vandewalle, and M R Popoff.
6 Botulinum neurotoxin type B uses a distinct entry pathwaymediated by CDC42 into intestinal cells versus
7 neuronal cells. *Cell. Microbiol.*, 19:e12738, 2017.
- 8 G Coop, J K Pickrell, J Novembre, S Kudaravalli, J Li, D Absher, R M Myers, L L Cavalli-Sforza, M W
9 Feldman, and J K Pritchard. The Role of Geography in Human Adaptation. *PLoS Genet.*, 5:e1000500,
10 2009.
- 11 A D Cutter and B A Payseur. Genomic signatures of selection at linked sites: unifying the disparity among
12 species. *Nat. Rev. Genet.*, 14:262–274, 2013.
- 13 H D Daetwyler, A Capitan, H Pausch, P Stothard, R van Binsbergen, R F Brøndum, X Liao, A Djari, S C
14 Rodriguez, C Grohs, et al. Whole-genome sequencing of 234 bulls facilitates mapping of monogenic and
15 complex traits in cattle. *Nat. Genet.*, 46:858–865, 2014.
- 16 M DeGiorgio, M Jakobsson, and N A Rosenberg. Explaining worldwide patterns of human genetic variation
17 using a coalescent-based serial founder model of migration outward from Africa. *Proc. Natl. Acad. Sci.*
18 U.S.A., 106:16057–16062, 2009.
- 19 M DeGiorgio, K E Lohmueller, and R Nielsen. A Model-Based Approach for Identifying Signatures of
20 Ancient Balancing Selection in Genetic Data. *PLoS Genet.*, 10:e1004561, 2014.
- 21 T Derrien, J Estellé, S M Sola, D G Knowles, E Rainieri, R Guigó, and P Ribeca. Fast Computation and
22 Applications of Genome Mappability. *PLoS ONE*, 7:e30377, 2012.
- 23 C Drury, K E Dale, J M Panlilio, S V Miller, D Lirman, E A Larson, E Bartels, D L Crawford, and
24 M F Oleksiak. Genomic variation among populations of threatened coral: *Acropora cervicornis*. *BMC*
25 *Genomics*, 17:286, 2011.
- 26 R M Durbin and The 1000 Genomes Project Consortium. A map of human genome variation from population-
27 scale sequencing. *Nature*, 467:1061–1073, 2011.
- 28 C J Edwards, C Ginja, J Kantanen, L Pérez-Pardal, A Tresset, F Stock, European Cattle Genetic Diversity
29 Consortium, L T Gama, M C T Penedo, D G Bradley, J A Lenstra, and I J Nijman. Dual Origins of

- 1 Dairy Cattle Farming – Evidence from a Comprehensive Survey of European Y-Chromosomal Variation.
2 *PLoS ONE*, 6:e15922, 2011.
- 3 R J Elshire, J C Glaubitz, Q Sun, J A Poland, K Kawamoto, E S Buckler, and S E Mitchell. A Robust,
4 Simple Genotyping-by-Sequencing (GBS) Approach for High Diversity Species. *PLoS ONE*, 6:e19379,
5 2011.
- 6 D Enard, P W Messer, and D A Petrov. Genome-wide signals of positive selection in human evolution.
7 *Genome Res.*, 24:885–895, 2014.
- 8 L Ermini, C D Sarkissian, E Willerslev, and L Orlando. Major transitions in human evolution revisited: A
9 tribute to ancient DNA. *J. Hum. Evol.*, 79:4–20, 2015.
- 10 C Evers, B Maas, K A Koch, A Jauch, J W G Janssen, C Sutter, M J Parker, K Hinderhofer, and U Moog.
11 Mosaic Deletion of EXOC6B: Further Evidence for An Important Role of the Exocyst Complex in the
12 Pathogenesis of Intellectual Disability. *Am. J. Med. Genet. Part A*, 164:3088–3094, 2014.
- 13 M Fagny, E Patin, D Enard, L B Barreiro, L Quintana-Murci, and G Laval. Exploring the Occurrence of
14 Classic Selective Sweeps in Humans Using Whole-Genome Sequencing Data Sets. *Mol. Biol. Evol.*, 31:
15 1850–1868, 2014.
- 16 A Ferrer-Admetlla, M Liang, T Korneliussen, and R Nielsen. On Detecting Incomplete Soft or Hard Selective
17 Sweeps Using Haplotype Structure. *Mol. Biol. Evol.*, 31:1275–1291, 2014.
- 18 R A Fisher. *The Genetical Theory of Natural Selection*. Oxford University Press, Inc., Clarendon, Oxford,
19 1st edition, 1930.
- 20 A Fujimoto, R Kimura, J Ohashi, K Omi, R Yuliwulandari, L Batubara, M S Mustofa, U Samakkarn,
21 W Settheetham-Ishida, T Ishida, Y Morishita, T Furusawa, M Nakazawa, R Ohtsuka, and K Tokunaga.
22 A scan for genetic determinants of human hair morphology: *EDAR* is associated with Asian hair thickness.
23 *Hum. Mol. Genet.*, 17:835–843, 2007.
- 24 M Fumagalli, F G Vieira, T Linderroth, and R Nielsen. *ngsTools*: methods for population genetics analyses
25 from next-generation sequencing data. *Bioinformatics*, 30:1486–1487, 2014.
- 26 A A Fushan, A A Turanov, S Lee, E B Kim, A V Lobanov, S H Yim, R Buffenstein, S Lee, K Chang,
27 H Rhee, J Kim, K Yang, and V N Gladyshev. Gene expression defines natural changes in mammalian
28 lifespan. *Aging Cell*, 14:352–365, 2015.

- 1 N R Garud and N A Rosenberg. Enhancing the mathematical properties of new haplotype homozygosity
2 statistics for the detection of selective sweeps. *Theor. Popul. Biol.*, 102:94–101, 2015.
- 3 N R Garud, P W Messer, E O Buzbas, and D A Petrov. Recent Selective Sweeps in North American
4 *Drosophila melanogaster* Show Signatures of Soft Sweeps. *PLoS Genet.*, 11:e1005004, 2015.
- 5 T J Gauvin, L E Young, and H N Higgs. The formin *FMNL3* assembles plasma membrane protrusions that
6 participate in cell–cell adhesion. *Mol. Biol. Cell*, 26:467–477, 2014.
- 7 P Gerbault, C Moret, M Currat, and A Sanchez-Mazas. Impact of Selection and Demography on the
8 Diffusion of Lactase Persistence. *PLoS ONE*, 4:e6369, 2009.
- 9 J H Gillespie. *Population Genetics: A Concise Guide*. The Johns Hopkins University Press, Baltimore, MD,
10 2nd edition, 2004.
- 11 S R Grossman, I Shylakhter, E K Karlsson, E H Byrne, S Morales, G Frieden, E Hostetter, E Angelino,
12 M Garber, O Zuk, E S Lander, S F Schaffner, and P C Sabeti. A Composite of Multiple Signals Distinguishes Causal Variants in Regions of Positive Selection. *Science*, 327:883–886, 2010.
- 14 B C Haller and P W Messer. SLiM 2: Flexible, Interactive Forward Genetic Simulations. *Mol. Biol. Evol.*,
15 34:230–240, 2017.
- 16 D L Hartl and A G Clark. *Principles of Population Genetics*. Sinauer Associates, Inc., Sunderland MA, 4th
17 edition, 2007.
- 18 L He, J Pitkäniemi, A Sarin, V Salomaa, M J Sillanpää, and S Ripatti. Hierarchical Bayesian Model for
19 Rare Variant Association Analysis Integrating Genotype Uncertainty in Human Sequence Data. *Genet. Epidemiol.*, 39:89–100, 2014.
- 21 J Hermisson and P S Pennings. Soft Sweeps: Molecular Population Genetics of Adaptation From Standing
22 Genetic Variation. *Genetics*, 169:2335–2352, 2005.
- 23 J Hermisson and P S Pennings. Soft sweeps and beyond: understanding the patterns and probabilities of
24 selection footprints under rapid adaptation. *Methods Ecol. Evol.*, 8:700–716, 2017.
- 25 C Hetheridge, A N Scott, R K Swain, J W Copeland, H N Higgs, R Bicknell, and H Mellor. The formin
26 *FMNL3* is a cytoskeletal regulator of angiogenesis. *J. Cell Sci.*, 125:1420–1428, 2012.
- 27 C D Huber, M DeGiorgio, I Hellmann, and R Nielsen. Detecting recent selective sweeps while controlling
28 for mutation rate and background selection. *Mol. Ecol.*, 25:142–156, 2016.

- 1 R R Hudson. Generating samples under a Wright-Fisher neutral model of genetic variation. *Bioinformatics*,
2 18:337–338, 2002.
- 3 Y Itan, A Powell, M A Beaumont, J Burger, and M G Thomas. The Origins of Lactase Persistence in
4 Europe. *PLoS Comput. Biol.*, 5:e1000491, 2009.
- 5 M Jakobsson, S W Scholz, P Scheet, J R Gibbs, J M VanLiere, H Fung, Z A Szpiech, J H Degnan, K Wang,
6 R Guerreiro, J M Bras, J C Schymick, D G Hernandez, B J Traynor, J Simon-Sanchez, M Matarin,
7 A Britton, J van de Leemput, I Rafferty, M Bucan, H M Cann, J A Hardy, N A Rosenberg, and A B
8 Singleton. Genotype, haplotype and copy-number variation in worldwide human populations. *Nature*,
9 451:998–1003, 2008.
- 10 B L Jones, T O Raga, A Liebert, P Zmarz, E Bekele, E T Danielson, A K Olsen, N Bradman, J T Troelsen,
11 and D M Swallow. Diversity of Lactase Persistence Alleles in Ethiopia: Signature of a Soft Selective
12 Sweep. *Am. J. Hum. Genet.*, 93:538–544, 2013.
- 13 J K Kelly. A Test of Neutrality Based on Interlocus Associations. *Genetics*, 146:1197–1206, 1997.
- 14 A D Kern and D R Schrider. diploS/HIC: An Updated Approach to Classifying Selective Sweeps. *G3-Genes*
15 *Genom. Genet.*, 8:1959–1970, 2018.
- 16 Y Kim and R Nielsen. Linkage Disequilibrium as a Signature of Selective Sweeps. *Genetics*, 167:1513–1524,
17 2004.
- 18 Y Kim and W Stephan. Detecting a Local Signature of Genetic Hitchhiking Along a Recombining Chromo-
19 some. *Genetics*, 160:765–777, 2002.
- 20 T S Korneliussen, A Albrechtsen, and R Nielsen. ANGSD: Analysis of Next Generation Sequencing Data.
21 *BMC Bioinformatics*, 15:356, 2014.
- 22 R L Lamason, M P K Mohideen, J R Mest, A C Wong, H L Norton, M C Aros, M J Juryne, X Mao, V R
23 Humphreville, J E Humbert, S Sinha, J L Moore, P Jagadeeswaran, W Zhao, G Ning, I Makalowska, P M
24 McKeigue, D O’Donnell, R Kittles, E J Parra, N J Mangini, D J Grunwald, M D Shriver, V A Canfield,
25 and K C Cheng. SLC24A5, a Putative Cation Exchanger, Affects Pigmentation in Zebrafish and Humans.
26 *Science*, 310:1782–1786, 2005.
- 27 T W Laver, R C Caswell, K A Moore, J Poschmann, M B Johnson, M M Owens, S Ellard, K H Paszkiewicz,
28 and M N Weedon. Pitfalls of haplotype phasing from amplicon-based long-read sequencing. *Sci. Rep.-U.K.*,
29 6:21746, 2016.

- ¹ M D Lee and E Wagenmakers. *Bayesian Cognitive Modeling: A Practical Course*. Cambridge University
² Press, Cambridge U.K., 1st edition, 2013.
- ³ K Lin, H Li, C Schlötterer, and A Futschik. Distinguishing Positive Selection From Neutral Evolution:
⁴ Boosting the Performance of Summary Statistics. *Genetics*, 187:229–244, 2011.
- ⁵ X Liu, R T Ong, E N Pillai, A M Elzein, K S Small, T G Clark, D P Kwiatowski, and Y Teo. Detecting and
⁶ Characterizing Genomic Signatures of Positive Selection in Global Populations. *Am. J. Hum. Genet.*, 92:
⁷ 866–881, 2013.
- ⁸ K E Lohmueller, C D Bustamante, and A G Clark. Methods for Human Demographic Inference Using
⁹ Haplotype Patterns From Genomewide Single-Nucleotide Polymorphism Data. *Genetics*, 182:217–231,
¹⁰ 2009.
- ¹¹ C B Mallick, F M Iliescu, M Möls, S Hill, R Tamang, G Chaubey, R Goto, S Y W Ho, I G Romero,
¹² F Crivellaro, G Hudjashov, N Rai, M Metspalu, C G N Mascie-Taylor, R Pitchappan, L Singh, M Mirazon-
¹³ Lahr, K Thangaraj, R Villemans, and T Kivisild. The Light Skin Allele of SLC24A5 in South Asians and
¹⁴ Europeans Shares Identity by Descent. *PLoS Genet.*, 9:e1003912, 2013.
- ¹⁵ J Marchini and B Howie. Genotype imputation for genome-wide association studies. *Nat. Rev. Genet.*, 11:
¹⁶ 499–511, 2010.
- ¹⁷ B J Maron, K P Carney, H M Lever, J F Lewis, I Barac, S A Casey, and M V Sherrid. Relationship of Race
¹⁸ to Sudden Cardiac Death in Competitive Athletes With Hypertrophic Cardiomyopathy. *J. Am. Coll.
Cardiol.*, 41:974–980, 2003.
- ²⁰ J Maynard Smith and J Haigh. The hitch-hiking effect of a favorable gene. *Genet. Res.*, 23:23–35, 1974.
- ²¹ G McVicker, D Gordon, C Davis, and P Green. Widespread Genomic Signatures of Natural Selection in
²² Hominid Evolution. *PLoS Genet.*, 5:e1000471, 2009.
- ²³ I Mendizabal, U M Marigorta, O Lao, and D Comas. Adaptive evolution of loci covarying with the human
²⁴ African Pygmy phenotype. *Hum. Genet.*, 131:1305–1317, 2012.
- ²⁵ P W Messer. SLiM: Simulating Evolution with Selection and Linkage. *Genetics*, 194:1037–1039, 2013.
- ²⁶ F Mignone, C Gissi, S Liuni, and G Pesole. Untranslated regions of mRNAs. *Genome Biol.*, 3:reviews0004–1,
²⁷ 2002.
- ²⁸ M R Mughal and M DeGiorgio. Localizing and classifying adaptive targets with trend filtered regression.
²⁹ *BioRxiv*, 2018. doi: 10.1101/320523.

- 1 M W Nachman and S L Crowell. Estimate of the mutation rate per nucleotide in humans. *Genetics*, 156:
2 297–304, 2000.
- 3 V M Narasimhan, R Rahbari, A Scally, A Wuster, D Mason, Y Xue, J Wright, R C Trembath, E R Maher,
4 D A van Heel, A Auton, M E Hurles, C Tyler-Smith, and R Durbin. Estimating the human mutation rate
5 from autozygous segments reveals population differences in human mutational processes. *Nat. Commun.*,
6 8, 2017. doi: 10.1038/s41467-017-00323-y.
- 7 J Neyman and E S Pearson. On the Use and Interpretation of Certain Test Criteria for Purposes of Statistical
8 Inference: Part I. *Biometrika*, 20A:175–240, 1928.
- 9 L E Nicolaisen and M M Desai. Distortions in Genealogies due to Purifying Selection and Recombination.
10 *Genetics*, 195:221–230, 2013.
- 11 R Nielsen, S Williamson, Y Kim, M J Hubisz, A G Clark, and C Bustamante. Genomic scans for selective
12 sweeps using SNP data. *Genome Res.*, 15:1566–1575, 2005.
- 13 R Nielsen, J S Paul, A Albrechtsen, and Y S Song. Genotype and SNP calling from next-generation
14 sequencing data. *Nat. Rev. Genet.*, 12:443–451, 2011.
- 15 J O’Connell, D Gurdasani, O Delaneau, N Pirastu, S Ulivi, M Cocca, M Traglia, J Huang, J E Huffman,
16 I Rudan, R McQuillan, R M Fraser, H Campbell, O Polasek, G Asiki, K Ekoru, C Hayward, A F Wright,
17 V Vitart, P Navarro, J Zagury, J F Wilson, D Toniolo, P Gasparini, N Soranzo, M S Sandhu, and
18 J Marchini. A General Approach for Haplotype Phasing across the Full Spectrum of Relatedness. *PLoS
19 Genet.*, 10:e1004234, 2014.
- 20 J Ohashi, I Naka, and N Tsuchiya. The Impact of Natural Selection on an *ABCC11* SNP Determining
21 Earwax Type. *Mol. Biol. Evol.*, 28:849–857, 2010.
- 22 P Pavlidis, J D Jensen, and W Stephan. Searching for Footprints of Positive Selection in Whole-Genome
23 SNP Data From Nonequilibrium Populations. *Genetics*, 185:907–922, 2010.
- 24 B A Payseur and M W Nachman. Microsatellite Variation and Recombination Rate in the Human Genome.
25 *Genetics*, 156:1285–1298, 2000.
- 26 T J Pemberton, D Absher, M W Feldman, R M Myers, N A Rosenberg, and J Z Li. Genomic Patterns of
27 Homozygosity in Worldwide Human Populations. *Am. J. Hum. Genet.*, 91:275–292, 2012.
- 28 P S Pennings and J Hermisson. Soft Sweeps II: Molecular Population Genetics of Adaptation from Recurrent
29 Mutation or Migration. *Mol. Biol. Evol.*, 23:1076–1084, 2006a.

- ¹ P S Pennings and J Hermisson. Soft Sweeps III: The Signature of Positive Selection from Recurrent Mutation.
- ² *PLoS Genet.*, 2:e186, 2006b.
- ³ N Petousi, Q P P Croft, G L Cavalleri, H Cheng, F Formenti, K Ishida, D Lunn, M McCormack, K V Shianna, N P Talbot, P J Ratcliffe, and P A Robbins. Tibetans living at sea level have a hyporesponsive hypoxia-inducible factor system and blunted physiological responses to hypoxia. *J. Appl. Physiol.*, 116: 893–904, 2013.
- ⁷ J K Pickrell, G Coop, J Novembre, S Kudaravalli, J Z Li, D Absher, B S Srinivasan, G S Barsh, R M Myers, M W Feldman, and J K Pritchard. Signals of recent positive selection in a worldwide sample of human populations. *Genome Res.*, 19:826–837, 2009.
- ¹⁰ D Pierron, H Razafindrazaka, L Pagani, F Ricaut, T Antao, M Capredon, C Sambo, C Radimilahy, J Rakotoarisoa, R M Blench, T Letellier, and T Kivisild. Genome-wide evidence of Austronesian-Bantu admixture and cultural reversion in a hunter-gatherer group of Madagascar. *Proc. Natl. Acad. Sci. U.S.A.*, 111: 936–941, 2014.
- ¹⁴ M Przeworski. The Signature of Positive Selection at Randomly Chosen Loci. *Genetics*, 160:1179–1189, 2002.
- ¹⁶ M Przeworski, G Coop, and J D Wall. The Signature of Positive Selection on Standing Genetic Variation. *Evolution*, 59:2312–2323, 2005.
- ¹⁸ M Pybus, G M Dall’Olio, P Luisi, M Uzkudun, A Carreño-Torres, P Pavlidis, H Laayouni, J Bertranpetti, and J Engelken. 1000 Genomes Selection Browser 1.0: a genome browser dedicated to signatures of natural selection in modern humans. *Nucleic Acids Res.*, 42:D903–D909, 2014.
- ²¹ M Pybus, P Luisi, G M Dall’Olio, M Uzkudun, H Laayouni, J Bertranpetti, and J Engelken. Hierarchical boosting: a machine-learning framework to detect and classify hard selective sweeps in human populations. *Bioinformatics*, 31:3946–3952, 2015.
- ²⁴ F Racimo. Testing for Ancient Selection Using Cross-population Allele Frequency Differentiation. *Genetics*, 202:733–750, 2016.
- ²⁶ R Ronen, N Udpa, E Halperin, and V Bafna. Learning Natural Selection from the Site Frequency Spectrum. *Genetics*, 195:181–193, 2013.
- ²⁸ R Ronen, G Tesler, A Akbari, S Zakov, N A Rosenberg, and V Bafna. Predicting Carriers of Ongoing Selective Sweeps without Knowledge of the Favored Allele. *PLoS Genet.*, 11:e1005527, 2015.

- 1 P C Sabeti, D E Reich, J M Higgins, H Z P Levine, D J Richter, S F Schaffner, S B Gabriel, J V Platko,
2 N J Patterson, G J McDonald, H C Ackerman, S J Campbell, D Altshuler, R Cooper, D Kwiatkowski,
3 R Ward, and E S Lander. Detecting recent positive selection in the human genome from haplotype
4 structure. *Nature*, 419:832–837, 2002.
- 5 P C Sabeti, P Varilly, B Fry, J Lohmueller, E Hostetter, C Cotsapas, X Xie, E H Byrne, S A McCarroll,
6 R Gaudet, S F Schaffner, E S Lander, and The International HapMap Consortium. Genome-wide detection
7 and characterization of positive selection in human populations. *Nature*, 449:913–918, 2007.
- 8 M K Sakharkar, V T K Chow, and P Kangueane. Distributions of exons and introns in the human genome.
9 *In Silico Biol.*, 4:387–393, 2004.
- 10 F Schlamp, J van der Made, R Stambler, L Chesebrough, A R Boyko, and P W Messer. Evaluating the
11 performance of selection scans to detect selective sweeps in domestic dogs. *Mol. Ecol.*, 25:342–356, 2016.
- 12 D R Schrider and A D Kern. S/HIC: Robust Identification of Soft and Hard Sweeps Using Machine Learning.
13 *PLoS Genet.*, 12:e1005928, 2016.
- 14 D R Schrider and A D Kern. Soft Sweeps Are the Dominant Mode of Adaptation in the Human Genome.
15 *Mol. Biol. Evol.*, 34:1863–1877, 2017.
- 16 H Schulze, M Dose, M Korpal, I Meyer, J E Jr Italiano, and R A Shivdasani. RanBP10 Is a Cytoplasmic
17 Guanine Nucleotide Exchange Factor That Modulates Noncentrosomal Microtubules. *J. Biol. Chem.*, 283:
18 14109–14119, 2008.
- 19 J Schweinsberg and R Durrett. Random Partitions Approximating the Coalescence of Lineages During a
20 Selective Sweep. *Ann. Appl. Probab.*, 15:1591–1651, 2005.
- 21 J Seger, W A Smith, J J Perry, J Hunn, Z A Kaliszewska, L La Sala, L Pozzi, V J Rountree, and F R
22 Adler. Gene Genealogies Strongly Distorted by Weakly Interfering Mutations in Constant Environments.
23 *Genetics*, 184:529–545, 2010.
- 24 S Sheehan and Y S Song. Deep Learning for Population Genetic Inference. *PLoS Comput. Biol.*, 12:e1004845,
25 2016.
- 26 P Sinnock. The Wahlund Effect For The Two-Locus Model. *Am. Nat.*, 109:565–570, 1975.
- 27 M Slatkin. Linkage disequilibrium — understanding the evolutionary past and mapping the medical future.
28 *Nat. Rev. Genet.*, 9:477–485, 2008.

- ¹ N Takahata, Y Satta, and J Klein. Divergence Time and Population Size in the Lineage Leading to Modern
² Humans. *Theor. Popul. Biol.*, 48:198–221, 1995.
- ³ J Terhorst, J A Kamm, and Y S Song. Robust and scalable inference of population history from hundreds
⁴ of unphased whole genomes. *Nat. Genet.*, 49:303–309, 2017.
- ⁵ The International HapMap Consortium. A second generation human haplotype map of over 3.1 million
⁶ SNPs. *Nature*, 449:851–861, 2007.
- ⁷ K R Veeramah, D Wegmann, A Woerner, F L Mendez, J C Watkins, G Destro-Bisol, H Soodyall, L Louie,
⁸ and M F Hammer. An Early Divergence of KhoeSan Ancestors from Those of Other Modern Humans Is
⁹ Supported by an ABC-Based Analysis of Autosomal Resequencing Data. *Mol. Biol. Evol.*, 29:617–630,
¹⁰ 2011.
- ¹¹ B F Voight, S Kudaravalli, X Wen, and J K Pritchard. A Map of Recent Positive Selection in the Human
¹² Genome. *PLoS Biol.*, 4:e72, 2006.
- ¹³ H M T Vy and Y Kim. A Composite-Likelihood Method for Detecting Incomplete Selective Sweep from
¹⁴ Population Genomic Data. *Genetics*, 200:633–649, 2015.
- ¹⁵ G A Watterson. On the Number of Segregating Sites in Genetical Models without Recombination. *Theor.
Popul. Biol.*, 7:256–276, 1975.
- ¹⁷ S Wright. Evolution in Mendelian Populations. *Genetics*, 16:97–159, 1931.
- ¹⁸ F Zhang, L Christiansen, J Thomas, D Pokholok, R Jackson, N Morrell, Y Zhao, M Wiley, E Welch, E Jaeger,
¹⁹ A Granat, S J Norberg, A Halpern, M C Rogert, M Ronaghi, J Shendure, N Gormley, K L Gunderson,
²⁰ and F J Steemers. Haplotype phasing of whole human genomes using bead-based barcode partitioning in
²¹ a single tube. *Nat. Biotechnol.*, 35, 2017.
- ²² F Zhu, Q Cui, and Z Hou. SNP discovery and genotyping using Genotyping-by-Sequencing in Pekin ducks.
²³ *Sci. Rep.-U.K.*, 6, 2016.

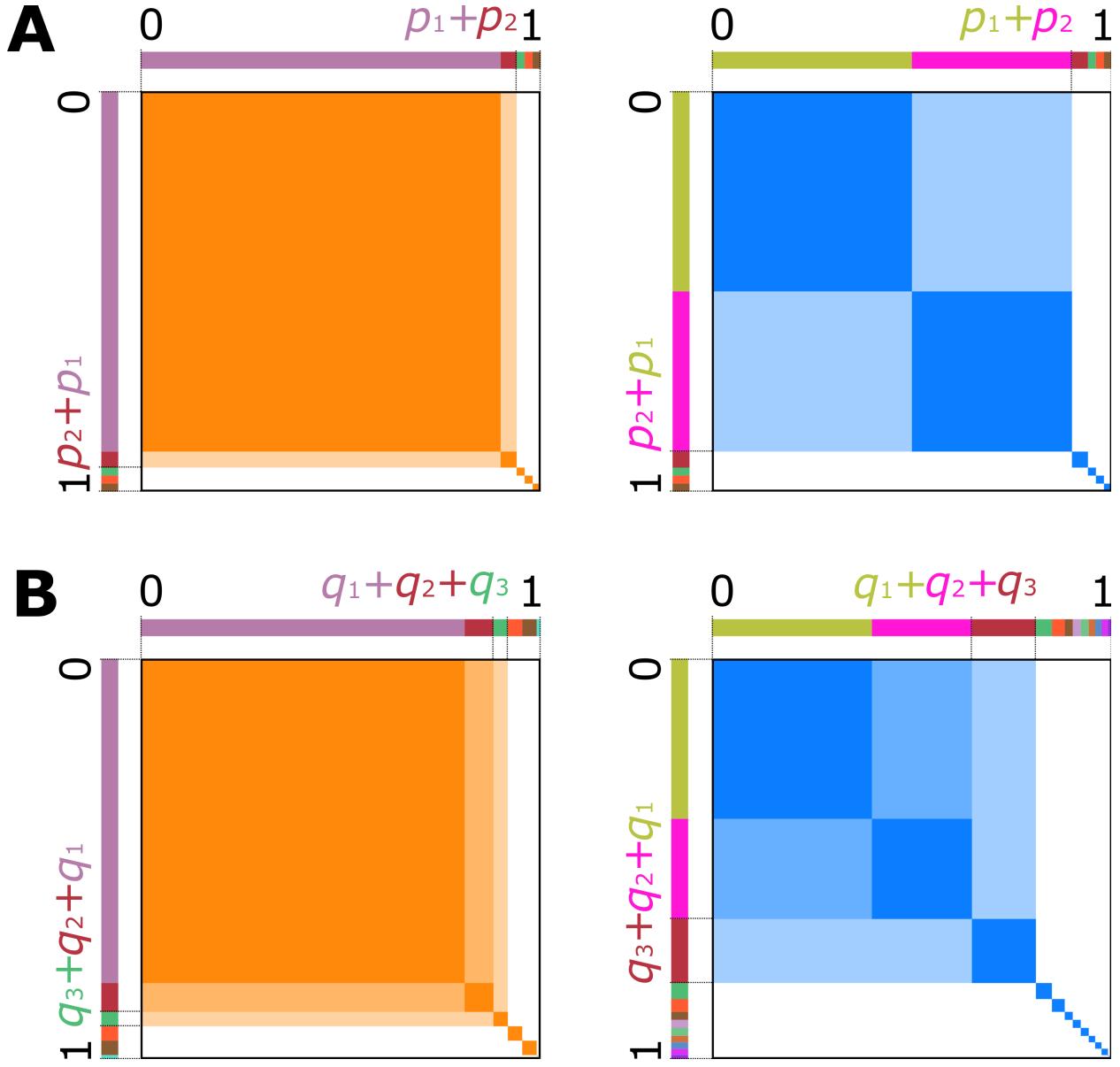


Figure 1: Visual representation of expected homozygosity statistics. For all panels, total area of the orange or blue squares within a panel represents the value of expected homozygosity statistics. Hard sweep scenarios are in orange, and soft sweeps are in blue. (A) Under a hard sweep (left), a single haplotype rises to high frequency, p_1 , so the probability of sampling two copies of that haplotype is p_1^2 . Choosing p_1 as the largest frequency yields H1 (dark orange area), while pooling $p_1 + p_2$ as the largest frequency yields H12 (total orange area). Under a soft sweep (right), pooling the largest haplotype frequencies results in a large shaded area, and therefore H12 has a similar value for both hard and soft sweeps. (B) Under Hardy Weinberg equilibrium, a single high-frequency haplotype produces a single high-frequency MLG (frequency q_1). Pooling frequencies up to q_3 has little effect on the value of the statistic, thus G1, G12, and G123 have similar values. When two haplotypes exist at high frequency, three MLGs exist at high frequency. Under a soft sweep, pooling the largest two MLGs (G12) may provide greater resolution of soft sweeps than not pooling (G1), and pooling the largest three creates a statistic (G123) truly analogous to H12.

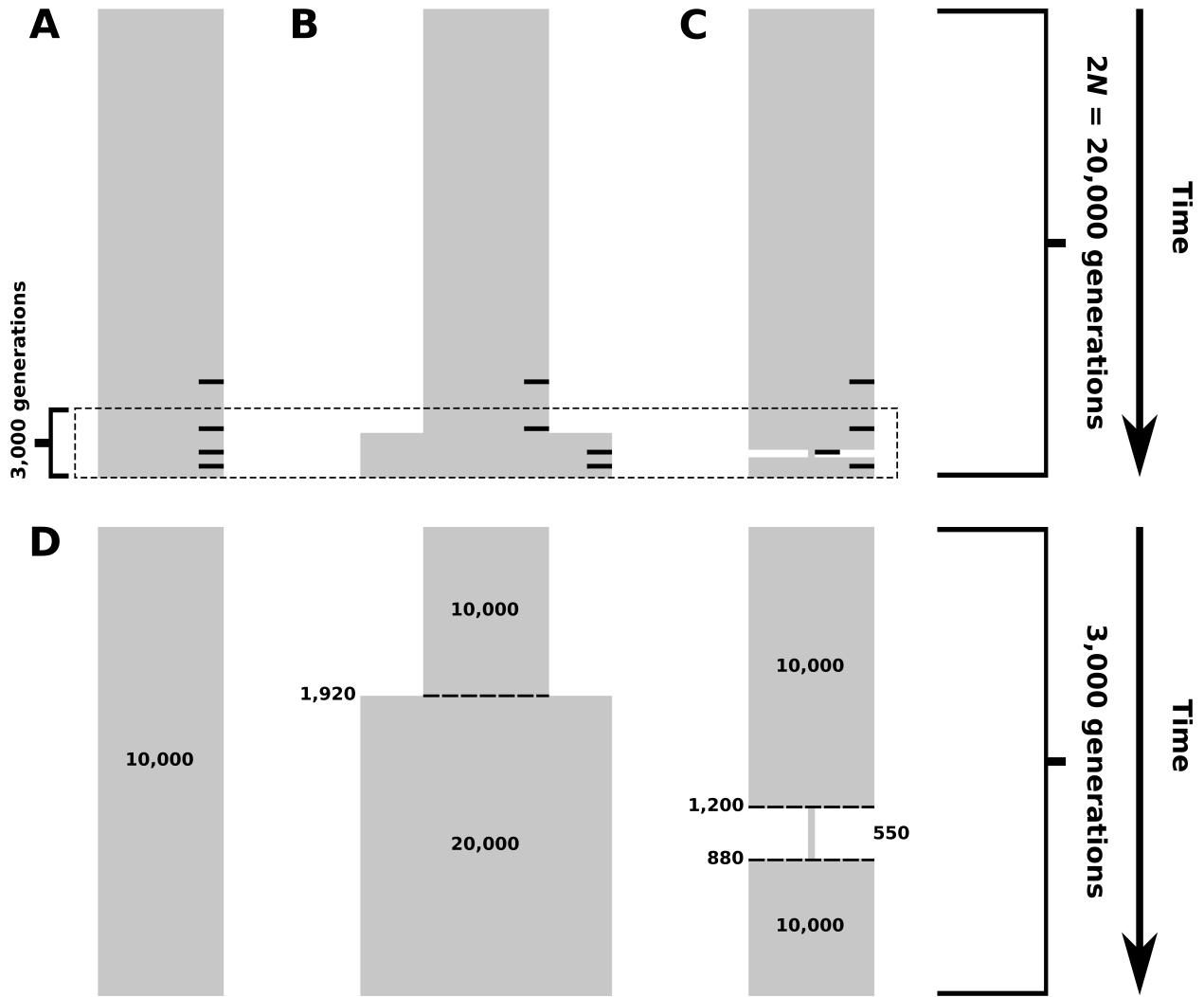


Figure 2: Simulated demographic models. Selection events, where applicable, occurred within $2N$ generations of sampling, indicated by small black bars on the right side of panels A-C corresponding to selection 4,000, 2,000, 1,000, and 400 generations before sampling. (A) Constant-size model. Diploid population size is 10^4 individuals throughout the time of simulation. (B) Model of recent population expansion. Diploid population size starts at 10^4 individuals and doubles to 2×10^4 individuals 1,920 generations ago. (C) Model of a recent strong population bottleneck. Diploid population size starts at 10^4 individuals and contracts to 550 individuals 1,200 generations ago, and subsequently expands 880 generations ago to 10^4 individuals. (D) View of the final 3,000 generations across demographic models, highlighting the effects of changing demographic factors on simulated populations.

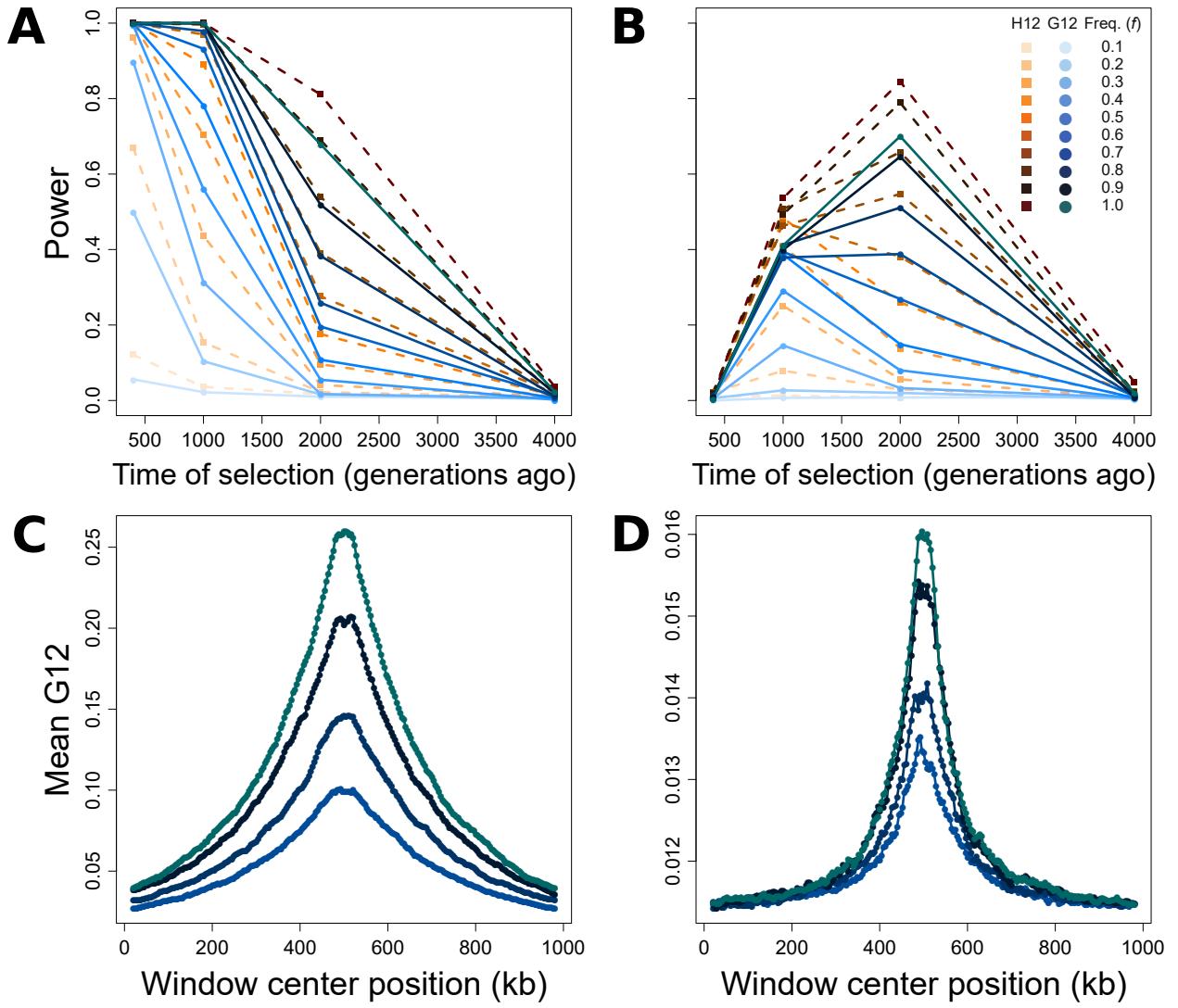


Figure 3: Capabilities of H12 (orange) and G12 (blue) to detect hard sweeps from simulated chromosomes, sample size $n = 100$ diploids, and window size of 40 kb for selection across four time points (400, 1,000, 2,000, and 4,000 generations before sampling) and 10 sweep frequencies (f , frequency to which the selected allele rises before becoming selectively neutral). Selection simulations conditioned on the beneficial allele not being lost. (A) Powers at a 1% false positive rate (FPR) of H12 and G12 to detect strong sweeps ($s = 0.1$) in a 100 kb chromosome. (B) Powers at a 1% FPR of H12 and G12 to detect moderate sweeps ($s = 0.01$) in a 100 kb chromosome. (C) Spatial G12 signal across a one Mb chromosome for strong sweeps occurring 400 generations prior to sampling. (D) Spatial G12 signal across a one Mb chromosome for moderate sweeps occurring 2,000 generations prior to sampling. Lines in (C) and (D) are mean values generated from the same set of simulations as panels A and B, and contain only results for $f \geq 0.7$. Note that vertical axes in panels C and D differ.

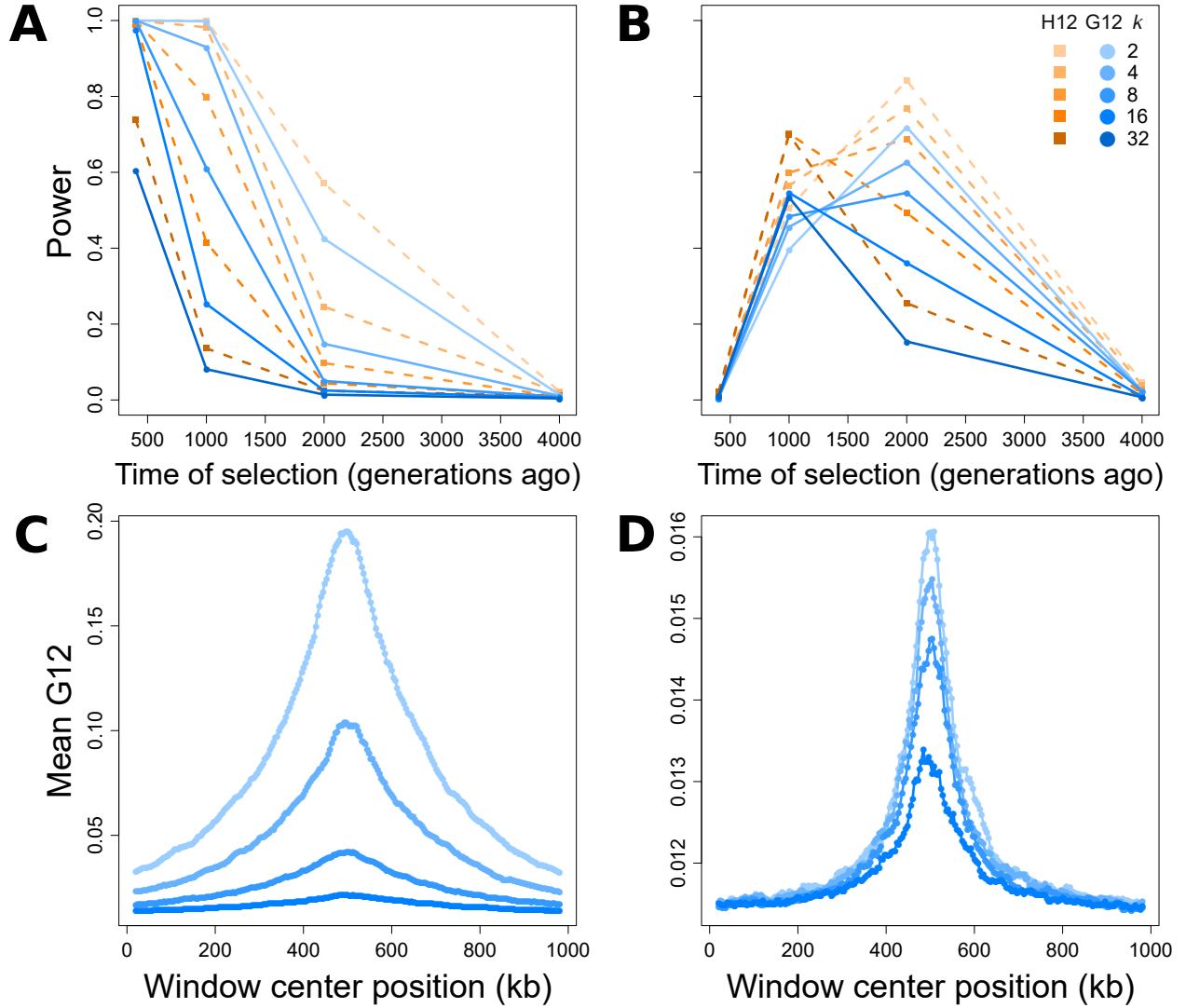


Figure 4: Capabilities of H12 (orange) and G12 (blue) to detect soft sweeps (SSV) from simulated chromosomes generated for selection times, sample size, and window size as in Figure 3, and five initially-selected haplotype values (k , number of haplotypes on which the selected allele arises at time of selection). Selection simulations conditioned on the beneficial allele not being lost. (A) Powers at a 1% false positive rate (FPR) of H12 and G12 to detect strong sweeps ($s = 0.1$) **in a 100 kb chromosome**. (B) Powers at a 1% FPR of H12 and G12 to detect moderate sweeps ($s = 0.01$) **in a 100 kb chromosome**. (C) Spatial G12 signal **across a one Mb chromosome** for strong sweeps occurring 400 generations prior to sampling. (D) Spatial G12 signal **across a one Mb chromosome** for moderate sweeps occurring 2,000 generations prior to sampling. Lines in (C) and (D) are mean values generated from the same set of simulations as panels A and B, and contain only results for $k \leq 16$. Note that vertical axes in panels C and D differ.

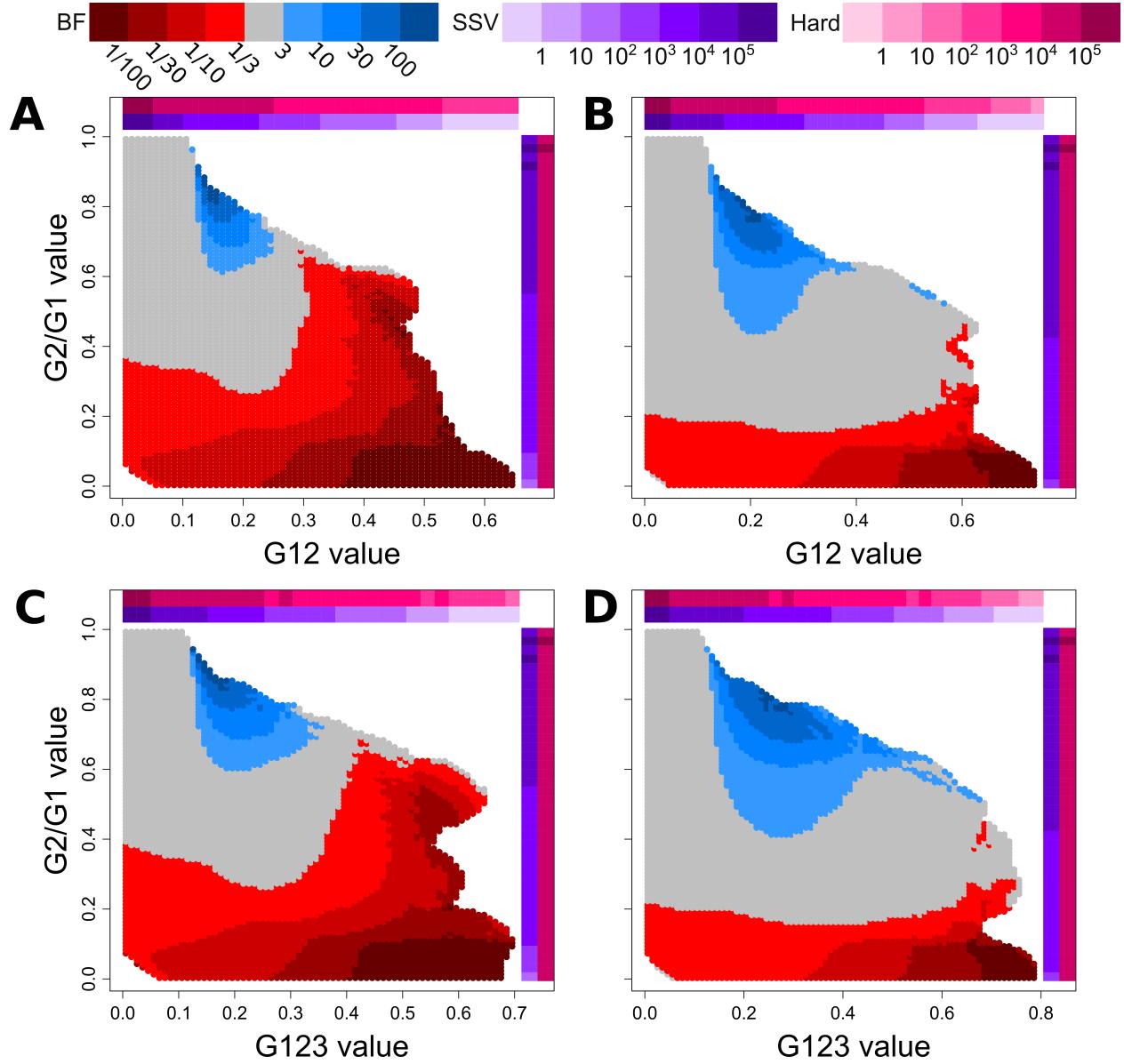


Figure 5: Assignment of Bayes factors (BFs) to tested paired values of (G12, G2/G1) and (G123, G2/G1). Plots represent the relative probability of obtaining a paired (G12, G2/G1) or (G123, G2/G1) value within a Euclidean distance of 0.1 from a test point for **hard** versus **soft sweeps**, determined as described in the *Materials and methods*. Selection coefficients (s) and times (t) were drawn as described in the *Materials and methods*. Red regions represent a higher likelihood for hard sweeps, while blue regions represent a higher likelihood for soft sweeps. Colored bars along the axes indicate the density of G12 or G123 (horizontal) and G2/G1 (vertical) observations within consecutive intervals of size 0.025 for hard sweep (magenta) and SSV (purple) simulations. (A) BFs of paired (G12, G2/G1) values for hard sweep scenarios and SSV scenarios ($k = 5$). (B) BFs of paired (G12, G2/G1) values for hard sweep scenarios and SSV scenarios ($k = 3$). (C) BFs of paired (G123, G2/G1) values for hard sweep scenarios and SSV scenarios ($k = 5$). (D) BFs of paired (G123, G2/G1) values for hard sweep scenarios and SSV scenarios ($k = 3$). Only test points for which at least one simulation of each type was within a Euclidean distance of 0.1 were counted (and therefore colored).

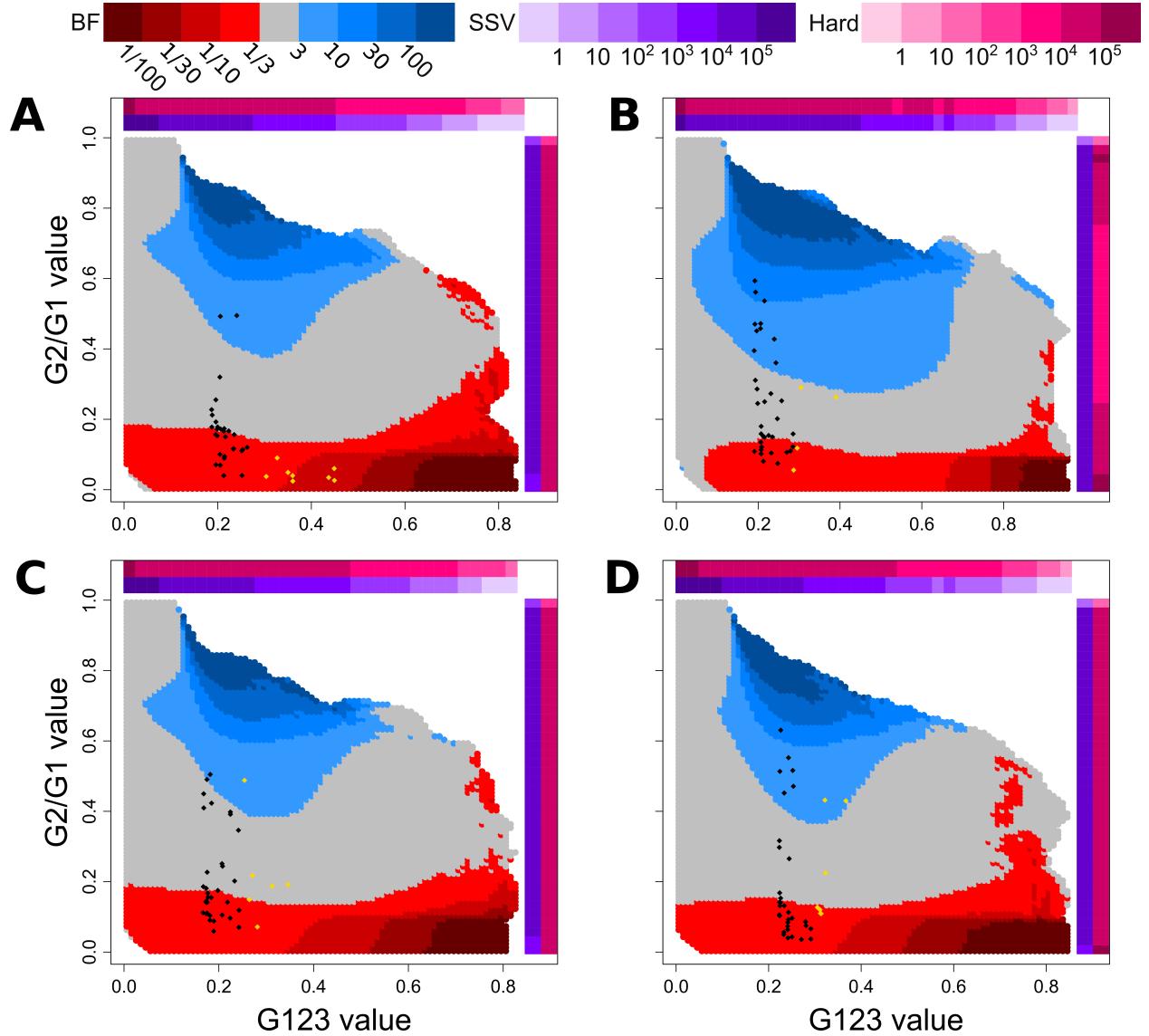


Figure 6: (G123, G2/G1) values used to distinguish hard (red) and soft (blue) sweeps in human empirical data using demographic models inferred with `smc++` [Terhorst et al., 2017]. Points representing the top 40 G123 selection candidates (Tables S4, S7, S10, and S13) for the (A) CEU, (B) YRI, (C) GIH, and (D) CHB populations are overlayed onto each population's specific (G123, G2/G1) distribution. Candidates exceeding the significance threshold (Table S1; different for each population) are colored in gold. Colored bars along the horizontal (G123) and vertical (G2/G1) axes are defined as in Figure 5.

G2/G1

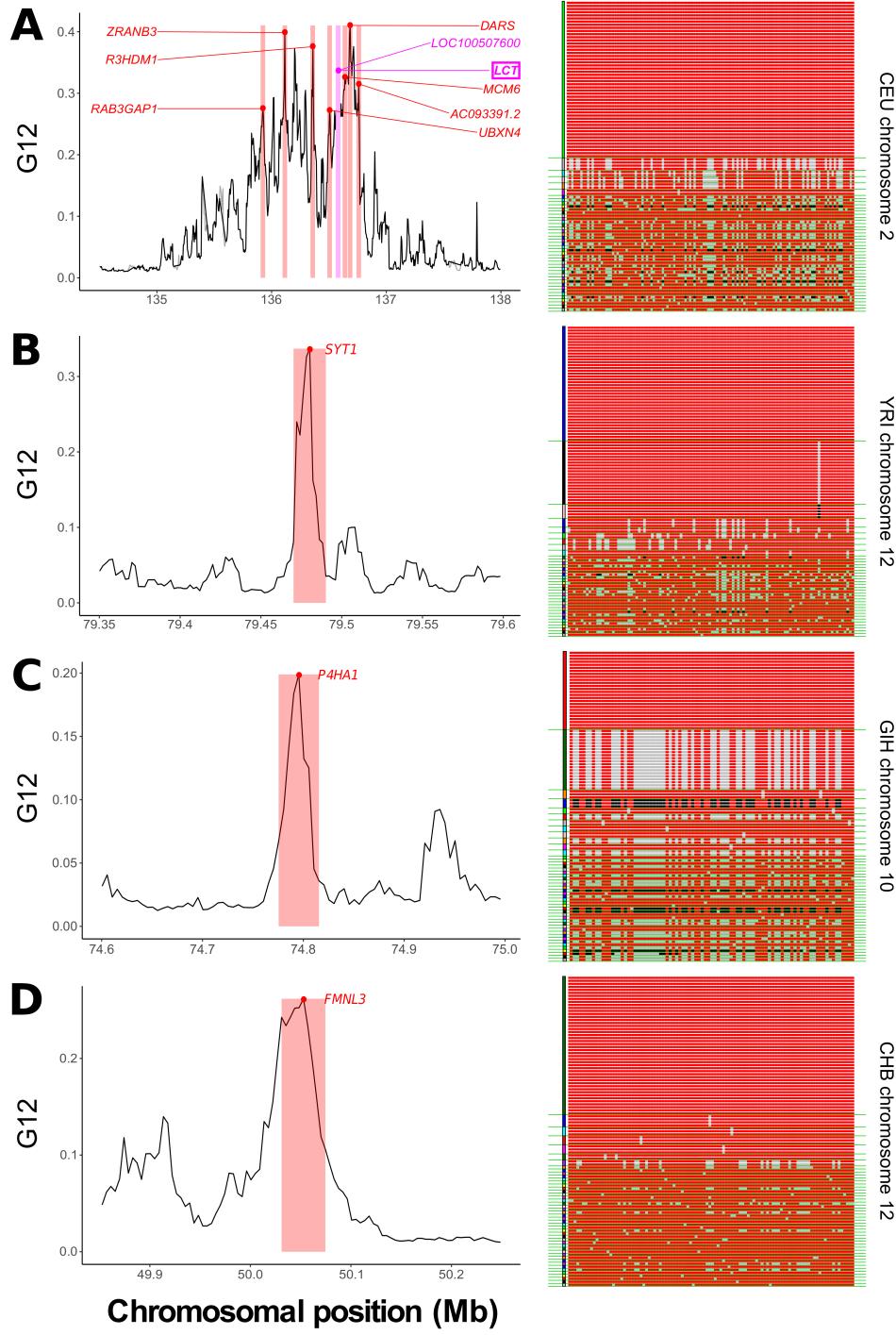


Figure 7: Outlying G12 signals in human genomic data. For each population, we show a top selection candidate and display its sampled MLGs within the genomic window of maximum signal. Red and black sites are homozygous genotypes at a SNP within the MLG, while gray are heterozygous. **Green lines separate MLG classes in the sample.** (A) CEU chromosome 2, centered around *LCT*, including other outlying loci (labeled). *LOC100507600* is nested within *LCT* (left). A single MLG exists at high frequency, consistent with a hard sweep (right). (B) YRI chromosome 12, centered on *SYT1* (left). This signal is associated with two elevated-frequency MLGs (right). (C) GIH chromosome 10, centered on *P4HA1* (left). Two MLGs exist at high frequency (right). (D) CHB chromosome 12, centered on *FMNL3* (left). A single MLG predominates in the sample (right).