




iHDSel software: The price equation and the population stability index to detect genomic patterns compatible with selective sweeps. An example with SARS-CoV-2

Antonio Carvajal-Rodríguez ^{1,*}

¹Centro de Investigación Mariña (CIM), Departamento de Bioquímica, Genética e Inmunología, Universidade de Vigo, Vigo, 36310 Spain

*Corresponding author. Centro de Investigación Mariña (CIM), Departamento de Bioquímica, Genética e Inmunología, Universidade de Vigo, 36310 Vigo, Spain. E-mail: acraaj@uvigo.es

Abstract

A large number of methods have been developed and continue to evolve for detecting the signatures of selective sweeps in genomes. Significant advances have been made, including the combination of different statistical strategies and the incorporation of artificial intelligence (machine learning) methods. Despite these advances, several common problems persist, such as the unknown null distribution of the statistics used, necessitating simulations and resampling to assign significance to the statistics. Additionally, it is not always clear how deviations from the specific assumptions of each method might affect the results. In this work, allelic classes of haplotypes are used along with the informational interpretation of the Price equation to design a statistic with a known distribution that can detect genomic patterns caused by selective sweeps. The statistic consists of Jeffreys divergence, also known as the population stability index, applied to the distribution of allelic classes of haplotypes in two samples. Results with simulated data show optimal performance of the statistic in detecting divergent selection. Analysis of real severe acute respiratory syndrome coronavirus 2 genome data also shows that some of the sites playing key roles in the virus's fitness and immune escape capability are detected by the method. The new statistic, called J_{HAC} , is incorporated into the iHDSel (informed HacDivSel) software available at <https://acraaj.webs.uvigo.es/iHDSel.html>.

Keywords: Price equation; information theory; population stability index; Jeffreys divergence; haplotype allelic class; selective sweep

Introduction

Evolutionary biology studies the factors that affect genetic variability in populations and species. The main processes that influence the evolution of this variability include mutation and recombination, genetic drift, migration, and natural selection. Natural selection, in addition to affecting the allele carrying a beneficial mutation, impacts the neutral alleles of loci linked to the selective one, producing what is known as genetic hitchhiking [1, 2], which leads to a selective sweep [3, 4], meaning a loss of diversity around the selected site. These sweeps can be complete or incomplete, strong or soft, and they can even overlap [5]. Regarding the detection of the footprint left by selective sweeps in genomes, from the earliest methods that explored haplotype patterns, whether by studying homozygosity [6], its diversity [7], or interpopulation differentiation [8], among others, a great number of methods have been developed and continue to be developed. Significant advancements have been made, including the use of summary statistics, the combination of different statistical strategies, and the incorporation of artificial intelligence-based methods [4, 9–14].

Most methods for detecting selective sweeps require the existence of haplotypic data, although, as discussed in [15], summary

statistics calculated from unphased genotypes are used, which, in a supervised machine learning context, allows the classification of genomic windows subject to selection. However, supervised methods are computationally expensive and are highly dependent on training data, and their performance with data from other species, genome types, and in general, outside the scenarios for which they have been trained is unclear [16].

Despite improvements in the efficiency and accuracy of methods for estimating haplotypes [17–19], in non-model species (understood as those in which, whether or not a genome has been sequenced, it is poorly annotated and has not traditionally been a model species in the pre-genomic era), haplotype-based detection methods are still not widely used. Instead, it is more common to use interpopulation methods based on detecting molecular markers with excessively high differentiation values, known as “outliers.” But even in the case of model species, the use of haplotype-based methods to detect selective sweeps presents the problem that the same genomic pattern that could be produced by a selective sweep could also be explained under different scenarios related to factors as diverse as the quality and characteristics of the sampled data, biological characteristics related to mutation and recombination rates, as well as

Received: 1 September 2024. Revised: 19 November 2024. Editorial decision: 21 November 2024. Accepted: 25 November 2024

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

This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial License (<https://creativecommons.org/licenses/by-nc/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact journals.permissions@oup.com

demographic history and the effects of purifying and background selection [20–22].

Part of this problem arises from the lack of knowledge of the null distribution of the statistics used, which requires simulating the neutral biological scenario. But overall, it is clear that although a statistical tool can detect a specific genomic pattern in the data, it is unlikely that that pattern could be due solely to the effect of a selective scan. It may do so in some scenarios, but not in others. Therefore, to validate a candidate single nucleotide polymorphism (SNP) or region as a result of a selective process, it is first necessary to prove that the statistic does not generate false positives in realistic scenarios in terms of demography and other evolutionary parameters of interest. Subsequently, functional validation of these candidate loci will always be necessary [20]. This does not preclude that the development of statistical tools to detect genomic patterns that may be related to selective sweeps remains of great interest. It would also be interesting if that statistic had a known null distribution.

When studying a selective sweep, we can trace its effect over time (directional selection) or across space (divergent selection). Therefore, if we use two samples to compare the effect of the sweep, they can be separated by time or space. Detecting the footprint of natural selection in genomes, in general, and specifically divergent selection, is important for studying speciation processes [23] and climate adaptation [24], but also for more immediate effects, such as resistance to infections in commercially important marine species [25, 26].

In this work, I propose a statistic that uses the population stability index, also known as Jeffreys divergence, to compare the distribution of allelic classes of haplotypes [27, 28] between two populations or samples. To develop the statistic, I use the informational interpretation of the Price equation [29, 30] defined for the haplotype allelic class (HAC) trait. The advantage of this statistic is that it follows a chi-square distribution when the null hypothesis (equal distribution of HACs among samples) is true. This not only increases computational efficiency by several orders of magnitude but also allows for the testing of biological models expected to deviate from this hypothesis, including the presence of local selection and its corresponding selective sweep. Below, I will present the development of the statistic and then demonstrate its behavior with both simulated and real genomic data from various samples of the severe acute respiratory syndrome (SARS) coronavirus 2 (SARS-CoV-2) virus.

The Price equation and the population stability index for comparing population genomes

Price equation

The Price equation in its most general formulation describes the change between two populations at any scale, spatial or temporal [30, 31]. The equation partitions the change into a part due to natural selection and another part due to other effects. We compare two populations or frequency distributions which can be separated by space and/or time. Natural selection causes populations to accumulate information, which is measured in relation to the logarithm of biological fitness $m = \log(\omega)$, where ω is the relative fitness [30, 32].

Therefore, let z be a character that takes different values z_i with associated frequency p_i in population P and with frequency q_i in population Q . If we consider the logarithm of fitness as the character, $z = m$, we have that the mean change in m due to the effect of natural selection in one or the other population is [30]

$$\Delta_s \bar{m} = J(p, q) = \beta_{mw} \frac{V_w}{\bar{w}} \quad (1)$$

where J is the Jeffreys divergence or population stability index, p and q the frequency of the different values of m in the populations P and Q , respectively, and β_{mw} is the regression of m on the absolute fitness w .

However, it is possible to use scales other than the fitness logarithm to measure information, with the key element being the regression of values in the new scale on fitness [33]. Therefore, to detect the effect of natural selection from genomic data, it will be necessary to measure those genomic patterns with high regression values on biological fitness. In this work, I propose using HACs as a suitable pattern to capture the increase in information generated by natural selection, whether in temporal comparisons (directional selection) or spatial comparisons (divergent selection).

HAC

HACs were initially introduced in Labuda *et al.* [27] and later used to detect genomic patterns caused by selective sweeps [28] and divergent selection [34].

Consider a sample of sequences and compute the reference haplotype R as the one formed by the major allele of each site. Now, consider for the same or another sample of sequences, the haplotypes of length $L + 1$ centered in a given candidate SNP c and define the mutational distance between any haplotype and the reference R as the Hamming distance between the haplotype and the reference, that is, the number h of sites in the haplotype carrying an allele different to the one in R (i.e. a minor allele if the system is biallelic). Each group of haplotypes having the same h will constitute an HAC [27, 28]. The HAC distribution is estimated from the distribution of the h values in a sample.

Thus, in a given haplotype with the candidate SNP position c in the middle, for each position other than c we count the outcome $X_k = I(s_k \neq r_k)$ where s_k is the allele in the position k of the haplotype, r_k is the allele in the reference, and $I(A)$ is the indicator variable taking 1 if A is true and 0 otherwise. Therefore, the h value of a haplotype of length $L + 1$ is

$$h = \sum_{k=1}^{L+1} X_k \text{ where } k \neq c, X_k = I(s_k \neq r_k) \text{ and } h \in [0, L] \quad (2)$$

The idea behind using h -values to detect selective sweeps is that if one allele increases in frequency due to the effect of selection, the higher frequency alleles from adjacent sites will be swept along with the selected allele so that these haplotypes will have many common alleles with the reference configuration, that is an h -value close to zero.

Information for HACs: the population stability index

Let h_i be the HAC value that satisfies $h = i$ with $i \in [0, L]$ then for a sample of n_1 sequences in P , the frequency of h_i is

$$P_i = \#h_i/n_1 \text{ with } \sum P_i = 1$$

where $\#h_i$ is the number of occurrences of h_i .

Similarly, for a sample of n_2 sequences in Q , the frequency of h_i is

$$Q_i = \#h_i/n_2 \text{ with } \sum Q_i = 1$$

In previous works, studying the distribution of alleles around a candidate site in both samples P and Q has been performed by comparing, in several ways, the HAC variances of the partitions that have the reference allele or not in the different samples [34, 35]. There are some problems with this type of approach as the unknown distribution of the defined statistics or a loss of power when using homogeneity variance tests. Here, I rely on the abstract model of the Price equation as proposed by Frank [30, 31, 33, 36] to calculate, using Jeffreys divergence, the change caused by selection in the distribution of HAC values between two populations.

Number of classes and smoothing

For a total of $L + 1$ different classes, the Jeffreys divergence is [37]

$$J_{HAC} = \frac{n_1 n_2}{n_1 + n_2} \sum_{i=0}^L (P_i - Q_i) \ln \frac{P_i}{Q_i}$$

However, computing J_{HAC} in this way could suffer from the curse of dimensionality [38] if eventually $L > n_1 + n_2$, which will cause the presence of the different classes to be scarce. To alleviate this problem, we will group the values in K ($K \leq L + 1$) HAC classes. The number of classes K is an important parameter because too many classes have the dimensionality issue but too few classes will have low power for the distribution comparison. A conservative heuristic guess is $K = (L + 1)/2$ when $L \geq 15$ or $K = L$ otherwise, since we have empirically verified that less than 15 classes implies a low detection power, possibly because a smaller number of classes implies very few SNPs, which may be due to very short sequences, and/or very low sample size, and/or very homogeneous samples.

Given K , we will group uniformly the h values into K groups so that the first group indicates classes with equal or less than $(100/K)\%$ of minor (non-major) alleles, the next corresponds to classes with more than $(100/K)\%$ but equal or less than $2 \times (100/K)\%$, until the last group with more than $(K - 1) \times (100/K)\%$ but equal or less than 100% . If necessary ($L + 1$ not divisible by K), the class with 100% of minor alleles is included in this last group. For simplicity, we consider K as a divisor of $L + 1$.

Thus, for population P , the frequency P'_i of each group of classes is

$$\begin{aligned} P'_i &= \sum_{j=u}^U \#h_j/n_1 \text{ where } i \in [0, K - 1], u = \frac{L+1}{K}i \text{ and } U \\ &= \frac{L+1}{K}(i+1) - 1 \end{aligned} \quad (3)$$

However, note that the Jeffreys divergence is defined only if P and Q have no zeros. To avoid zeros, we use additive smoothing [39] with a pseudocount $\alpha = 0.5$ for each possible outcome so that P'_i in Equation (3) becomes

$$P'_i = \sum_{j=u}^U (\#h_j + \alpha)/(n_1 + \alpha K)$$

So, for K groups of HAC classes, the Jeffreys divergence for comparing the HAC distribution between populations P and Q finally is [c.f. Equation (5.10) in [37] p. 130]

$$J_{HAC} = \frac{n_1 n_2}{n_1 + n_2} \sum_{i=0}^{K-1} (P'_i - Q'_i) \ln \frac{P'_i}{Q'_i} \quad (4)$$

with values in $[0, +\infty)$

Figure 1 shows an example of a J_{HAC} calculation for two samples with eight haplotypes each ($n_1 = n_2 = 8$). Haplotypes have a length of 9, so discounting the candidate site, we have eight sites and nine possible HAC classes (from 0 to 8). Note that if a class does not appear in either of the two samples, the contribution value to J_{HAC} for that class is 0. When the class is present only in one sample, to correct the problem of zeros, pseudocounting is applied with $\alpha = 0.5$, so that the frequency value of the class that does not appear will be $0.5/(8 + 9/2) = 0.5/12.5 = 1/25$ in that sample. If the count in a class is 2, the corrected frequency value will be $2.5/12.5 = 5/25$ (Fig. 1).

The advantage of using Equation (4) in the context of studying the genomic footprint of selection is that, contrary to other statistics, it can be approached by a chi-square distribution providing a faster approach as we can avoid performing computationally expensive simulations or resampling.

Phenotypic scale, linkage disequilibrium, and window size

Phenotypic scale

The gain in information caused by the effect of natural selection as expressed in Equation (1) depends on the log-fitness m and if we measure the frequency of the h_i classes instead of fitness classes, the relationship between the average change in the h distribution and the gain in information will depend on the regression of h -values on fitness as follows [33]

$$\Delta_s \bar{h} = \beta_{hw} \frac{V_w}{\bar{w}}$$

thus, if we use the HAC values to compute J we obtain J_{HAC}

$$J_{HAC} = \beta_{hw} D_w = \frac{\beta_{hw}}{\beta_{mw}} J$$

The quantity β_{hw}/β_{mw} is the change in phenotype (HAC values) relative to the change in information [33]. Therefore, if there is a perfect fit between $\ln(P/Q)$ and m , then $J_{HAC} = J$.

The regression of h on w will be high when it is fitness that is distributing the classes of h , which requires that there are indeed one or more sites under selection within the haplotype window. However, this is a necessary but not sufficient condition. Price's equation for total change indicates that the average variation in phenotype h has two components: one due to selection and the other due to other causes, including changes in the components of the phenotype that are transmitted (Δh)

$$\Delta \bar{h} = \Delta_s \bar{h} + q' \Delta h$$

In our context, the change in h not caused by selection may be due to, besides mutation, the effect of recombination on haplotypes, which in turn will depend on the window size. Therefore, we are interested in using window sizes that correspond to haplotype blocks in order to minimize Δh .

Window size

The program computes haplotype blocks and sets the candidate position c in the middle of each block. A haplotype block is computed as a sequence of reference SNPs with length W that satisfies $r^2(c - W/2, c - W/2 + 1), \dots, r^2(c - 1, c), r^2(c, c + 1), r^2(c + x - 1, c + x), \dots, r^2(c + W/2 - 1, c + W/2), \dots$ where r is the correlation coefficient calculated from the sample of size n , so that $Pr(nr^2) \leq \alpha$,

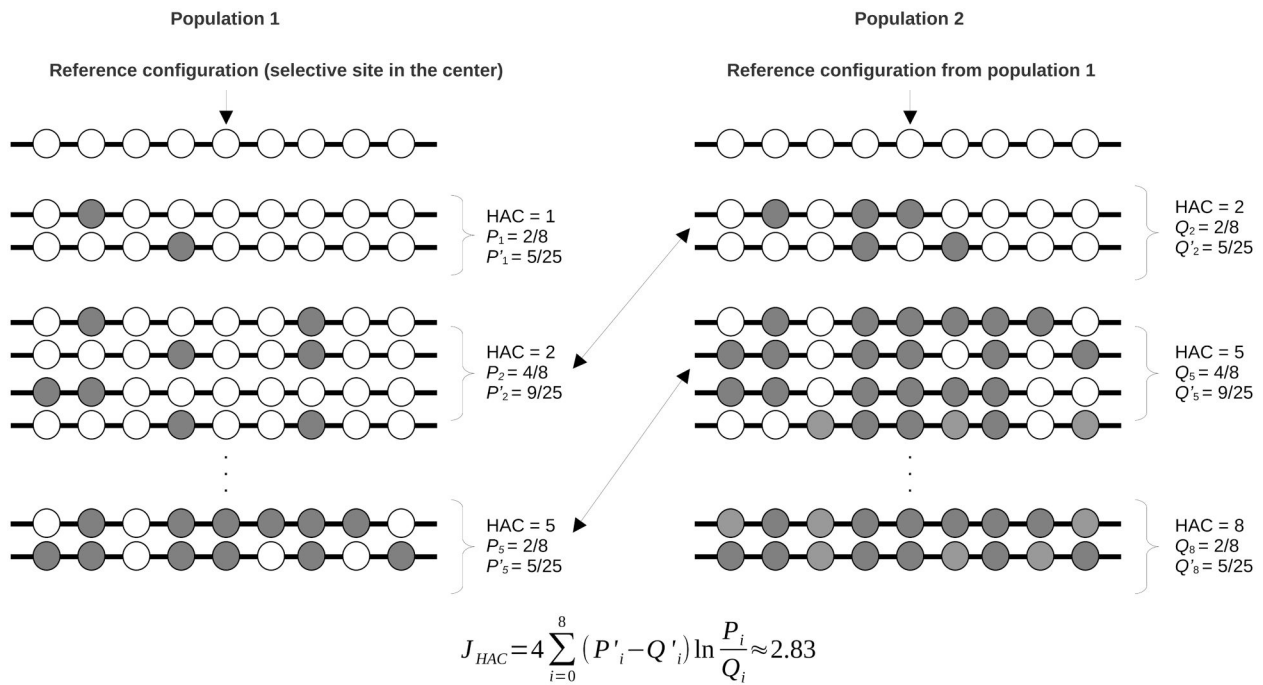


Figure 1. Reference configuration and the distribution of HAC values in two samples and calculation of the J_{HAC} statistic. The haplotype length is 9 and the number of haplotypes is 8 in each sample. White circles represent major alleles and gray circles represent any other allele (usually the minor). Note that J_{HAC} is also known as the population stability index and is asymptotically distributed as chi-square with $K-1$ degrees of freedom.

and nr^2 has a chi-square distribution. Furthermore, for a given SNP $c+1$ to be included in the block, it is also required that $D'(c, c+1) \geq 0.4$, where D' is the normalized linkage disequilibrium [40]. The block is extended until any of both conditions is rejected that is, $\Pr(nr_{c+x-1, c+x}^2) > \alpha$ or $D'(c+x-1, c+x) < 0.4$.

Optionally, the program can use an outlier as the putative center of a block and build the block around it. In this case, the condition for defining a block is more liberal, allowing blocks that have a mean normalized linkage disequilibrium value greater than zero. The reason is that the outliers may have been part of older blocks, so we use the minimum condition that the average linkage of reference alleles is greater than zero assuming that, if they are not the product of selective sweep, the distribution of HACs will not be affected, the latter will be checked in the next section by simulation.

The windows used by the program do not overlap when calculating automatic blocks, but may overlap if the outlier option is used or if the user passes specific candidate sites to the program.

Simulations

To check the performance of the method, its power, and its control of false positives, we will perform two types of simulations, namely, with diploid and haploid genomes. Simulations were carried out with the GenomePop2 program [41].

Diploid genomes

For diploid genomes, the same simulated data as in [34, 35] were used. Two populations of $N=1000$ facultative hermaphrodites were simulated under divergent selection and different conditions about mutation, recombination, migration, and selection. Each individual consisted of a diploid chromosome of length 1 Mb. In Tables 1–4, we can appreciate the different cases with the corresponding parameter values. The population migration rate was $Nm = 10$ in all cases. The number of generations was 10^4 or 5×10^3 , the population mutation rate $\theta = 4N\mu$ was {12, 60}

where μ is the mutation rate per haploid genome, which implies mutation rates per site and per generation of $\{3 \times 10^{-9}, 1.5 \times 10^{-8}\}$, the population recombination rate $\rho = 4Nr$ was {0, 4, 12, 60} where r is the recombination rate per haploid genome, which implies recombination rates per pair of sites per generation of $\{0, 10^{-9}, 3 \times 10^{-9}, 1.5 \times 10^{-8}\}$ or noted as infinite when the segregation was independent. The selection coefficient s was ± 0.15 depending on if the mutant allele is deleterious (+0.15) or beneficial (−0.15). This implies that $4Ns = 600$, which can be considered a moderate selection associated with a locus with large effect as expected for the model that resembles the most favorable conditions for the formation of ecotypes under local adaptation, which is the context for which these data were generated [34].

iHDSel input settings for analyzing diploid simulation data

A minor allele frequency (MAF) value of 0.01 was used. As we have already seen, the program allows defining the window or haplotypic block size automatically, using the correlation between pairs of sites to define the block size and placing the central SNP as a candidate or, alternatively, it uses the detected outliers as candidate SNPs and then calculates the window size. Both methods were used. All other parameters were as defined by default (maximum window size 1000, minimum window size 11, significance level 0.05, etc., see the program manual). An example of the command line to launch case C1 (Table 1) and analyze the 1000 files located in subfolder C1 and using the automatic calculation of blocks (-useblocks 1) is:

```
./iHDSel0.5.2 -path/home/data/C1/-runs 1000 -input Om_SNPFile_Run
-format ms -sample 50 -minwin 11 -output JHAC_C1_ -maf0.01
-useblocks 1 -doEOS 1 &
```


Table 1. Percent power for detecting divergent selection by J_{Hac} in simulated data with the selective site in the middle.

Case	T	θ	ρ	s	% power	Dist Kb	W
C1	10^4	12	0	± 0.15	100 (98)	–	14 (13)
C2	10^4	12	4	± 0.15	100 (98)	42 (38)	14 (13)
C3	10^4	12	12	± 0.15	100 (96)	4 (10)	13 (12)
C7	5×10^3	60	0	± 0.15	100 (94)	–	13 (12)
C8	5×10^3	60	4	± 0.15	100 (85)	37 (14)	13 (12)
C9	5×10^3	60	60	± 0.15	98 (80)	14 (2)	12 (11)
C13	10^4	60	0	± 0.15	100 (100)	–	14 (=)
C14	10^4	60	4	± 0.15	99 (100)	126 (15)	14 (13)
C15	10^4	60	60	± 0.15	95 (91)	19 (2)	13 (12)
C15Indep	10^4	60	∞	± 0.15	0 (2 ^a)	– (-)	– (11)

The power was computed as 100× the number of replicates where selection was detected/1000. In parentheses, the corresponding value when the blocks were built around outliers instead of finding the blocks automatically, if the value is equal the = symbol appears. Genome size is 1 Mb. Population size $N = 1000$. T: number of generations. Population mutation rate $\theta = 4N\mu$. Population recombination rate $\rho = 4Nr$. s: selection coefficient. Dist: average distance in Kb from the detected position to the actual effect, given only when $\rho > 0$. W: average size, in number of SNPs, of the haplotypes analyzed. Significance level $\alpha = 0.05$. Each case was replicated 1,000 times.

^a Note that this 2% results from using the outlier-centered haplotype method. When directly inspecting outliers with the EOS method, the power was 78%.

Table 2. Percent power for detecting divergent selection by J_{Hac} in simulated data with the selective site in different locations

Case	T	θ	ρ	s	Loc	% power	Dist Kb	W
C13loc0	10^4	60	0	± 0.15	0.0	100 (99)	–	14 (13)
C13loc10	10^4	60	0	± 0.15	0.01	100 (99)	–	14 (13)
C13loc100	10^4	60	0	± 0.15	0.1	100 (98)	–	14 (13)
C13loc250	10^4	60	0	± 0.15	0.25	100 (99)	–	14 (13)
C14loc0	10^4	60	4	± 0.15	0.0	98 (93)	300 (262)	14 (13)
C14loc10	10^4	60	4	± 0.15	0.01	98 (96)	285 (292)	14 (13)
C14loc100	10^4	60	4	± 0.15	0.1	99 (96)	180 (229)	14 (13)
C14loc250	10^4	60	4	± 0.15	0.25	99 (98)	62 (114 ^a)	14 (14)
C15loc0	10^4	60	60	± 0.15	0.0	86 (79)	211 (189)	13 (11)
C15loc10	10^4	60	60	± 0.15	0.01	87 (80)	198 (170)	13 (11)
C15loc100	10^4	60	60	± 0.15	0.1	91 (89)	106 (70)	13 (12)
C15loc250	10^4	60	60	± 0.15	0.25	92 (89)	37 (14)	13 (12)

The power was computed as 100× the number of replicates where selection was detected/1000. In parentheses, the corresponding value when the blocks were built around outliers instead of finding the blocks automatically, if the value is equal the = symbol appears. Genome size is 1 Mb. Population size $N = 1000$. T: number of generations. Population mutation rate $\theta = 4N\mu$. Population recombination rate $\rho = 4Nr$. s: selection coefficient. Loc: true relative position of the selective site. Dist: average distance in Kb from the detected position to the actual effect, given only when $\rho > 0$. W: average size, in number of SNPs, of the haplotypes analyzed. Significance level $\alpha = 0.05$. Each case was replicated 1000 times.

^a Several runs with average $F_{ST} > 0.5$ and no outliers, so the 90th percentile was considered.

Table 3. Percent power for detecting divergent selection by J_{Hac} in simulated data for a polygenic model with five selective sites uniformly distributed in the chromosome.

Case	Candidate	% power	W
C15poly	1	99 (97)	15 (18)
C15poly	2	89 (0)	16
C15poly	3	75 (0)	16
C15poly	4	59 (0)	16
C15poly	5	44 (0)	16

The power was computed as the number of replicates where the selection was detected. In parentheses the corresponding % power when the blocks were built around outliers instead of finding the blocks automatically, if the value is equal the = symbol appears. Genome size is 1 Mb. Population size $N = 1000$. Number of generations $T = 10^4$. Population mutation rate $\theta = 4N\mu = 60$. Population recombination rate $\rho = 4Nr = 60$. Selection coefficient per site $s = \pm 0.032$. W: average size, in number of SNPs, of the haplotypes analyzed. Each case was replicated 100 times.

Table 4. Percent false-positive rate for detecting divergent selection in simulated neutral data.

Case	T	θ	ρ	% FPR	W
C4	10^4	12	0	0.1 (1)	11 (=)
C5	10^4	12	4	0.3 (2)	12 (11)
C6	10^4	12	12	0.1 (4)	12 (11)
C10	5×10^3	60	0	0 (0.4)	– (11)
C11	5×10^3	60	4	0 (2)	– (11)
C12	5×10^3	60	60	0.2 (3)	12 (11)
C16	10^4	60	0	0.3 (0.4)	12 (11)
C17	10^4	60	4	1 (2)	12 (11)
C18	10^4	60	60	1 (4)	13 (=)
C18Indep	10^4	60	∞	0 (0.2)	– (11)
C18Bottle	10^4	60	60	3 (13)	12 (11)
C18Bottle	10^4	60	60	2	26 ^a
C18Bottle	10^4	60	60	2	51 ^a

In parentheses the corresponding value when the blocks were built around outliers instead of finding the blocks automatically, if the value is equal the = symbol appears. Genome size is 1 Mb. Population size $N = 1000$. T: number of generations. Population mutation rate $\theta = 4N\mu$. Population recombination rate $\rho = 4Nr$. %FPR = 100× number of replicates with significant J_{Hac} test/1000. W: average size, in number of SNPs, of the haplotypes analyzed. Each case was replicated 1000 times.

^a Window size set to a specific value.

The $-doEOS$ tag indicates whether we want (1, default) or not (0) to run in addition to the EOS outlier test [34]. If the calculation without blocks is used ($-useblocks$ 0) the $doEOS$ tag must necessarily be set to 1.

Diploid simulation results

In the following tables, the results of power (Tables 1–3) and false-positive rate (Table 4) after analyzing 1000 replicates of each scenario are presented. In summary, for haplotypes with linkage and the selective site in the center of the chromosome, when using the automatic blocks system, the power is equal to or greater than 95%, regardless of mutation and recombination rates. As expected, if the sites are not linked, the method does not work because there is no selective sweep (Table 1). When the position of the selective site moves away from the center of the chromosome (Table 2), the power remains high. Localization improves as recombination increases and as the marker is located closer to the center. In the case of multiple selective sites (Table 3), the power to detect at least three is above 75% when using automatic blocks but only detects one (97% power) in the case of blocks centered on outliers. In general, for blocks centered on outliers, the power is slightly lower, but in some cases, the localization was considerably more accurate.

Finally, in the neutral simulations where there was no selective site (Table 4), the false-positive rate conservatively remains below the expected 5%, both using automatic blocks and those centered on outliers, with one exception corresponding to the effect of bottlenecks. When a bottleneck occurs, it can generate linkage disequilibrium that could resemble the effect of a selective sweep, thus increasing the possibility of false positives [42–44]. In our case, we observed that J_{Hac} becomes liberal with 13% when the blocks are centered around the outliers, which means an 8% excess over the expectation. The explanation for this happening with blocks centered on outliers but not with automatic ones is that, as previously indicated, the construction of blocks centered on outliers is somewhat more liberal, validating as blocks those regions that have an average disequilibrium greater than 0. A conservative option available for the above exception is

to set the window size to a higher value, say 25 or 50, which solves the problem and sets the false-positive rate to just 2%. While for the corresponding selective case when we run the program with these window sizes the power is 90%. The underlying logic is that if the positive is due to an increase in the frequency of a pattern randomly generated by the bottleneck, then some increase in the window size, say doubling it, will undo the effect. However, if the positive is real, it is relatively easy for it to remain unless the sweep has already been detected at the limit of its size.

Haploid genomes

For haploid genomes, we consider a scenario related to the SARS-CoV-2 virus, which belongs to Severe Acute Respiratory Syndrome (SARS)-related viruses and is a positive-sense single-stranded RNA virus with a 30kb genome. The model is based on average parameter values associated with the inpatient evolutionary dynamics of the virus [45, 46].

The virus's generation interval for 1 year corresponds to 861.64 viral cycles/year [45, 46], and the lower estimate of the sampling time per individual after infection onset is approximately 7 days, equivalent to 17 generations.

The simulated scenario corresponds to a genome of 30 000 nucleotides, an effective population size of 10 000, a mutation rate per site of 2×10^{-6} , resulting in a mutation rate $\mu = 0.06$ per genome, and a recombination rate with the same value or absent. The population starts in a mutation-drift equilibrium, where the effective number of alleles, n_e has been calculated, defining a mutant proportion of $1/n_e$ [47].

Two selective sites are defined in the Spike region at nucleotide positions 23 403 and 23 604, with a favorable selection coefficient of 0.2 for the mutation. The three nucleotides corresponding to the triplets carrying these mutations (23 402–23 404 and 23 603–23 605) are set at an initial frequency of 5% at the beginning of the simulation, after reaching mutation-drift equilibrium.

The populations to be compared consist of two samples: one taken after 7 days of infection (17 generations) and another taken at the end of the simulation after an additional 7–8 days (generation 35). Each sample size is 1000. We conducted 1000 runs for each simulated case.

Two different scenarios are considered. The first scenario, without a bottleneck, includes two cases: one with selection on the two indicated sites and a neutral one. The second scenario simulates a bottleneck that occurs after 7 days and corresponds to a new infection, starting with only 1, 2, or 5 viruses. The bottleneck is modeled as discrete logistic growth [48] with a growth rate of 2 and an initial value corresponding to the founder effect (1, 2, or 5). For each of these situations, we simulate a neutral and a selective case. Simulations were carried out with a modified version of the GenomePop2 program [41].

Tables 5 and 6 present a summary of the parameter values used, along with the results of the detection power and false-positive rate of the J_{Hac} statistic.

iHDSel input settings for analyzing haploid simulation data

An MAF value of 0.01 was used. For the haploid scenario, the program was not able to automatically detect blocks, possibly due to short evolutionary time and therefore we used the outliers as candidate SNPs to calculate the window size. All other parameters were as defined by default (see the program manual). An example of the command line to launch the base case and analyze

Table 5. Percent power for detecting directional selection by J_{Hac} in simulated data for haploid genomes.

θ	ρ	s	Bottleneck	% power
1200	0	-0.2	-	66
1200	1200	-0.2	-	64
1200	0	-0.2	1	3
1200	1200	-0.2	1	4
1200	0	-0.2	2	38
1200	1200	-0.2	2	56
1200	0	-0.2	5	53
1200	1200	-0.2	5	56

The number of generations was 35. Population size $N = 30\,000$. Population mutation rate $\theta = 2N\mu$. Population recombination rate $\rho = 2Nr$. s : selection coefficient. Bottleneck: initial bottleneck at generation 17, a dash implies that there is no bottleneck. Each case was replicated 1000 times.

Table 6. Percent false-positive rate for detecting directional selection in simulated neutral data for haploid genomes.

θ	ρ	s	Bottleneck	%FPR
1200	0	0.0	-	0
1200	1200	0.0	-	0
1200	0	0.0	1	0
1200	1200	0.0	1	0
1200	0	0.0	2	3
1200	1200	0.0	2	6
1200	0	0.0	5	4
1200	1200	0.0	5	5

The number of generations was 35. Population size $N = 30\,000$. Population mutation rate $\theta = 2N\mu$. Population recombination rate $\rho = 2Nr$. s : selection coefficient. Bottleneck: initial bottleneck at generation 17, a dash implies that there is no bottleneck. %FPR = $100 \times$ number of replicates with significant J_{Hac} test/1000. Each case was replicated 1000 times.

the 1000 files located in subfolder Serial_M1_sites2 and using the outlier-based windows (-useblocks 0) is:

```
./iHDSel0.5.2 -path/home/data/Serial_M1_sites2/-runs 1000
-input GP2msout_Run -format ms -sample 1000 -minwin 11
-output simSARSCv2_ -maf 0.01 -useblocks 0.
```

Haploid simulation results

With the simulation of haploid genomes, the power is lower due to the shorter evolutionary time but it is still around 65%. The impact of bottlenecks in this short period is to reduce the power, which is logical because it produces the loss of the selection signal. However, in the presence of recombination and with bottlenecks of two or more individuals, the power remains at 56% while the false-positive rate remains between 4 and 6% (Tables 5 and 6).

Real data analysis: SARS-CoV-2

SARS-CoV-2 virus genomes stored in the GISAID database [49] are indexed by both locality and the time period in which they were sampled, thus presenting a unique opportunity to apply iHDSel to time or spatially separated samples. Therefore, as an example of application, we are going to compare SARS-CoV-2 genomes sampled in Spain (SP), England (EN), and South Africa (SA) in periods corresponding to different waves. The findings of this section are based on data associated with 30 274 SARS-CoV-2 genomes available on GISAID up to 12 February 2024, gisaid.org/EN1, gisaid.org/EN2, gisaid.org/EN3, gisaid.org/EN4, gisaid.org/SP1, gisaid.org/SP2, and gisaid.org/SA.

The downloaded genomes were complete (>29 000 bp) and of high quality (<1% undefined bases and <0.05% unique amino acid mutations). These datasets were then processed using the

Nextclade CLI for quality control [50]. Briefly, the Nextclade CLI examines the completeness, divergence, and ambiguity of bases in each genome. Only genomes considered “good” by Nextclade CLI were selected.

The samples from England (EN1, EN2, EN3, and EN4) correspond to the period of March 2020, at the beginning of the first wave of the pandemic (EN1, 4820 genomes collapsed to 4227 after quality control), a second sample taken between 28 March and 31 March 2021, inclusive (EN2, 5966 genomes collapsed to 4152 after quality control), a third from 24 June to 26 June 2021, inclusive (EN3, 6886 genomes collapsed to 5844 after quality control), and from 1 October 2023, until 31 January 2024, inclusive (EN4, 3928 genomes collapsed to 3712 after quality control).

The samples from Spain (SP1 and SP2) correspond to the periods 24 June 2021 to 12 July 2021, inclusive (SP1, 6195 genomes collapsed to 4627 after quality control) and 1 October 2023 to 31 January 2024, inclusive (SP2, 1012 genomes collapsed to 221 after quality control).

Finally, the sample from South Africa corresponds to the same period as SP1, 24 June 2021 to 12 July 2021, inclusive (SA, 1467 genomes collapsed to 1327 after quality control).

These samples will allow us to compare population changes in space or time. We will compare genomes from different samples to study if there are genomic patterns that the J_{HAC} test identifies as potentially caused by selection (see below).

Rationale of the comparisons

Spatial comparisons: SP1–SA, EN3–SA, EN3–SP1

These comparisons involve samples from different countries obtained in the same time period of the pandemic. The interest in the comparison with South Africa is that from 24 June 2021 to 12 July 2021, vaccination rates were high in Spain and England but very low in South Africa. Virtually, 100% of the Spanish and English population was vaccinated with at least one dose and less than 10% of the South African population [51].

Temporal comparisons: EN1–EN2, EN2–EN3, and EN3–EN4

These comparisons affect the same country but across different periods of the pandemic, from the beginning of the first wave to the beginning of 2024 with virtually the entire population already vaccinated several times and the majority variant being Omicron and its subvariants [52–54].

Spatial comparisons: EN4–SP2

At the end of 2023, the JN.1 subvariant of Omicron, originating from the BA.2.86 lineage, began to spread. This subvariant already carried more than 30 mutations in the spike protein compared to previous subvariants. JN.1 includes the L455S mutation and, by the

end of 2023, exhibited a higher reproductive rate than previous sublineages in countries, such as Spain, France, and England, with the number of detected JN.1 sequences being higher in England than in Spain [55]. During this period, DV.7.1, a sublineage of BA.2.75, was highly prevalent in Spain, 50% compared to 5% in the UK [56], and was considered a variant to monitor, although it was later downgraded. Therefore, the comparison between EN4 and SP2, corresponding to October 2023 to January 2024, is of interest to study the potential patterns of divergent selection in the evolutionary dynamics of Omicron subvariants between these two countries.

Genome alignment and lineage classification

The pooled genomes for each comparison were aligned with the MAFFT FFT-NS-2 program [57] with the specific version for SARS-CoV-2 accessible online [58].

Sequences that had more than 5% ambiguous sites were removed and also, to keep the alignment length the same as the input, insertions were deleted. The remaining options were the default. After the alignment, and following the protocol recommended by NextStrain given the possibility of artifactual SNPs located at the beginning and end of the alignment [59], sites in the first 130 base pairs and the last 50 were removed using the program Mega X [60]. Lineages were identified with Nextclade CLI (Table 7).

Input settings for iHDSel

An MAF value of 0.01 was used. The two methods already mentioned were used to define the window size (automatic or outlier-centered blocks) and the results detected by either of the two methods are reported. All other parameters were the ones by default (see program manual). An example of the command line for the comparison between EN3 and SP1, where both samples are in the file EN3_SP1.fas located in the data folder and using the outlier-centered block calculation (-useblocks 0) is:

```
-path/home/data/-input EN3_SP1.fas -format fasta
-output EN3_SP1 -useblocks 0 -tag ENGLAND &
```

where *-tag* is the argument that defines the word included in the name from the England sequences and that allows to separate both samples.

Similarly, for the temporal comparison between EN2 and EN3

```
-path/home/data/-input EN2_EN3.fas -runs 1 -format fasta
-output EN2_EN3 -useblocks 0 -tag 2021-03 -reference2
```

where we have added the *-reference* tag to indicate that the EN3 sample should be used as a sample to calculate the blocks and the reference haplotype.

Table 7. Percentage of SARS-CoV-2 lineages in the analyzed data.

Data	%Alpha	%Beta	%Delta (%AY.4/AY.45)	%Gamma	%Omicron (%JN.1/FLIP/DV.7.1)	%Other (pre-Alpha, Lambda, Mu, recombinants, undefined)
SP1	24	2	70 (2/0)	2	0	2
SA	1	3	94 (0/57)	0	0	2
EN1	0	0	0	0	0	100 (pre-Alpha)
EN2	98	1	0.1	0.1	0	0.8
EN3	1	0.02	98.9 (72/0)	0	0	0.08
EN4	0	0	0	0	96 (39/6/1)	4 (recombinants)
SP2	0	0	0	0	97 (26/12/28)	3 (recombinants)

The imprint of selection in the SARS-CoV-2 genomes

Spatial comparisons: SP1–SA (summer 2021)

The SP1 sample has a majority Delta (70%) and Alpha (24%) composition, while SA is mostly (94%) Delta (Table 7). The pooled SP1–SA sample consists of 247 SNPs with a frequency greater than 1%. After genome-wide analysis, iHDSel did not find any significant haplotypic blocks in the automatic search nor when focusing on outliers.

Spatial comparisons: EN3–SA (summer 2021)

Both samples are mostly Delta (99% EN3 and 94% SA, Table 7). The pooled EN3–SA sample consists of 107 SNPs with a frequency higher than 1%. After the whole-genome analysis, iHDSel found one site with the automatic block method (28282) and five sites centered on outliers (sites 7,851; 13,812; 21,846; 21,987 and 25 413 in Table 8).

The first site is 7851, which corresponds to ORF1a 2529. In the SA sample, 100% of the sequences have the amino acid A, while in EN3, there is 27%A and 73%V, indicating the change A2529V. It is noteworthy that A2529V is one of the main SARS-CoV-2 mutations associated with virus fitness [61]. Moreover, in a recent study [62], analyzing the evolution of different lineages in relation to the progress of vaccination, the A2529V mutation in ORF1a showed a significant positive correlation between the prevalence of the mutation and vaccination in Norway during the first 9 months of 2021 (including the sampling period of EN3 and SA).

The second site is 13 812, which, after identifying the slippery region [63] and the start of ORF1b at 13 468, corresponds to amino acid 115 in ORF1b (NSP12). This site has 100%M in EN3 but 42%M and 58%I in SA. The change M115I is a characteristic mutation of the AY.45 lineage [64], which is present in SA with a frequency of 57% but is absent in EN3.

The third and fourth sites are mutations corresponding to amino acid changes in the Spike protein. Specifically, T95I represents the change observed between SA and EN3, with I at a frequency of only 8% in SA but 72% in EN3. The other mutation in Spike is G142D, with D present at 62% in SA and 97% in EN3 (Table 8). Both mutations are characteristic of the Delta variants and increase in frequency in Delta Plus [65–68].

The fifth site is position 25 413 of the genome, corresponding to amino acid 7 in ORF3a, with amino acid I in both samples being EN3 (ATC) and SA (ATT|50%C). Therefore, the existence of a significant signal due to different HAC distributions must be caused by accumulated variation in the surrounding sites. Similarly, the sixth and final site corresponds to amino acid 3 of the N protein, with the amino acid being D (GAT) in 99% of the cases in both samples, with practically 1% being L (CTA). Again, the existence of a significant signal due to different HAC

distributions is caused by accumulated variation in the surrounding sites.

Spatial comparisons: EN3–SP1 (summer 2021)

We already saw that the EN3 genomes are predominantly Delta (99%), while SP1 has 70% Delta genomes and 24% Alpha (Table 7). The combined EN3–SP1 sample consists of 154 SNPs with a frequency greater than 1%. After the whole-genome analysis, iHDSel found one significant site. The nucleotide site 7851 corresponds to amino acid 2529 in ORF1a, which was also significant in the EN3–SA comparison, and we saw that A2529V is one of the main SARS-CoV-2 mutations associated with virus fitness. In this comparison, the change is from 98%A in SP1 to 73%V (27%A) in EN3.

Therefore, regarding the spatial comparisons in the summer of 2021, we see that in the SA and SP1 samples, amino acid 2529 of ORF1a was still A in virtually 100% of the sequences analyzed, while in EN3, only 27% had A and the remaining 73% were already V. This mutation is associated with an advantage for the virus and in relation to vaccination, and indeed, the J_{HAC} statistic detects it as a site with a selective pattern.

Temporal comparisons: EN1–EN2 (March 2020 versus March 2021)

The comparison between the English genomes is between samples separated in time (different waves). These comparisons should be considered with caution as the differentiation between samples is very large. Indeed, the mean F_{ST} in all three comparisons (EN1–EN2, EN2–EN3, and EN3–EN4) is above 0.5. However, the sites detected in the three comparisons correspond to sites with a recognized impact on virus fitness.

The genomes in EN1 belong to pre-alpha variants, while the genomes in EN2 are Alpha. The combined EN1–EN2 sample consists of 77 SNPs with a frequency greater than 1%. After the whole-genome analysis, iHDSel found six significant sites for the J_{HAC} test. These sites correspond to six Spike mutations, namely amino acids 501, 570, 681, 716, 982, and 1118 (Table 9). All of them correspond to the characteristic Spike mutations of Alpha [64]. The only one missing is D614G, although it is included in the detected haplotypic regions. The fact that it does not come out as directly significant may be because the program did not use that position as the center of a haplotypic block, as it detected the other sites as more extreme outliers since 614G has a presence of 61%G in EN1 and 99.9% in EN2. However, when the program is run proposing the nucleotide positions corresponding to the amino acid 614 as candidates, the result is significant. Therefore, it seems that the haplotypic region including all these mutations has been detected.

Table 8. Significant J_{Hac} tests (p -val < 0.05) for EN3–SA comparison (with 107 SNPs and sample sizes $n_{EN3} = 5844$, $n_{SA} = 1327$).

EN3–SA		Gene (protein)	AA	%
Block size	Site (+1 +130)		(AA in EN3) (AA in SA)	(p1 p2 ... EN3) : (p1 p2 ... SA)
41	7851	ORF1a (NSP3)	(V A) 2529 (A)	(73 27):(- 100)
11	13 812	ORF1b (NSP12)	(M) 115 (M I)	(100):(42 58)
30	21 846	ORF2 (S)	(I T) 95 (I T)	(76 24):(8 92)
14	21 987	ORF2 (S)	(D G) 142 (D G)	(97 3): (62 38)
11	25 413	ORF3a	(I) 7 (I)	(100):(100)
14	28 282	ORF9 (N)	(D L) 3 (D L)	(99 1):(99 1)

(+1 +130): added to the program output position, the +1 to correct the program indexing to 0 and the +130 to correct the eliminated initial positions.

Table 9. Significant J_{Hac} tests ($p\text{-val} < 0.05$) for EN1–EN2 comparison (with 77 SNPs and sample sizes $n_{EN1} = 4224$, $n_{EN2} = 4152$).

EN1–EN2		Gene (protein)	AA	%
Block size	Site (+1 +130)		(AA in EN1) position (AA in EN2)	(p1 p2 ... EN1) : (p1 p2 ... EN2)
11	23 063	ORF2 (S)	N501Y	(100):(1 99)
11	23 271	ORF2 (S)	A570D	(100):(2 98)
11	23 604	ORF2 (S)	P681H	(100):(1 99)
11	23 709	ORF2 (S)	T716I	(100):(1 99)
11	24 506	ORF2 (S)	S982A	(100):(2 98)
11	24 914	ORF2 (S)	D1118H	(100):(1 99)

(+1 +130): added to the program output position, the +1 to correct the program indexing to 0 and the +130 to correct the eliminated initial positions.

Table 10. Significant J_{Hac} tests ($p\text{-val} < 0.05$) for EN2–EN3 comparison (with 105 SNPs and sample sizes $n_{EN2} = 4152$, $n_{EN3} = 5844$).

EN2–EN3		Gene (protein)	AA
Block size	Site (+1 +130)		(AA in EN2) position (AA in EN3)
18	22 917	ORF2 (S)	L452R
18	22 995	ORF2 (S)	T478K
13	23 604	ORF2 (S)	H681R
12	24 410	ORF2 (S)	D950N
12	28 461	ORF9 (N)	D63G
11	28 881	ORF9 (N)	K203M
13	29 402	ORF9 (N)	D377Y

(+1 +130): added to the program output position, the +1 to correct the program indexing to 0 and the +130 to correct the eliminated initial positions.

Temporal comparisons: EN2–EN3 (March 2021 versus June 2021)

This is a comparison of Alpha (EN2) with Delta (EN3) genomes. The pooled EN2–EN3 sample consists of 105 SNPs with a frequency greater than 1%. After whole-genome analysis, iHDSel found seven significant sites using blocks centered on outliers (Table 10).

These included the substitution of relevant Spike amino acids at sites such as 452, 478, 681, and 950 [69]. For example, the L452R substitution appears to be associated with evasion of the immune response [70]. As well as three sites in the N protein, 63, 203, and 377, which correspond to significant mutations of the delta variant, namely, D63G, R203M, and D377Y [71].

Temporal comparisons: EN3–EN4 (June 2021 versus January 2024)

This is a comparison of Delta genomes (EN3) with Omicron genomes (EN4). The pooled EN3–EN4 sample consists of 239 SNPs with a frequency greater than 1%. After whole-genome analysis, iHDSel identified several sites with F_{ST} greater than 0.99 and 14 of them were in the center of significant blocks (Table 11).

The first site occurs in ORF1a (NSP5) and corresponds to the amino acid change P132H, which is a mutation in a functionally important domain and characteristic of Omicron [72]. The remaining sites presented in Table 11 correspond to core Omicron mutations in Spike [73, 74] including some like S371F, S373P, and S375F, which are related to alterations in binding and entry preference [75, 76] and also the “Kraken” subvariant immune escape F486P [77]. Finally, the synonymous change L18L in ORF7b is within the same haplotypic block as the reversions A82V and I120T in ORF7a, which, when directly contrasted as candidates, were significant.

Spatial comparisons: EN4–SP2

The genomes of both samples are Omicron but the subvariant composition is different (Table 7). The pooled EN4–SP2 sample consists of 218 SNPs with a frequency greater than 1%. After whole-genome analysis, iHDSel identified four significant sites (Table 12).

The change A427V in ORF1a is characteristic mutation of the DV.7.1 Omicron sublineage [64] which is virtually absent in EN4 (0.6%) but has a 28% in SP2 (Table 7) which explains the absence of 427 V in EN4 and the 29%V in SP2. The same scenario applies to A520V in ORF1b. The other two significant sites belong to Spike. The mutation at 445 would be related to the V445H and V445P changes that seem to favor immune evasion of the virus [74, 78] with the presence of 445 V being 30% in SP2 but only 1% in EN4 (Table 12). Finally, L858I is also a characteristic mutation of DV.7.1.

Discussion

In this work, a new statistic called J_{Hac} is proposed to detect genomic patterns compatible with selective sweeps. The statistic is constructed from the interpretation in terms of information from the Price equation [29, 30] and consists of the population stability index applied to the distribution of HACs in two samples. The iHDSel program incorporates the statistic, along with the calculation of haplotype blocks in such a way that each candidate site is located at the center of a block. J_{Hac} appears to work optimally with simulated data, where two diploid populations are subjected to divergent selection under different mutation and recombination conditions. However, if using the program mode that places the outlier sites at the center of the blocks, care must be taken because the false-positive rate increases in bottleneck scenarios. A possible correction in these scenarios is to repeat the calculation with a slightly larger window size.

Real SARS-CoV-2 data have also been used to test J_{Hac} in both spatial and temporal comparisons. Some sites known to impact virus fitness and its ability to promote immune escape have been detected.

The Price equation for comparing genomic patterns

The general formulation of the Price equation describes a change between two populations at any scale, spatial or temporal [31]. The Price equation has been proposed as a unifying principle in evolutionary biology, allowing the formulation and systematization of different evolutionary models and motivating the development of equations and models that reveal invariances and general principles [79, 80]. Here, we have used the selective component of the Price equation, specifically its interpretation in

Table 11. Significant J_{Hac} tests ($p\text{-val} < 0.05$) for the EN3–EN4 comparison (with 239 SNPs and sample sizes $n_{EN3} = 5844$, $n_{EN4} = 3712$).

EN3–EN4		Gene (protein)	AA
Block size	Site (+1 + 130)		(AA in EN3) position (AA in SP2)
11	10 447	ORF1a (NSP5)	P132H
11	22 674	ORF2 (S)	S371F
11	22 679	ORF2 (S)	S373P
11	22 686	ORF2 (S)	S375F
11	22 775	ORF2 (S)	D405N
11	22 786	ORF2 (S)	R408S
11	22 813	ORF2 (S)	K417N
11	22 898	ORF2 (S)	G446S
11	22 992	ORF2 (S)	S477N
11	23 019	ORF2 (S)	F486P
11	23 055	ORF2 (S)	Q498R
11	23 075	ORF2 (S)	Y505H
11	25 000	ORF2 (S)	D1146D
11	27 807	ORF7b	L18L (ORF7a A82V, ORF7a I120T)

(+1 + 130): added to the program output position, the +1 to correct the program indexing to 0 and the +130 to correct the eliminated initial positions.

Table 12. Significant J_{Hac} tests ($p\text{-val} < 0.05$) for the EN4–SP2 comparison (with 218 SNPs and sample sizes $n_{EN4} = 3712$, $n_{SP2} = 221$). Only amino acids with a frequency equal or greater than 1% are indicated.

EN4–SP2		Gene (protein)	AA
Block size	Site (+1 + 130)		(EN4 AA) position (SP2 AA)
11	1545	ORF1a (NSP2)	A427(71%A 29%V)
13	15 026	ORF1b (NSP12)	A520(71%A 29%V)
11	22 895	ORF2 (S)	(51%H 47%P 1%V) 445(31%H 39%P 30%V)
11	24 134	ORF2 (S)	L858(71%L 29%I)

(+1 + 130): added to the program output position, the +1 to correct the program indexing to 0 and the +130 to correct the eliminated initial positions.

terms of information theory [30], which allows the expression of the covariance between fitness and the trait under study in terms of Jeffreys divergence or population stability index. We have defined as a trait the HAC and used Jeffreys divergence to compare the distribution of the trait between two populations. The change in trait distribution would be compatible with the effect of selective sweeps, whether due to divergent or directional selection, depending on whether we are comparing populations in space or time.

Limitations of the J_{Hac} method

The detection of selective sweeps is affected by different evolutionary and demographic scenarios. Throughout the space of the various parameters (mutation, recombination, background and deleterious selection, etc.), it is not difficult to find scenarios that generate an excess of false positives [20, 21]. In our case, we have seen that some evolutionary scenarios, such as bottlenecks, can generate interpopulation genomic patterns that increase the false-positive rate when using automatic window sizes centered on outliers. Although increasing the window size restores control over the false-positive rate, it is possible that other scenarios without positive selection could also alter the HAC patterns. Furthermore, an excessively large window size will cause the recombination effect to dilute the swept signal.

Moreover, as we have already indicated, the method proposed here arises from the informational interpretation of the selective component of the Price equation. However, it is a statistical decomposition based on covariance, and we know that correlation does not imply causation. There is also no *a priori* guarantee that the partition between selection and transmission is additive [81]. Therefore, J_{Hac} is an indirect method that detects a genomic

pattern possibly related to selection but which can also be generated under other circumstances. Hence, the detected sites should be verified through direct methods such as the study of gene function, fitness, etc.

Finally, some genomic patterns of selection correlate with environmental variables, making it difficult to separate both effects [24]. The method proposed here could be combined with other methods that take this correlation into account.

Concluding remarks

There are many statistics for identifying regions of selective sweeps in genomes, see for example [4, 9, 10, 13, 82]. The use of machine learning-based methods to detect selection patterns has been increasing due to their accuracy and ability to handle large amounts of complex data. The underlying idea of all these methods is to use classification algorithms trained with known response data (simulations). That is, if we aim to detect a selection pattern, we train the algorithm with data that we know contains that pattern and with other data without the pattern. Different types of algorithms have been applied: neural networks, extremely randomized trees, and boosting algorithms [9, 13]. A major advantage of these methods is their power and flexibility, partly due to the ease of incorporating new statistics with minimal changes to the structure of the method. Two recent machine learning methods have been designed to detect genomic signatures caused by natural selection, using a supervised multi-statistic machine learning approach [12, 83]. In this work, we have developed a new statistic, J_{Hac} , which, due to its known null distribution, allows us to efficiently and quite accurately test for the existence of genomic patterns compatible with selective sweeps. Therefore, J_{Hac} could be an additional measure to

consider for future AI-based selection detection methods. In addition, J_{Hac} has been incorporated into the iHDSel program (<https://acraaj.webs.uvigo.es/iHDSel.html>) along with an automatic haplotype block detection system, so it can be run independently or in conjunction with the heuristic EOS outlier detection method [34].

Acknowledgements

I gratefully acknowledge all data contributors, that is the Authors and their Originating laboratories responsible for obtaining the specimens, and their Submitting laboratories for generating the genetic sequence and metadata and sharing via the GISAID Initiative, on which the real data example in this article is based.

Funding

This work was supported by Xunta de Galicia (Grupo de Referencia Competitiva, ED431C 2024/22), Ministerio de Ciencia e Innovación (PID2022-137935NB-I00), Centro singular de investigación de Galicia accreditation 2024-2027 (ED431G 2023/07) and ERDF A way of making Europe, and the Marine Science Programme (ThinkInAzul) supported by the Ministerio de Ciencia e Innovación and Xunta de Galicia with funding from the European Union NextGenerationEU (PRTR-C17.11) and European Maritime and Fisheries Fund. Funding for open access charge: Universidade de Vigo/CISUG.

Data availability

The simulated data underlying this article are available in [https://gisaid.org/EPI_SET_240212oy](https://zenodo.org/records/14269530?token=eyJhbGciOiJIUzUxMiJ9.eYpZC16ImU2YzU3YmM2LThiZG5EtdNDQ1My1hZWQ0LTdiZWYxZWZiN2EzYSIsImRhGEiOnt9LCljYyW5kb20iOiIzZmMmViMjk1Yjc4ODllNTFkNjM5MWFfjMzMxZWVjNjg2NyJ9.4VSQkfoj-GdPaQ75yKhr_Hnoc2LyR5F19TpjUm7HwBjyNXeWtk5s5afca_R5ruiKIh16zYAadLE3XFlUyWywData associated with 30,274 SARS-CoV-2 genomes are available on GISAID at the following links: <a href=), https://gisaid.org/EPI_SET_240213ze, https://gisaid.org/EPI_SET_240212vc, https://gisaid.org/EPI_SET_240213oc, https://gisaid.org/EPI_SET_240212xr, https://gisaid.org/EPI_SET_240213br, https://gisaid.org/EPI_SET_240212sa and in the links indicated within the article.

References

- Smith JM, Haigh J. The hitch-hiking effect of a favourable gene. *Genet Res* 1974;**23**:23–35.
- Kaplan NL, Hudson RR, Langley CH. The “hitchhiking effect” revisited. *Genetics* 1989;**123**:887–99.
- Berry AJ, Ajioka J, Kreitman M. Lack of polymorphism on the drosophila fourth chromosome resulting from selection. *Genetics* 1991;**129**:1111–7.
- Stephan W. Selective sweeps. *Genetics* 2019;**211**:5–13. <https://doi.org/10.1534/genetics.118.301319>.
- Johri P, Stephan W, Jensen JD. Soft selective sweeps: addressing new definitions, evaluating competing models, and interpreting empirical outliers. *PLoS Genet* 2022;**18**:e1010022. <https://doi.org/10.1371/journal.pgen.1010022>.
- Sabeti PC, Varilly P, Fry B, International HapMap Consortium. Genome-wide detection and characterization of positive selection in human populations. *Nature* 2007;**449**:913–8.
- Kimura R, Fujimoto A, Tokunaga K et al. A practical genome scan for population-specific strong selective sweeps that have reached fixation. *PLoS One* 2007;**2**:e286. <https://doi.org/10.1371/journal.pone.0000286>.
- Chen H, Patterson N, Reich D. Population differentiation as a test for selective sweeps. *Genome Res* 2010;**20**:393–402. <https://doi.org/10.1101/gr.100545.109>.
- Horscroft C, Ennis S, Pengelly RJ et al. Sequencing era methods for identifying signatures of selection in the genome. *Brief Bioinform* 2019;**20**:1997–2008. <https://doi.org/10.1093/bib/bby064>.
- Abondio P, Cilli E, Luiselli D. Inferring signatures of positive selection in whole-genome sequencing data: an overview of haplotype-based methods. *Genes (Basel)* 2022;**13**:926. <https://doi.org/10.3390/genes13050926>.
- Amin MR, Hasan M, Arnab SP et al. Tensor decomposition-based feature extraction and classification to detect natural selection from genomic data. *Mol Biol Evol* 2023;**40**:msad216. <https://doi.org/10.1093/molbev/msad216>.
- Arnab SP, Amin MR, DeGiorgio M. Uncovering footprints of natural selection through spectral analysis of genomic summary statistics. *Mol Biol Evol* 2023;**40**:msad157. <https://doi.org/10.1093/molbev/msad157>.
- Panigrahi M, Rajawat D, Nayak SS et al. Landmarks in the history of selective sweeps. *Anim Genet* 2023;**54**:667–88. <https://doi.org/10.1111/age.13355>.
- Whitehouse LS, Schrider DR. Timesweeper: accurately identifying selective sweeps using population genomic time series. *Genetics* 2023;**224**:iyad084. <https://doi.org/10.1093/genetics/iyad084>.
- Kern AD, Schrider DR. DiploS/HIC: an updated approach to classifying selective sweeps. *G3 (Bethesda)* 2018;**8**:1959–70. <https://doi.org/10.1534/g3.118.200262>.
- Lourenço VM, Ogotu JO, Rodrigues RAP et al. Genomic prediction using machine learning: a comparison of the performance of regularized regression, ensemble, instance-based and deep learning methods on synthetic and empirical data. *BMC Genomics* 2024;**25**:152. <https://doi.org/10.1186/s12864-023-09933-x>.
- Delaneau O, Zagury JF, Robinson MR et al. Accurate, scalable and integrative haplotype estimation. *Nat Commun* 2019;**10**:5436. <https://doi.org/10.1038/s41467-019-13225-y>.
- Meier JI, Salazar PA, Kučka M. et al. Haplotype tagging reveals parallel formation of hybrid races in two butterfly species. *Proc Natl Acad Sci USA* 2021;**118**:e2015005118. <https://doi.org/10.1073/pnas.2015005118>.
- Shipilina D, Pal A, Stankowski S et al. On the origin and structure of haplotype blocks. *Mol Ecol* 2023;**32**:1441–57. <https://doi.org/10.1111/mec.16793>.
- Johri P, Aquadro CF, Beaumont M et al. Recommendations for improving statistical inference in population genomics. *PLoS Biol* 2022;**20**:e3001669. <https://doi.org/10.1371/journal.pbio.3001669>.
- Soni V, Johri P, Jensen JD. Evaluating power to detect recurrent selective sweeps under increasingly realistic evolutionary null models. *Evolution* 2023;**77**:2113–27. <https://doi.org/10.1093/evo/lut/qpad120>.
- Soni V, Jensen JD. Temporal challenges in detecting balancing selection from population genomic data. *G3 (Bethesda)* 2024;**14**:jkae069. <https://doi.org/10.1093/g3journal/jkae069>.

23. Galindo J, Carvalho J, Sotelo G et al. Genetic and morphological divergence between littorina *Fabalis* ecotypes in Northern Europe. *J Evol Biol* 2021;**34**:97–113. <https://doi.org/10.1111/jeb.13705>.
24. Folkertsma R, Charbonnel N, Henttonen H et al. Genomic signatures of climate adaptation in bank voles. *Ecol Evol* 2024;**14**: e10886. <https://doi.org/10.1002/ece3.10886>.
25. Pampín M, Casanova A, Fernández C et al. Genetic markers associated with divergent selection against the parasite *Marteilia* *Cochillia* in common cockle (*Cerastoderma Edule*) using transcriptomics and population genomics data. *Front Mar Sci* 2023; **10**:1057106.
26. Vera M, Wilmes SB, Maroso F et al. Heterogeneous microgeographic genetic structure of the common cockle (*Cerastoderma Edule*) in the Northeast Atlantic Ocean: biogeographic Barriers and environmental factors. *Heredity (Edinb)* 2023;**131**:292–305. <https://doi.org/10.1038/s41437-023-00646-1>.
27. Labuda D, Labbé C, Langlois S et al. Patterns of variation in DNA segments upstream of transcription start sites. *Hum Mutat* 2007; **28**:441–50. <https://doi.org/10.1002/humu.20463>.
28. Hussin J, Nadeau P, Lefebvre JF et al. Haplotype allelic classes for detecting ongoing positive selection. *BMC Bioinformatics* 2010; **11**:65.
29. Price GR. Extension of covariance selection mathematics. *Ann Hum Genet* 1972;**35**:485–90.
30. Frank SA. Natural selection. V. How to read the fundamental equations of evolutionary change in terms of information theory. *J Evol Biol* 2012;**25**:2377–96.
31. Frank SA. Universal expressions of population change by the price equation: natural selection, information, and maximum entropy production. *Ecol Evol* 2017;**7**:3381–96. <https://doi.org/10.1002/ece3.2922>.
32. Frank SA. Natural selection. IV. The price equation. *J Evol Biol* 2012;**25**:1002–19.
33. Frank SA. Natural selection. VI. Partitioning the information in fitness and characters by path analysis. *J Evol Biol* 2013; **26**:457–71.
34. Carvajal-Rodríguez A. HacDivSel: two new methods (haplotype-based and outlier-based) for the detection of divergent selection in Pairs of populations. *PLoS One* 2017;**12**:e0175944. <https://doi.org/10.1371/journal.pone.0175944>.
35. Gabián M, Morán P, Saura M et al. Detecting local adaptation between North and South European Atlantic Salmon populations. *Biology (Basel)* 2022;**11**:933. <https://doi.org/10.3390/biology11060933>.
36. Frank SA. Simple unity among the fundamental equations of science. *Philos Trans R Soc Lond B Biol Sci* 2020;**375**:20190351. <https://doi.org/10.1098/rstb.2019.0351>.
37. Kullback S. *Information Theory and Statistics; New Edition*. Mineola, NY: Dover Publications, 1997.
38. Hastie T, Tibshirani R, Friedman JH. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd edn. Berlin, Germany: Springer Science & Business Media, 2009.
39. Manning CD, Raghavan P, Schütze H. *Introduction to Information Retrieval*. Cambridge: Cambridge University Press, 2008.
40. Lewontin RC. The interaction of selection and linkage. I. General considerations; heterotic models. *Genetics* 1964;**49**:49–67.
41. Carvajal-Rodríguez A. GENOMEPOP: a program to simulate genomes in populations. *BMC Bioinformatics* 2008;**9**:223.
42. Barton NH. The effect of hitch-hiking on neutral genealogies. *Genet Res* 1998;**72**:123–33. <https://doi.org/10.1017/S0016672398003462>.
43. Thornton KR, Jensen JD. Controlling the false-positive rate in Multilocus genome scans for selection. *Genetics* 2007; **175**:737–50.
44. Harris RB, Sackman A, Jensen JD. On the unfounded enthusiasm for soft selective sweeps II: examining recent evidence from humans, flies, and viruses. *PLoS Genet* 2018;**14**:e1007859. <https://doi.org/10.1371/journal.pgen.1007859>.
45. Terbot JW, Johri P, Liphardt SW et al. Developing an appropriate evolutionary baseline model for the study of SARS-CoV-2 patient samples. *PLoS Pathog* 2023;**19**:e1011265. <https://doi.org/10.1371/journal.ppat.1011265>.
46. Terbot JW, Cooper BS, Good JM et al. A simulation framework for modeling the within-patient evolutionary dynamics of SARS-CoV-2. *Genome Biol Evol* 2023;**15**:evad204. <https://doi.org/10.1093/gbe/evad204>.
47. Crow JF, Kimura M. *An Introduction to Population Genetics Theory*. New York, NY: Harper & Row, 1970.
48. Roughgarden J. *Theory of Population Genetics and Evolutionary Ecology: An Introduction*. New York, NY: Macmillan, 1996.
49. Khare S, Gurry C, Freitas L et al. GISAID's role in pandemic response. *China CDC Wkly* 2021;**3**:1049–51. <https://doi.org/10.46234/ccdcw2021.255>.
50. Aksamentov I, Roemer C, Hodcroft EB et al. Nextclade: clade assignment, mutation calling and quality control for viral genomes. *Joss* 2021;**6**:3773. <https://doi.org/10.21105/joss.03773>.
51. Mathieu E, Ritchie H, Rodés-Guirao L et al. Coronavirus pandemic (COVID-19). *Our World in Data*. 2020. Published online at OurWorldinData.org. Retrieved from: <https://ourworldindata.org/coronavirus> [Online Resource]
52. Brüssow H. COVID-19: omicron—the latest, the least virulent, but probably not the last variant of concern of SARS-CoV-2. *Microb Biotechnol* 2022;**15**:1927–39. <https://doi.org/10.1111/1751-7915.14064>.
53. Wang X, Lu L, Jiang S. SARS-CoV-2 omicron subvariant BA.2.86: limited potential for global spread. *Signal Transduct Target Ther* 2023;**8**:439–3. <https://doi.org/10.1038/s41392-023-01712-0>.
54. Wang X, Lu L, Jiang S. SARS-CoV-2 evolution from the BA.2.86 to JN.1 variants: unexpected consequences. *Trends Immunol* 2024; **45**:81–4. <https://doi.org/10.1016/j.it.2024.01.003>.
55. Kaku Y, Okumura K, Padilla-Blanco M et al.; Genotype to Phenotype Japan (G2P-Japan) Consortium. Virological characteristics of the SARS-CoV-2 JN.1 variant. *Lancet Infect Dis* 2024; **24**:e82. [https://doi.org/10.1016/S1473-3099\(23\)00813-7](https://doi.org/10.1016/S1473-3099(23)00813-7).
56. O'Toole Á, Hill V, Pybus OG et al. Tracking the international spread of SARS-CoV-2 lineages B.1.1.7 and B.1.351/501Y-V2 with Grinch [version 2; peer review: 3 approved]. *Wellcome Open Res* 2021;**6**:121. <https://doi.org/10.12688/wellcomeopenres.16661.2>.
57. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol* 2013;**30**:772–80. <https://doi.org/10.1093/molbev/mst010>.
58. Katoh K, Rozewicki J, Yamada KD. MAFFT online service: multiple sequence alignment, interactive sequence choice and visualization. *Brief Bioinform* 2019;**20**:1160–6. <https://doi.org/10.1093/bib/bbx108>.
59. van Dorp L, Acman M, Richard D et al. Emergence of genomic diversity and recurrent mutations in SARS-CoV-2. *Infect Genet Evol* 2020;**83**:104351. <https://doi.org/10.1016/j.meegid.2020.104351>.
60. Kumar S, Stecher G, Li M et al. MEGA X: molecular evolutionary genetics analysis across computing platforms. *Mol Biol Evol* 2018;**35**:1547–9. <https://doi.org/10.1093/molbev/msy096>.
61. Jankowiak M, Obermeyer FH, Lemieux JE. Inferring selection effects in SARS-CoV-2 with Bayesian viral allele selection. *PLoS*

- Genet 2022;**18**:e1010540. <https://doi.org/10.1371/journal.pgen.1010540>.
62. Garcia I, Lee Y, Brynildsrud O et al. Tracing the adaptive evolution of SARS-CoV-2 during vaccine roll-out in Norway. *Virus Evol* 2024;**10**:vead081. <https://doi.org/10.1093/ve/vead081>.
 63. Kelly JA, Woodside MT, Dinman JD. Programmed -1 ribosomal frameshifting in coronaviruses: a therapeutic target. *Virology* 2021;**554**:75–82. <https://doi.org/10.1016/j.virol.2020.12.010>.
 64. Gangavarapu K, Latif AA, Mullen JL et al.; GISAID Core and Curation Team. Outbreak. Info genomic reports: scalable and dynamic surveillance of SARS-CoV-2 variants and mutations. *Nat Methods* 2023;**20**:512–22. <https://doi.org/10.1038/s41592-023-01769-3>.
 65. Cai HY, Cai A. SARS-CoV2 spike protein gene variants with N501T and G142D mutation-dominated infections in mink in the united States. *J Vet Diagn Invest* 2021;**33**:939–42. <https://doi.org/10.1177/10406387211023481>.
 66. Dhawan M, Sharma A, Thakur N et al. Delta variant (B.1.617.2) of SARS-CoV-2: mutations, impact, challenges and possible solutions. *Hum Vaccin Immunother* 2022;**18**:2068883. <https://doi.org/10.1080/21645515.2022.2068883>.
 67. Kannan SR, Spratt AN, Sharma K et al. Omicron SARS-CoV-2 variant: unique features and their impact on pre-existing antibodies. *J Autoimmun* 2022;**126**:102779. <https://doi.org/10.1016/j.jaut.2021.102779>.
 68. Mahmood TB, Hossain MI, Mahmud S et al. Missense mutations in spike protein of SARS-CoV-2 delta variant contribute to the alteration in viral structure and interaction with HACE2 receptor. *Immun Inflamm Dis* 2022;**10**:e683. <https://doi.org/10.1002/iid3.683>.
 69. Kannan SR, Spratt AN, Cohen AR et al. Evolutionary analysis of the delta and delta plus variants of the SARS-CoV-2 viruses. *J Autoimmun* 2021;**124**:102715. <https://doi.org/10.1016/j.jaut.2021.102715>.
 70. He P, Liu B, Gao X et al. SARS-CoV-2 delta and omicron variants evade population antibody response by mutations in a single spike epitope. *Nat Microbiol* 2022;**7**:1635–49. <https://doi.org/10.1038/s41564-022-01235-4>.
 71. Bhattacharya M, Chatterjee S, Sharma AR et al. Delta variant (B.1.617.2) of SARS-CoV-2: current understanding of infection, transmission, immune escape, and mutational landscape. *Folia Microbiol (Praha)* 2023;**68**:17–28. <https://doi.org/10.1007/s12223-022-01001-3>.
 72. Hossain A, Akter S, Rashid AA et al. Unique mutations in SARS-CoV-2 omicron subvariants' non-spike proteins: potential impacts on viral pathogenesis and host immune evasion. *Microb Pathog* 2022;**170**:105699. <https://doi.org/10.1016/j.micpath.2022.105699>.
 73. Basheer A, Zahoor I, Yaqub T. Genomic architecture and evolutionary relationship of BA.2.75: a Centaurus subvariant of omicron SARS-CoV-2. *PLoS One* 2023;**18**:e0281159. <https://doi.org/10.1371/journal.pone.0281159>.
 74. Chen S, Huang Z, Guo Y et al. Evolving spike mutations in SARS-CoV-2 omicron variants facilitate evasion from breakthrough infection-acquired antibodies. *Cell Discov* 2023;**9**:86. <https://doi.org/10.1038/s41421-023-00584-6>.
 75. Hu B, Chan JF-W, Liu H et al. Spike mutations contributing to the altered entry preference of SARS-CoV-2 omicron BA.1 and BA.2. *Emerg Microbes Infect* 2022;**11**:2275–87.
 76. Zheng B, Xiao Y, Tong B et al. S373P mutation stabilizes the receptor-binding domain of the spike protein in omicron and promotes binding. *JACS Au* 2023;**3**:1902–10. <https://doi.org/10.1021/jacsau.3c00142>.
 77. Parums DV. Editorial: the XBB.1.5 ('Kraken') subvariant of omicron SARS-CoV-2 and its rapid global spread. *Med Sci Monit* 2023;**29**:e939580. <https://doi.org/10.12659/MSM.939580>.
 78. Ao D, He X, Hong W et al. The rapid rise of SARS-CoV-2 omicron subvariants with immune evasion properties: XBB.1.5 and BQ.1.1 subvariants. *MedComm (2020)* 2023;**4**:e239. <https://doi.org/10.1002/mco2.239>.
 79. Luque VJ. One equation to rule them all: a philosophical analysis of the price equation. *Biol Philos* 2017;**32**:97–125. <https://doi.org/10.1007/s10539-016-9538-y>.
 80. Luque VJ, Baravalle L. The mirror of physics: on how the price equation can unify evolutionary biology. *Synthese* 2021;**199**:12439–62. <https://doi.org/10.1007/s11229-021-03339-6>.
 81. Okasha S, Otsuka J. The price equation and the causal analysis of evolutionary change. *Philos Trans R Soc Lond B Biol Sci* 2020;**375**:20190365. <https://doi.org/10.1098/rstb.2019.0365>.
 82. Horscroft C, Pengelly R, Sluckin TJ et al. Zalpha: an R package for the identification of regions of the genome under selection. *Joss* 2020;**5**:2638.
 83. Lauterbur ME, Munch K, Enard D. Versatile detection of diverse selective sweeps with flex-sweep. *Mol Biol Evol* 2023;**40**:msad139. <https://doi.org/10.1093/molbev/msad139>.