

1 **Detecting ancient positive selection in humans using** 2 **extended lineage sorting**

3 Stéphane Peyrégne*, Michael James Boyle, Michael Dannemann, Kay Prüfer*

4 Department of Evolutionary Genetics, Max Planck Institute for Evolutionary Anthropology, 04103
5 Leipzig, Germany.

6 *Corresponding authors: stephanepeyregne@gmail.com; pruefer@eva.mpg.de

7 Key words: human evolution; archaic hominins; positive selection.

8 **ABSTRACT**

9 Natural selection that affected modern humans early in their evolution has likely shaped some of the
10 traits that set present-day humans apart from their closest extinct and living relatives. The ability to
11 detect ancient natural selection in the human genome could provide insights into the molecular basis
12 for these human-specific traits. Here, we introduce a method for detecting ancient selective sweeps by
13 scanning for extended genomic regions where our closest extinct relatives, Neandertals and
14 Denisovans, fall outside of the present-day human variation. Regions that are unusually long indicate
15 the presence of lineages that reached fixation in the human population faster than expected under
16 neutral evolution. Using simulations we show that the method is able to detect ancient events of
17 positive selection and that it can differentiate those from background selection. Applying our method
18 to the 1000 Genomes dataset, we find evidence for ancient selective sweeps favoring regulatory
19 changes and present a list of genomic regions that are predicted to underlie positively selected human
20 specific traits.

21

1 INTRODUCTION

2 Modern humans differ from their closest extinct relatives, Neandertals, in several aspects, including
3 skeletal and skull morphology (Weaver 2009), and may also differ in other traits that are not preserved
4 in the archeological record (Laland et al. 2010; Varki et al. 2008). Natural selection may have played a
5 role in fixing these traits on the modern human lineage. However, the selection events driving the
6 fixation would have been restricted to a specific timeframe, extending from the split between archaic
7 and modern humans ca. 650,000 years ago to the split of modern human populations from each other
8 around 100,000 years ago (Prüfer et al. 2014). While methods exist, that can be used to scan the
9 genome for the remnants of past or ongoing positive selection (Lemey et al. 2009; Nielsen et al. 2007),
10 current methods have limited power to detect positive selection on the human lineage that acted during
11 this older timeframe (see Sabeti et al. 2006 for a review on detection methods and their timeframes):
12 an unusually high ratio of functional changes to non-functional changes, such as the dn/ds test,
13 requires millions of years and often multiple events of selection to generate detectable signals
14 (Kryazhimskiy and Plotkin 2008), while unusual patterns of genetic diversity between individuals and
15 populations (e.g. extended homozygosity, Tajimas D , F_{st}) are most powerful during the selective
16 sweep or shortly after (Oleksyk et al. 2010; Sabeti et al. 2006). A favorable substitution is not
17 expected to leave a mark on linked neutral variation beyond 250,000 years in humans (Przeworski
18 2002, 2003).

19 The genome sequencing of archaic humans (Neandertals and Denisovans) to high coverage (Meyer et
20 al. 2012; Prüfer et al. 2014) has spawned new methods to investigate the genetic basis of modern
21 human traits that are not shared by the archaics (Pääbo 2014). One method, called 3P-CLR, models
22 allele frequency changes before and after the split of two populations using the archaic genomes as an
23 outgroup (Racimo 2016). 3P-CLR outperforms previous methods in the detection of older event of
24 selection (up to 150,000 years ago, Figure 2 from Racimo 2016) but has little power to detect events
25 older than 200,000 years ago in modern humans. A second method applied an approximate Bayesian
26 computation on patterns of homozygosity and haplotype diversity around alleles that reach fixation

1 (Racimo et al. 2014). Although, this approach expands our ability to investigate older time frames, this
2 signal of selection also fades over time and events of positive selection older than 300kya become
3 undetectable.

4 Based on a method introduced by Green et al. (2010), Prüfer et al. (2014) presented a hidden Markov
5 model that identifies regions in the genome where the Neandertal and Denisovan individuals fall
6 outside of present-day human variation (i.e. the archaic lineages fall basal compared to all present-day
7 humans), and applied the model to detect selective sweeps on the modern human lineage. Regions that
8 are unusually long are candidates for ancient selective sweeps as variants are likely to have swept
9 rapidly to fixation, dragging along with them large parts of the chromosomes that did not have time to
10 be broken up by recombination. While this method is, in principle, expected to be able to detect events
11 as old as the modern human split from Neandertals and Denisovans, this power was never formally
12 tested and it has several other shortcomings. First, the method was limited to modern human
13 polymorphisms, ignoring the additional information given by fixed substitutions. Second, the method
14 does not fit parameters to the data, but requires these parameters to be estimated through coalescent
15 simulations.

16 Here, we introduce a refined version of this method, called ELS (Extended Lineage Sorting), that
17 models explicitly the longer regions produced under selection, and includes the fixed differences
18 between archaic and modern human genomes as an additional source of information. The ELS method
19 also takes advantage of an Expectation-Maximization algorithm to estimate the model parameters
20 from the data itself, making it free from assumptions regarding human demographic history.

21 To evaluate the power of the ELS method to detect ancient selective sweeps we tested its performance
22 under scenarios of background selection and neutrality. Finally, we present an updated list of
23 candidate regions that likely underwent positive selection on the modern human lineage since the split
24 from the common ancestor with Neandertals and Denisovans.

1 RESULTS

2 Selection causes extended lineage sorting between closely related populations

3 The ancestors of modern humans split from the ancestors of Neandertals and Denisovans between
4 450,000 and 750,000 years ago (Prüfer et al. 2014). Because the two newly formed descendant groups
5 sampled the genetic variation from the ancestral population, a derived variant can be shared between
6 some members of both groups, while other individuals show the ancestral variant. At these positions,
7 some lineages from one group share a more recent common ancestor with some lineages in the other
8 group than within the same group (Rosenberg 2002), a phenomenon called incomplete lineage sorting
9 (Figure 1A).

10 Eventually, a derived allele may reach fixation as part of a region that has not been unlinked by
11 recombination. In these regions all descendants will derive from one common ancestor and any lineage
12 from the other population will constitute an out-group, i.e. all lineages are sorted. Because of
13 recombination, the human genome is a mosaic of independent evolutionary histories and the process
14 of lineage sorting is expected to randomly affect regions, until, ultimately, all lineages will be sorted.
15 In the case of modern humans, only a fraction of the regions in the genome are expected to show
16 lineage sorting (Prüfer et al. 2014), and the genome can be partitioned into regions where an archaic
17 lineage falls either within the variation of modern humans (internal region) or outside of the human
18 variation (external region) (Figure 1B).

19 While lineage sorting can occur under neutrality, selection on the modern human branch is expected to
20 always lead to external regions as long as the selective sweep finished. In cases where the selective
21 sweep is sufficiently strong, there will not be sufficient time for recombination to break the linkage
22 with neighboring sites and a large region will reach fixation (extended lineage sorting, ELS, Figure
23 1C). In contrast, selection on standing variation may fail to generate such large regions, since
24 recombination can act on the haplotype(s) with the prospective advantageous variant before selection

1 sets in. We note that neither demography nor selection on the archaic lineage affect the lineage sorting
2 within modern humans and thus the power to detect selective sweeps.

3 [Expected Incomplete Lineage Sorting among Humans to Archaics](#)

4 We used coalescent simulations to determine the incidence and expected length of regions resulting
5 from incomplete lineage sorting in modern humans. Using a model of human demographic history
6 (Yang et al. 2014), we estimated the fraction of lineage sorting in modern humans in regards to
7 Neandertals and Denisovans. In simulations with 370 African chromosomes, and assuming a uniform
8 recombination rate, about 10% of the archaic genome is more divergent than the time to the most
9 recent common ancestor of all sampled human variation. The length of the external regions is expected
10 to be about 0.0016 cM (95%-CI: 0.001-0.0095 cM; e.g. 1-9.5kb for a recombination rate of 1cM/Mb)
11 with the longest regions in the order of 0.02 cM. In contrast, internal regions are expected to be 0.012
12 cM long (95%-CI: 0.0097-0.07 cM).

13 [Minimum Strength of Selection to Produce Detectable Sweep Signals](#)

14 We investigated the range of selection coefficients that could have led to the fixation of a lineage after
15 the split with the Archaic hominins, but before the differentiation of genetically modern humans about
16 100–120 kyr ago (Li and Durbin 2011) by simulating mutations occurring at different times and
17 evolving with different selection coefficients. While the simulations show that completed selective
18 sweeps could have occurred with selection coefficients as low as 0.0005 (Figure 2A), the length
19 distribution of haplotypes reaching fixation is indistinguishable from neutrality for selection
20 coefficients under 0.001 (Figure 2, B and C). Under neutrality, the average length of external regions
21 was 0.02 cM and remained below 0.03cM for most simulations with a selection coefficient of 0.001.
22 In contrast, external regions longer than 0.1cM were observed for selection coefficients above 0.05.
23 Therefore, detectable signals are expected to be biased towards strong events with a selection
24 coefficient larger than 0.001.

1 Hidden Markov Model to Detect Extended Lineage Sorting

2 To detect regions of Extended Lineage Sorting, we modeled the changes of local genealogies along the
3 genome with a hidden Markov model. We distinguish two types of genealogies, internal or external,
4 depending on whether the archaic lineage falls inside or outside of the human variation respectively
5 (Figure 3A). The model includes a third state corresponding to extended lineage sorting, and external
6 regions produced by this state are required to be longer, on average, than those produced by the
7 external state. The three states are inferred from the state of the archaic allele (ancestral or derived)
8 either at a polymorphic position in modern humans or at a position where modern humans carry a
9 fixed derived variant. In the following, we describe the different statistical properties expected for
10 each type of genealogy.

11 We first consider external regions. At modern human polymorphic sites, the archaic genome is
12 expected to carry the ancestral variant since the derived variant would indicate incomplete lineage
13 sorting. To account for sequencing errors or misassignment of the ancestral state, we allow a
14 probability of 0.01 for carrying the derived allele (see Material and Methods). At sites where the
15 derived allele is fixed, the archaic genome will often carry the derived state, if the fixation event
16 occurred before the split of the archaic from the modern human lineage, or, occasionally, the ancestral
17 state, if the fixation event is more recent and occurred after the split.

18 For internal regions, the archaic is expected to share the derived allele at modern human fixed derived
19 sites, but can carry the ancestral allele in our model to accommodate errors, albeit with low
20 probability. In contrast, at sites that are polymorphic in modern humans, the probabilities of observing
21 the ancestral or the derived allele in the archaic genome will depend on the age of the derived variant,
22 with young variants being less likely to be shared compared to older variants. The frequency of the
23 derived variant in the modern human population can be used as a proxy for its age and the emission
24 probabilities in our model take the modern human derived allele frequency into account (see Material
25 and Methods).

1 We modeled the transition probabilities between internal and external regions (related to the length of
2 the regions) by exponential distributions. The extended lineage sorting state has the same chance of
3 emitting derived alleles as the other external state but is required to have a larger average length. We
4 used the Baum-Welch algorithm (Durbin et al. 1998), an Expectation-Maximization algorithm, to
5 estimate the emission probabilities, and estimate the transition probabilities with a likelihood
6 maximization algorithm.

7 Accuracy of Parameter Estimates and Inferred Genealogies

8 We first investigated the performance of the parameter inference on simulated data under neutral
9 evolution. We found that the estimated probabilities for encountering ancestral/derived alleles in
10 external and internal regions fit the simulated parameters well (on average less than ± 0.08 from
11 simulated under all tested conditions) (Supplemental Figures S1 and S2), while the estimated length of
12 internal and external regions deviate more from the simulated lengths (around 15% overestimate of the
13 mean length, Supplemental Figure S3). However, we found that the model exhibits better accuracy in
14 labelling the correct genealogies with the estimated length parameters compared to the simulated true
15 values (Supplemental Figure S4). This difference seems to originate from the difficulty in accurately
16 detecting very short external regions or internal regions with very few informative sites. We note that
17 detecting selection is not affected by this problem since we are primarily interested in detecting long
18 external regions. Including fixed differences improves the power to assign the correct genealogies
19 compared to a version of the method without this additional source of information (Supplemental
20 Figure S4).

21 We do not expect ELS regions to be detected in our neutral simulations and indeed we found that
22 either the estimated proportion of ELS converged to zero or the maximum likelihood estimate for the
23 length of ELS and external regions converge to the same value (49% and 51% of simulations
24 respectively). A likelihood ratio test comparing a model without the ELS state to the full model with
25 the ELS state also showed no significant improvement with the additional state in almost all neutral

1 simulations (only one likelihood ratio test out of 100 simulations showed a significant improvement
2 after Bonferroni correction for multiple testing).

3 We then evaluated the accuracy of the ELS method to assign the correct genealogy to regions based on
4 sequences obtained through coalescent simulations with selection (Figure 3, B and C). In these
5 simulations, the underlying genealogy at each site along the sequences is known and can be compared
6 to the estimates. To be conservative, we only focus on results with the smallest selection coefficient
7 ($s=0.005$) that produces regions long enough to be detectable. In Figure 3B we show the accuracy for
8 labelling the extended lineage sorting regions dependent on the posterior probability cutoff for the
9 ELS state. The results demonstrate that the model has sufficient power to accurately label sites that
10 experienced selection with a coefficient $s \geq 0.005$ and an occurrence of the beneficial mutation as long
11 as 600,000 years ago.

12 We also used the simulations of positive selection events ($s=0.005$) with two different times at which
13 the beneficial mutation occurred, 300kya and 600kya, to test how often the beneficial simulated
14 variant fall within a detected ELS region (Supplemental Table S1). To put this rate of true positives
15 into perspective, we also counted how many ELS regions did not overlap the selected variant (false
16 positives). A large fraction of selected mutations were detected (87-92%). However, we also found a
17 substantial fraction of false positive ELS regions (10-11%). When restricting detected ELS regions to
18 those that are longer than 0.025cM, we find less than 0.1% false positives compared to 65-68% true
19 positives. Not all simulated regions with a selection coefficient of 0.005 produce ELS regions of this
20 size, so that the rate of true positives for truly long regions is expected to be higher. For all following
21 analysis, we used this minimal length cutoff of 0.025 cM.

22 Role of Background selection

23 Background selection is defined as the constant removal of neutral alleles due to linked deleterious
24 mutations (Charlesworth et al. 1993). In regions of the genome that undergo background selection, a
25 fraction of the population will not contribute to subsequent generations, causing a reduced effective

1 population size. As a consequence, remaining neutral alleles can reach fixation faster than under
2 neutrality, potentially producing unusually long external regions that could be mistaken as signals of
3 positive selection. We investigated the effects of background selection by running forward simulations
4 with parameters that mimic the strength and extent of background selection estimated for the human
5 genome (Messer 2013). While background selection simulations did produce some long outlier
6 regions that fall outside the distribution observed in neutral simulations, most regions are still smaller
7 than regions simulated with positive selection at a conservative selection coefficient of 0.005 (Figure
8 4A). Indeed, among the 1160 external regions detected in our simulations of background selection
9 ($s=0.05$, Figure 4A) only six were labeled as ELS and only three passed the minimal length filter of
10 0.025 cM.

11 Candidate Regions of Positive Selection on the Human Lineage

12 To identify ancient events of positive selection on the human lineage, we applied the ELS method to
13 African genomes from the 1000 genomes project (The 1000 Genomes Project Consortium 2012). We
14 disregarded non-African populations since Neandertal introgression in these populations could mask
15 selective sweeps and lead to false negatives. A model with ELS fits the data significantly better than a
16 model without the ELS state for all chromosomes and for both tested recombination maps (p-value <
17 $1e-8$, Supplemental Table S2).

18 We identified 81 regions of human extended lineage sorting for which both recombination maps
19 support a genetic length greater than 0.025cM (average length: 0.05 cM). Depending on the
20 recombination map, the longest overlap between the maps is 0.12 (African-American map) or 0.17
21 (deCODE map) cM long, which is three to four times longer than the longest regions produced under
22 background selection in our simulations. An additional 233 regions are longer than 0.025cM according
23 to only one recombination map, with 71% of those additional regions showing support for the ELS
24 state using both recombination maps. This suggests that the variation in the candidate set mostly stems
25 from uncertainty about recombination rates. We will refer to the set of 81 regions as the core set

1 (Supplemental File S1) and the set including the 233 putatively selected regions found with just one
2 recombination map as the extended set (314 regions, Supplemental File S2).

3 For completeness, we also ran our model on the X Chromosome and identified 12 additional
4 candidates (43 if we consider candidates found with at least one recombination map), applying a more
5 stringent length cutoff of 0.035 cM to account for the stronger effects of random drift on this
6 chromosome (cf. Material and Methods). Interestingly, we also found a significant increase of
7 posterior probabilities for selection within previously reported regions under potential recurrent
8 selective sweeps in apes (Dutheil et al. 2015; Nam et al. 2015) (Mann-Whitney U one-sided test, P -
9 value $< 2.2e-16$, Supplemental Table S3).

10 The detected selection candidate regions on the autosomes do not show a decrease in B scores
11 (McVicker et al. 2009), a local measure of background selection strength, compared with random
12 regions (Figure 4B; Wilcoxon rank sum test comparing the average B-scores with permuted regions,
13 P -value=0.565, or comparing the lowest B-scores in our regions to permuted regions, P -value=0.504).
14 This suggests that candidate regions are not primarily generated by strong background selection.

15 We compared our candidate regions to the top candidates of 8 previous scans for selection, including
16 iHS, Fst, XP-CLR and HKA (Cagan et al. 2016; Pybus et al. 2014). Using the estimated time to the
17 most recent common ancestor among Africans for each identified region/site, we found that our ELS
18 scan identified significantly older events than other screens (Figure 5, Mann-Whitney U tests,
19 Supplemental Table S4). We found 23 regions from the core set (detected by both recombination
20 maps) overlapping with candidates from previous scans and 68 for the extended set (detected by at
21 least one recombination map); neither overlap is more than expected at random (P -values are 0.06 and
22 0.595 respectively). In contrast, our candidate regions overlap more often candidate regions from 3P-
23 CLR (Racimo 2016) and the ABC approach for detecting ancient selection (Racimo et al. 2014) than
24 expected by chance (P -values <0.05 ; Supplemental Table S5).

1 Biological functions of the candidate regions

2 Since positive selection acts on advantageous phenotypes that are caused by changes to functional
3 elements in the genome, we would expect that our candidate regions overlap functional elements in the
4 genome more often than expected.

5 We first tested this hypothesis by counting the overlap between sweep candidate regions and protein
6 coding genes (Ensembl release 82)(Aken et al. 2016). We find no statistically significant overlap of
7 ELS regions with protein coding genes compared to randomly placed regions of the same size (P -
8 value = 0.671 and 0.124, for core and extended set, respectively; Figure 6A). Previous work has
9 identified 96 proteins that carry human fixed derived non-synonymous changes compared to
10 Neandertal and Denisova, which constitute a particularly interesting subset of potentially functional
11 changes to genes that may have been caused by selective sweeps (Prüfer et al. 2014). We found no
12 overlap between these genes and the core set of sweep candidate regions that were identified by both
13 recombination maps. However, when considering the extended set of sweep candidate regions, 11
14 regions overlapped such genes: *ADSL*, *BBIP1*, *ENTHD1*, *HERC5*, *KATNA1*, *KIF18A*, *NCOA6*,
15 *PRDM10*, *SCAP*, *SLITRK1* and *ZNHIT2*. This overlap is significantly larger than expected by chance
16 (only 2 genes are expected on average; P -value $< 10^{-3}$). In all instances, the candidate regions
17 contained at least one fixed amino acid change. Since fixed changes are part of the information used to
18 infer external regions, it stands to reason that the presence of such a change may bias towards
19 observing an overlap with candidate regions (72/81 core regions and 275/314 regions from the
20 extended set contain fixed changes). However, we note that the overlap with fixed amino acid changes
21 is also significantly larger than the overlap with other fixed changes (963 of 20347 fixed changes fall
22 within candidate regions from the extended set; binomial P -value=0.006).

23 Phenotype may also be influenced by regulatory changes that affect gene expressions. Interestingly,
24 we found a significant enrichment for regions overlapping enhancers and promoters (P -value <0.001

1 and P -value=0.002, respectively; see Figure 6A) when considering the extended set of 314 candidate
2 regions. However, this enrichment was not significant for the smaller core set of candidates.

3 To further investigate the biological function of our regions, we tested for Gene Ontology enrichment
4 in genes within the extended set of regions. No category showed significant enrichment when
5 comparing to randomly placed regions of identical sizes in the genome (see Supplemental Methods).
6 We also assigned genes that overlap our extended dataset to tissues in which they show the
7 significantly highest expression and found again no enrichment (Supplemental Table S6). In an
8 attempt to include potential regulatory changes in the enrichment test, we assigned genes to candidate
9 regions when a region fell upstream or downstream of a gene (see Supplemental Methods). Although
10 many candidate genes that were annotated in this way were expressed highest in the brain or the heart
11 (Odds ratio=2.10 for both tissues), this enrichment is not significant when correcting for gene length
12 and multiple testing (Family-wise error rate=0.336 and 0.997 respectively, Supplemental Table S7).

13 Additional work will be required to investigate the phenotypic consequences of changes in candidate
14 regions for selection. To facilitate this work, we provide an annotated list of fixed or nearly fixed sites
15 on the human lineage that fall within our candidate regions (Supplemental File S3).

16 [Overlap with Neandertal Introgression](#)

17 Introgression from Neandertals and Denisovans into modern humans occurred approximately 37,000
18 to 86,000 years ago (Fu et al. 2014, 2015; Sankararaman et al. 2012, 2016). For those advantageous
19 derived variants that arose on the modern human lineage prior to introgression, we would expect that
20 selection may have acted against the re-introduction of the ancestral variant through admixture. We
21 tested whether this selection may have affected the distribution of Neandertal introgressed DNA
22 around fixed changes in candidate sweep regions. Out of a total of 963 fixed derived variants in
23 Africans overlapping the extended set of sweep regions, 240 (25%) show the ancestral allele in non-
24 Africans and show evidence for re-introduction by admixture using a map of Neandertal introgression
25 (Vernot and Akey 2014). This level of Neandertal ancestry is comparable to the genome-wide fraction

1 of out-of-Africa ancestral alleles at African fixed derived sites (~26%; bootstrap P -value=0.583). We
2 also find no significant reduction in frequency of Neandertal ancestry around candidate substitutions
3 in sweep regions, when comparing one randomly sampled fixed African substitution per region against
4 random regions matched for size and distance to genes (Supplemental Figure S5 and S6).

5 If selection against the re-introduction of an ancestral variant were very strong, selection may have
6 depleted Neandertal ancestry in a large region surrounding the selected allele. Interestingly we find
7 some of our sweep candidate regions that fall within the longest deserts of both Neandertal and
8 Denisova ancestry (Supplemental Table S8) (Vernot et al. 2016). A significantly high number of the
9 core set of regions fall in these deserts (5/81 regions, P -value=0.024), while the extended set shows no
10 significant enrichment (9/314 regions, P -value=0.205).

11 DISCUSSION

12 Many genetic changes set modern humans apart from Neandertals and Denisovans but their functions
13 remain elusive. Most of these changes probably resulted in either no change to the phenotype or to a
14 selectively neutral change. However, in rare instances selection may have favored changes modifying
15 the appearance, behavior and abilities of present-day humans. Unfortunately, current methods to
16 identify selection have limited power to detect such old events of positive selection (Przeworski 2002,
17 2003; Sabeti et al. 2006).

18 Here, we introduce a hidden Markov model to detect ancient selective sweeps based on a signal of
19 extended lineage sorting. Using simulations we were able to show that the method can detect older
20 events of selection as long as the selected variant was sufficiently advantageous. The power to detect
21 older events is due to the fact that the method increases in power with the number of mutations that
22 accumulated after the sweep finished. We also showed that background selection can cause false
23 signals and have chosen a minimum length cutoff on candidate regions. While this cutoff reduces the
24 number of false positives due to background selection, we note that this cutoff is expected to exclude
25 *bona fide* events of positive selection, too.

1 We applied the ELS method to 185 African genomes, the Altai Neandertal genome and the Denisovan
2 genome, and detected 81 candidate regions of selection when requiring a minimum genetic length
3 supported by two independent recombination maps. The uncertainty in the recombination maps has a
4 large effect on our results, as shown by the much larger number of 314 regions identified by either
5 recombination map. Recombination rates over the genome are known to evolve rapidly (Lesecque et
6 al. 2014) and of particular concern are recent changes in recombination rates that make some regions
7 appear larger in genetic length than they were in the past. By comparing the current recombination
8 rates in our regions to recombination rates in the ancestral population of both chimpanzee and humans
9 (Munch et al. 2014), we identified some candidate regions that may have increased in recombination
10 rates (Supplemental Table S9). However, it is currently impossible to date the change in
11 recombination rates confidently and these candidate sweeps may post-date the change.

12 A particular strength of our screen for selective sweeps is the ability to detect older events, as
13 indicated by the estimated power to detect simulated events of positive selection of old age and
14 moderate strength. This sets the ELS method apart from previous approaches that made use of archaic
15 genomes, which were geared towards detecting younger events with an age of less than 300,000 years
16 ago (Racimo 2016; Racimo et al. 2014). Despite this difference, we found significant overlap between
17 the ELS candidates and the candidates identified by these other approaches, while the overlap with
18 other types of positive selection scans is smaller. Among our candidates, 55 are novel candidates (234
19 if considering the extended set) that were not detected in any of the previous screens, including
20 previous versions of the screen without fixed differences (Supplemental Figure S7).

21 While we find no difference in the fraction of genes in selected regions compared to randomly placed
22 regions, we detect an enrichment for enhancers and promoter regions. This result is in agreement with
23 the hypothesis that regulatory changes may play an important role in human-specific phenotypes
24 (Carroll 2003; Enard et al. 2014; King and Wilson 1975), maybe more so than amino-acid changes
25 (Hernandez et al. 2011; see also Enard et al. 2014 and Racimo et al. 2014). Interestingly, several gene
26 candidates falling within sweep regions play a role in the function and development of the brain. A

1 particularly interesting observation is the potential selection on both the ligand *SLIT2* and its receptor
2 *ROBO2*, which reside on Chromosome 4 and 3 respectively (see Supplemental File S3 for an
3 annotated list of changes in those genes). Members of the Roundabout (ROBO) gene family play an
4 important role in guiding developing axons in the nervous system through interactions with the ligands
5 SLITs. SLITs proteins act as attractive or repulsive signals for axons expressing different ROBO
6 receptors. *ROBO2* has been further associated with vocabulary growth (St Pourcain et al. 2014),
7 autism (Suda et al. 2011), and dyslexia (Fisher and DeFries 2002) and is involved in the development
8 of neural circuits related to vocal learning in birds (Wang et al. 2015). Interestingly, *ROBO2* is also in
9 a long desert of both Denisovan and Neandertal ancestry in non-Africans.

10 We also identified interesting brain-related candidates on the X Chromosome, among them *DCX*, a
11 protein controlling neuronal migration by regulating the organization and stability of microtubules
12 (Gleeson et al. 1999). Mutations in this gene can have consequences for the expansion and folding of
13 the cerebral cortex, leading to the “double cortex” syndrome in females and “smooth brain” syndrome
14 in males (Gleeson et al. 1998).

15 We have presented a new approach to detect ancient selective sweeps based on a signal of extended
16 lineage sorting. Applying this approach to modern human data revealed that selection may have acted
17 primarily on regulatory changes. With population level sequencing of non-human species becoming
18 more readily available we anticipate that this approach will help to reveal the targets of ancient
19 selection in other species.

20 MATERIALS AND METHODS

21 Data

22 We used single nucleotide polymorphisms (SNPs) from 185 unrelated Luhya and Yoruba individuals
23 from the 1000 Genomes Project phase I (The 1000 Genomes Project Consortium 2012) together with
24 four ape reference genome assemblies (chimpanzee (panTro3) (Mikkelsen et al. 2005), bonobo
25 (panPan1.1) (Prufer et al. 2012), gorilla (gorGor3) (Scally et al. 2012) and orangutan (ponAbe2)

1 (Locke et al. 2011)) to compile a list of polymorphic and fixed derived changes in Luhya and Yoruba.
 2 Neandertal and Denisova alleles at these positions were extracted from published VCFs (Danecek et
 3 al. 2011) using recommended filters (Prüfer et al. 2014) (see Supplemental Material for further
 4 details). Sites where either Neandertal or Denisova carried a third allele were disregarded.

5 Genetic distances between those positions were calculated using the African-American (Hinch et al.
 6 2011) and the deCODE (Kong et al. 2010) recombination maps (available in Build 37 from
 7 <http://www.well.ox.ac.uk/~anjali/>). Both maps were chosen since they estimate recombination rates
 8 from events that occurred within a few generations before present. Recombination maps based on
 9 older events (i.e. LD based map) can underestimate recombination rates in regions that underwent
 10 recent selective sweeps, potentially masking true signals.

11 Hidden Markov model

12 We would like to estimate for each informative position the probabilities for the three possible
 13 genealogies external (*E*), internal (*I*) and extended lineage sorting (*ELS*) given the observed data.
 14 Formally, and following the notation from Durbin et al. 1998, we calculate $P(\pi_i = k|x)$ where *i*
 15 denotes the position, $k \in \{E, I, ELS\}$ and *x* is the sequence of observations with the *i*th observation
 16 denoted x_i . With the genetic distance *d* between consecutive sites and l_k , the average genetic length of
 17 a region in state *k*, we specify the transition probabilities between identical states as $t_{k,k} = e^{-\frac{d}{l_k}}$.
 18 Transitions from *I* to the states *ELS* and *E* depend on an additional parameter *p*, the proportion of
 19 transitions from *I* to *ELS*, and their probability is given by $t_{I,ELS} = p \left(1 - e^{-\frac{d}{l_I}}\right)$ and $t_{I,E} = (1 -$
 20 $p) \left(1 - e^{-\frac{d}{l_I}}\right)$. Lastly, transitions from the two external states to internal have the probability
 21 $t_{j,I} = 1 - e^{-\frac{d}{l_j}}$, with $j \in \{E, ELS\}$. By construction, transitions between *E* and *ELS* genealogies are
 22 not allowed: it would not be possible to detect such transitions as those two states have the same
 23 statistical properties.

1 The inference further requires the probability for observing an ancestral or derived allele in the archaic
 2 at a site i with a derived allele frequency $f_i > 0$ in modern humans (noted x_i) given that the true
 3 genealogy is $k \in \{I, E, ELS\}$: $e_k(x_i) = P(x_i | \pi_i = k)$. We assume that $\forall x: e_{ELS}(x) = e_E(x)$, i.e. that
 4 both external states give rise to ancestral and derived alleles in the archaic with equal probabilities
 5 given the same observation. Since external regions are not expected to give rise to derived sites when
 6 the derived allele is segregating in modern humans, the only sources for such an observation can be
 7 errors or independent coinciding identical mutations and we define an error rate for external regions:
 8 $\epsilon_E = e_E(x_i = \text{derived}, f_i < 1)$. Similarly fixed derived sites are expected to show the derived allele
 9 in the archaics if the local genealogy is internal and we define an error rate for internal regions:
 10 $\epsilon_I = e_I(x_i = \text{derived}, f_i = 1)$.

11 We compute the posterior probability $P(\pi_i = k | x)$ that an observation x_i came from state k given the
 12 observed sequence x as: $P(\pi_i = k | x) = \frac{P(x, \pi_i = k)}{P(x)}$. $P(x, \pi_i = k) = f_k(i) b_k(i)$ where $f_k(i) =$
 13 $P(x_1 \dots x_i, \pi_i = k)$ and $b_k(i) = P(x_{i+1} \dots x_L | \pi_i = k)$ are the output of the Forward and Backward
 14 algorithms respectively (Durbin et al. 1998; Rabiner 1989). $P(x)$ corresponds to the likelihood of the
 15 data given our model and was also calculated from the Forward algorithm.

16 Parameter estimate

17 We used the Baum-Welch algorithm to estimate all emission probabilities with the exception of ϵ_E ,
 18 the proportion of segregating sites derived in the archaic genome in external regions, due to limited
 19 accuracy in the estimates. We set this last parameter to a value of 0.01, a conservative upper limit on
 20 contamination and sequencing error in the two high-coverage archaic genomes. The Baum-Welch
 21 algorithm was run for a maximum of 40 iterations and the convergence criteria was set to a log-
 22 likelihood maxima difference of less than 10^{-4} .

23 We estimated the remaining parameters (average lengths of regions and the proportion of transitions to
 24 the ELS state) using the derivative free optimization method COBYLA (Powell 1994) as implemented
 25 in the nlopt library (Steven G. Johnson, The NLOpt nonlinear-optimization package) to maximize the

1 log-likelihood values calculated by the Forward algorithm. Convergence was attained in a maximum
2 of 1000 evaluations and the log-likelihood maximization accuracy was set to 10^{-4} . To test for
3 convergence to local maxima, we ran the algorithm twice with different starting points and used the
4 parameters of the run with the highest likelihood to run the re-estimation algorithm a third time
5 starting with those parameters. All three runs gave similar results on all chromosomes.

6 Post-processing

7 The HMM was executed independently on all chromosomes for both Denisova and Neandertal and
8 using the African-American and deCODE recombination maps. An external region was defined as a
9 stretch of high posterior probabilities ($p \geq 0.7$) for the extended lineage sorting state that was
10 uninterrupted by sites with a low probability ($p \leq 0.1$). The two cutoffs on the posterior probabilities
11 were determined by simulating sequences with positive selection ($s=0.005$, 500kya, see below). Sites
12 that were simulated external in both Archaics were labeled as 1 and the remaining sites as 0. The
13 HMM was then run on the simulations. By running a grid-search over possible cutoffs (step-sizes of
14 0.05 for the two parameters) and labeling the HMM output accordingly, we identified the set of chosen
15 parameters by minimizing the root mean square error $\sqrt{\frac{\sum_i (t_i - o_i)^2}{n}}$ with n the number of labelled sites, t_i
16 the true label and o_i the observed label.

17 Simulations

18 We simulated sequences using a model of recent human demography to test the performance of our
19 HMM under different scenarios of neutral evolution, positive selection or background selection. Each
20 simulation consisted of one chimpanzee chromosome, one chromosome from each archaic hominin
21 and 370 human chromosomes, matching the 185 Luhya and Yoruba individuals used in our analysis.
22 For all simulations in this study, a constant mutation rate of 1.45×10^{-8} bp⁻¹.generation⁻¹, a constant
23 recombination rate of 1cM.Mb⁻¹.generation⁻¹ and a generation time of 29 years were assumed. We
24 used estimates of population sizes from (Yang et al. 2014) and population split estimates from (Prüfer
25 et al. 2014) as parameters for the simulated demography (Supplemental Information 1 and 2).

1 Neutral simulations were generated with the coalescent simulator *scrm* (Staab et al. 2014) and give a
2 good match to our observed data when plotting derived allele frequency in modern humans against the
3 proportion of derived alleles in the outgroup (Supplemental Figure S8). Simulations with positive
4 selection were generated with the coalescent simulator *msms* (Ewing and Hermisson 2010) and
5 background selection was explored using forward in time simulations generated by SLiM (Messer
6 2013). Further details on simulation parameters are given in the Supplemental Material.

7 [Age Comparison with other Scans for Selection](#)

8 To compare our sweep screen with previous scans, we downloaded candidate regions from the 1000G
9 positive selection database (Pybus et al. 2014). Only candidates with a *P*-value lower than 0.001 were
10 considered. We added to this set of regions the top reported regions from a HKA scan (Cagan et al.
11 2016). Allele age estimates were obtained from ARGweaver (Rasmussen et al. 2014).

12 *F_{st}*, *iHS* and XP-EHH are site-based statistics which localise sites that may have been selected (Sabeti
13 et al. 2007; Malécot 1948; Voight et al. 2006; Wright 1951), whereas selective scans such as CLR,
14 XP-CLR, Tajima's *D*, Fay & Wu's *H* and HKA identify candidate regions (Chen et al. 2010; Fay and
15 Wu 2000; Hudson et al. 1987; Kim and Stephan 2002; Tajima 1989). In order to compare the age of
16 the selection events, we assumed that the selected variant in candidate regions was the site with the
17 highest frequency. We note that this procedure will underestimate the age of events if the true selected
18 site reached fixation, as often expected for our method; the comparison is thus conservative.

19 [Annotations](#)

20 We annotated candidate regions using protein coding genes from Ensembl (release 82), promoters and
21 enhancers mapped by GenoSTAN (Zacher et al. 2016), a measure of background selection (B-scores)
22 (McVicker et al. 2009). Candidate regions were also overlapped with regions previously suggested to
23 have experienced recurrent selective sweeps in apes on the X Chromosome (Dutheil et al. 2015; Nam
24 et al. 2015), regions of Neandertal ancestry (Sankararaman et al. 2014; Vernot et al. 2014) and long
25 regions devoid of Neandertal and Denisova ancestry (Vernot et al. 2016).

1 To statistically test the overlap of our regions with these annotations, we randomly placed regions of
2 similar physical sizes in the parts of the genome that passed our quality filters. Quality filtered regions
3 that were smaller than the longest gap present in our candidate ELS regions were regarded as
4 sufficiently short to not prohibit the placement of regions.

5 Changes of recombination rates along the human lineage could limit our power to detect selected
6 regions, and we used an ancestral recombination map of the human-chimpanzee ancestor to annotate
7 top candidate regions (Supplemental Table S9) (Munch et al. 2014).

8 Finally, we further characterized fixed or nearly fixed human-specific changes within the candidate
9 regions using annotations of histone marks (enhancers, promoters), eQTLs, transcription factor
10 binding sites and conservation scores (Supplemental File S3).

11 SOFTWARE AVAILABILITY

12 The software and input files used in this study have been made available through the website
13 <http://bioinf.eva.mpg.de/ELS/> and <https://github.com/StephanePeyregne/ELS/>. A version of the source
14 code is also available as Supplemental Code in the online version of this article.

15 ACKNOWLEDGMENTS

16 We would like to thank Michael Lachmann for early discussions about the design of the study, Janet
17 Kelso and Mark Stoneking for many useful comments on the manuscript, Svante Pääbo, Udo Stenzel,
18 Fernando Racimo, Amy Ko and Adam Siepel for helpful discussions, Julien Peyrégne for help in
19 implementing the software, and Matthias Ongyerth, Manjusha Chintalapati, Steffi Grote and Christoph
20 Theunert for help with earlier analysis. We are also grateful for the comments and suggestions of three
21 anonymous reviewers that helped to improve this manuscript. This research was funded by the Max
22 Planck Society, the Paul G. Allen Family foundation and the European Research Council (grant
23 agreement no. 694707).

1 DISCLOSURE DECLARATION

2 The authors declare no competing financial interests.

3 AUTHOR CONTRIBUTIONS

4 SP implemented the method. SP, MJB and MD analyzed data. SP, MJB, MD and KP interpreted the
5 results. KP designed the study. SP and KP wrote the manuscript with input from all authors.

6 REFERENCES

- 7 The 1000 Genomes Project Consortium. 2012. An Integrated Map of Genetic Variation from 1,092
8 Human Genomes. *Nature* **491**: 56–65.
- 9 Anders S, Huber W. 2010. DESeq: Differential Expression Analysis for Sequence Count Data.
10 *Genome biology* **11**: R106. doi: 10.1186/gb-2010-11-10-r106.
- 11 Cagan A, Theunert C, Laayouni H, Santpere G, Pybus M, Casals F, Prüfer K, Navarro A, Marques-
12 Bonet T, Bertranpetit J, et al. 2016. Natural Selection in the Great Apes. *Molecular Biology and*
13 *Evolution* **33**: 3268-3283.
- 14 Carroll SB. 2003. Genetics and the Making of Homo Sapiens. *Nature* **422**: 849–857.
- 15 Charlesworth B, Morgan MT, Charlesworth D. 1993. The Effect of Deleterious Mutations on Neutral
16 Molecular Variation. *Genetics* **134**: 1289–1303.
- 17 Chen H, Patterson N, Reich D. 2010. Population Differentiation as a Test for Selective Sweeps.
18 *Genome Research* **20**: 393–402.
- 19 Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE, Lunter G,
20 Marth GT, Sherry ST, et al. 2011. The Variant Call Format and VCFtools. *Bioinformatics* **27**:
21 2156–2158.
- 22 Durbin R, Eddy S, Krogh A, Mitchison G. 1998. Biological Sequence Analysis: Probabilistic Models

- 1 of Proteins and Nucleic Acids. Cambridge University Press, Cambridge.
- 2 Dutheil JY, Munch K, Nam K, Mailund T, Schierup MH. 2015. Strong Selective Sweeps on the X
3 Chromosome in the Human-Chimpanzee Ancestor Explain Its Low Divergence. *PLoS Genetics*
4 **11**: e1005451. doi: 10.1371/journal.pgen.1005451.
- 5 Enard D, Messer PW, Petrov DA. 2014. Genome-Wide Signals of Positive Selection in Human
6 Evolution. *Genome Research* **24**: 885–895.
- 7 Ewing G, Hermisson J. 2010. MSMS: A Coalescent Simulation Program Including Recombination,
8 Demographic Structure and Selection at a Single Locus.” *Bioinformatics* **26**: 2064–2065.
- 9 Fay JC, Wu CI. 2000. Hitchhiking under Positive Darwinian Selection. *Genetics* **155**: 1405–1413.
- 10 Fisher SE, DeFries JC. 2002. Developmental Dyslexia: Genetic Dissection of a Complex Cognitive
11 Trait. *Nature Reviews Neuroscience* **3**: 767–780.
- 12 Fu Q, Li H, Moorjani P, Jay F, Slepchenko SM, Bondarev AA, Johnson PL, Aximu-Petri A, Prüfer K,
13 de Filippo C, et al. 2014. Genome sequence of a 45,000-year-old modern human from western
14 Siberia. *Nature* **514**: 445–449.
- 15 Fu Q, Hajdinjak M, Moldovan OT, Constantin S, Mallick S, Skoglund P, Patterson N, Rohland N,
16 Lazaridis I, Nickel B, et al. 2015. An early modern human from Romania with a recent
17 Neanderthal ancestor. *Nature* **524**: 216–219.
- 18 Gleeson JG, Allen KM, Fox JW, Lamperti ED, Berkovic S, Scheffer I, Cooper EC, Dobyns WB,
19 Minnerath SR, Ross ME, et al. 1998. Doublecortin, a Brain-Specific Gene Mutated in Human X-
20 Linked Lissencephaly and Double Cortex Syndrome, Encodes a Putative Signaling Protein. *Cell*
21 **92**: 63–72.
- 22 Gleeson JG, Lin PT, Flanagan LA, Walsh CA. 1999. Doublecortin Is a Microtubule-Associated
23 Protein and Is Expressed Widely by Migrating Neurons. *Neuron* **23**: 257–271.

- 1 Green RE, Krause J, Briggs AW, Maricic T, Stenzel U, Kircher M, Patterson N, Li H, Zhai W, Fritz
2 MH, et al. 2010. A draft sequence of the Neandertal genome. *Science* **328**: 710-722.
- 3 Hernandez RD, Kelley JL, Elyashiv E, Melton SC, Auton A, McVean G, 1000 Genomes Project, Sella
4 G, Przeworski M. 2011. Classic Selective Sweeps Were Rare in Recent Human Evolution.
5 *Science* **331**: 920–924.
- 6 Hinch AG, Tandon A, Patterson N, Song Y, Rohland N, Palmer CD, Chen GK, Wang K, Buxbaum
7 SG, Akylbekova EL, et al. 2011. The Landscape of Recombination in African Americans. *Nature*
8 **476**: 170–175.
- 9 Hudson RR, Kreitman M, Aguadé M. 1987. A Test of Neutral Molecular Evolution Based on
10 Nucleotide Data. *Genetics* **116**: 153–159.
- 11 Kim Y, Stephan W. 2002. Detecting a Local Signature of Genetic Hitchhiking along a Recombining
12 Chromosome. *Genetics* **160**: 765–777.
- 13 King MC, Wilson AC. 1975. Evolution at Two Levels in Humans and Chimpanzees. *Science* **188**:
14 107–116.
- 15 Kong A, Thorleifsson G, Gudbjartsson DF, Masson G, Sigurdsson A, Jonasdottir A, Walters GB,
16 Jonasdottir A, Gylfason A, Kristinsson KT, et al. 2010. Fine-Scale Recombination Rate
17 Differences between Sexes, Populations and Individuals. *Nature* **467**: 1099–1103.
- 18 Kryazhimskiy S, Plotkin JB. 2008. The Population Genetics of dN / dS . *PLoS Genetics* **4**: e1000304.
19 doi: 10.1371/journal.pgen.1000304.
- 20 Laland KN, Odling-Smee J, Myles S. 2010. How Culture Shaped the Human Genome: Bringing
21 Genetics and the Human Sciences Together. *Nature Reviews Genetics* **11**: 137–148.
- 22 Lesecque Y, Glémin S, Lartillot N, Mouchiroud D, Duret L. 2014. The Red Queen Model of
23 Recombination Hotspots Evolution in the Light of Archaic and Modern Human Genomes. *PLoS*
24 *Genetics* **10**: e1004790. doi: 10.1371/journal.pgen.1004790.

- 1 Li H, Durbin R. 2011. Inference of Human Population History from Individual Whole-Genome
2 Sequences. *Nature* **475**: 493–496.
- 3 Locke DP, Hillier LW, Warren WC, Worley KC, Nazareth LV, Muzny DM, Yang SP, Wang ZY,
4 Chinwalla AT, Minx P, et al. 2011. Comparative and demographic analysis of orang-utan
5 genomes. *Nature* **469**: 529-533.
- 6 Malécot G. 1948. Les mathématiques de l'hérédité. Masson & Cie, Paris, France.
- 7 McVicker G, Gordon D, Davis C, Green P. 2009. Widespread Genomic Signatures of Natural
8 Selection in Hominid Evolution. *PLoS Genetics* **5**: e1000471. doi:
9 10.1371/journal.pgen.1000471.
- 10 Messer PW. 2013. SLiM: Simulating Evolution with Selection and Linkage. *Genetics* **194**: 1037–
11 1039.
- 12 Meyer M, Kircher M, Gansauge MT, Li H, Racimo F, Mallick S, Schraiber JG, Jay F, Prüfer K, de
13 Filippo C, et al. 2012. A High-Coverage Genome Sequence from an Archaic Denisovan
14 Individual. *Science* **338**: 222–226.
- 15 Mikkelsen TS, Hillier LW, Eichler EE, Zody MC, Jaffe DB, Yang SP, Enard W, Hellmann I,
16 Lindblad-Toh K, Altheide TK et al. 2005. Initial sequence of the chimpanzee genome and
17 comparison with the human genome. *Nature* **437**: 69-87.
- 18 Munch K, Mailund T, Dutheil JY, Schierup MH. 2014. A Fine-Scale Recombination Map of the
19 Human-Chimpanzee Ancestor Reveals Faster Change in Humans than in Chimpanzees and a
20 Strong Impact of GC-Biased Gene Conversion. *Genome Research* **24**: 467–474.
- 21 Nam K, Munch K, Hobolth A, Dutheil JY, Veeramah KR, Woerner AE, Hammer MF, Great Ape
22 Genome Diversity Project, Mailund T, et al. 2015. Extreme Selective Sweeps Independently
23 Targeted the X Chromosomes of the Great Apes. *Proceedings of the National Academy of*
24 *Sciences* **112**: 6413–6418.

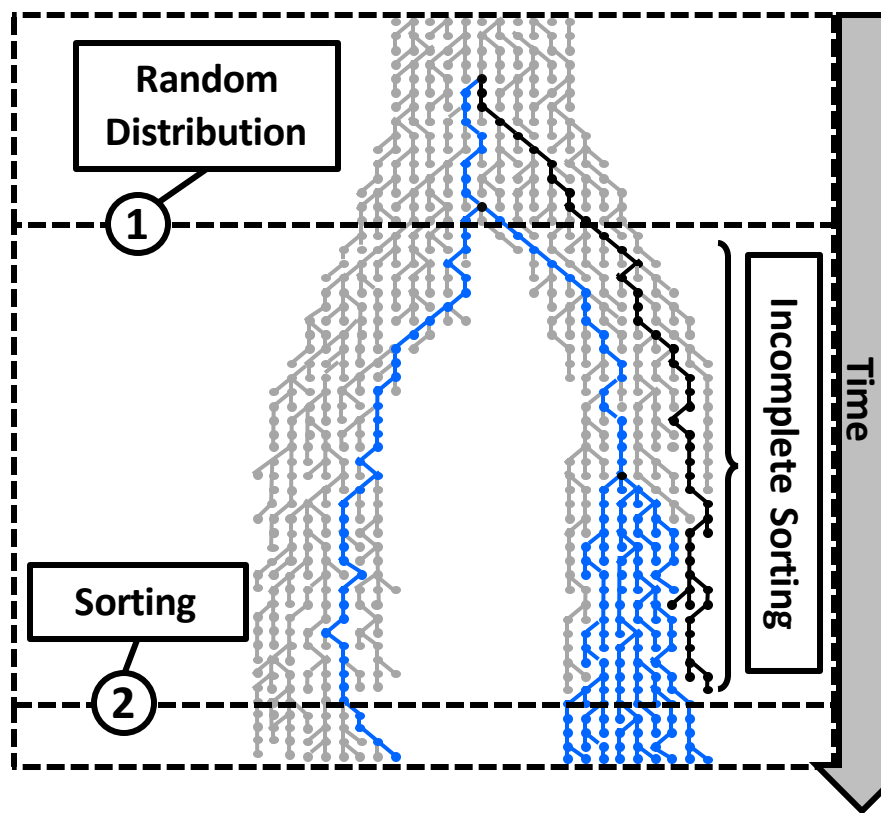
- 1 Nielsen R, Hellmann I, Hubisz M, Bustamante C, Clark AG. 2007. Recent and Ongoing Selection in
2 the Human Genome. *Nature Reviews Genetics* **8**: 857–868.
- 3 Oleksyk TK, Smith MW, O’Brien SJ. 2010. Genome-Wide Scans for Footprints of Natural Selection.
4 *Philosophical transactions of the Royal Society of London. Series B, Biological sciences* **365**:
5 185–205.
- 6 Pääbo S. 2014. The Human Condition—A Molecular Approach. *Cell* **157**: 216–226.
- 7 Powell MJD. 1994. A Direct Search Optimization Method That Models the Objective and Constraint
8 Functions by Linear Interpolation. *Advances in optimization and numerical analysis* **275**: 51–67.
- 9 Prüfer K, Muetzel B, Do HH, Weiss G, Khaitovich P, Rahm E, Pääbo S, Lachmann M, Enard W.
10 2007. FUNC: A Package for Detecting Significant Associations between Gene Sets and
11 Ontological Annotations. *BMC bioinformatics* **8**: 41.
- 12 Prüfer K, Munch K, Hellmann I, Akagi K, Miller JR, Walenz B, Koren S, Sutton G, Kodira C, Winer
13 R et al. 2012. The bonobo genome compared with the chimpanzee and human genomes. *Nature*
14 **486**: 527-531.
- 15 Prüfer K, Racimo F, Patterson N, Jay F, Sankararaman S, Sawyer S, Heinze A, Renaud G, Sudmant
16 PH, de Filippo C, et al. 2014. The Complete Genome Sequence of a Neanderthal from the Altai
17 Mountains. *Nature* **505**: 43–49.
- 18 Przeworski M. 2002. The Signature of Positive Selection at Randomly Chosen Loci. *Genetics* **160**:
19 1179–1189.
- 20 Przeworski M. 2003. Estimating the Time since the Fixation of a Beneficial Allele. *Genetics* **164**:
21 1667–1676.
- 22 Pybus M, Dall’Olio GM, Luisi P, Uzkudun M, Carreño-Torres A, Pavlidis P, Laayouni H, Bertranpetit
23 J, Engelken J. 2014. 1000 Genomes Selection Browser 1.0: A Genome Browser Dedicated to
24 Signatures of Natural Selection in Modern Humans. *Nucleic Acids Research* **42**: D903-D909.

- 1 Pybus OG, Shapiro B. 2009. Natural Selection and Adaptation of Molecular Sequences. The
2 Phylogenetic Handbook. Lemey P, Salemi M, Vamdamme AM. pp 415-417. Cambridge
3 University Press, Cambridge.
- 4 Rabiner LR. 1989. A Tutorial on Hidden Markov Models and Selected Applications in Speech
5 Recognition. *Proceedings of the IEEE* **77**: 257–286.
- 6 Racimo F. 2016. Testing for Ancient Selection Using Cross-Population Allele Frequency
7 Differentiation. *Genetics* **202**: 733–750.
- 8 Racimo F, Kuhlwilm M, Slatkin M. 2014. A Test for Ancient Selective Sweeps and an Application to
9 Candidate Sites in Modern Humans. *Molecular Biology and Evolution* **31**: 3344–3358.
- 10 Rasmussen MD, Hubisz MJ, Gronau I, Siepel A. 2014. Genome-Wide Inference of Ancestral
11 Recombination Graphs. *PLoS Genetics* **10**: e1004342. doi: 10.1371/journal.pgen.1004342.
- 12 Rosenberg NA. 2002. The Probability of Topological Concordance of Gene Trees and Species Trees.
13 *Theoretical Population Biology* **61**: 225–247.
- 14 Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, Cotsapas C, Xie X, Byrne EH, McCarroll SA,
15 Gaudet R et al. 2007. Genome-Wide Detection and Characterization of Positive Selection in
16 Human Populations. *Nature* **449**: 913–918.
- 17 Sabeti PC, Schaffner SF, Fry B, Lohmueller J, Varilly P, Shamovsky O, Palma A, Mikkelsen TS,
18 Altshuler D, Lander ES. 2006. Positive Natural Selection in the Human Lineage. *Science* **312**:
19 1614–1620.
- 20 Sankararaman S, Patterson N, Li H, Pääbo S, Reich D. 2012. The Date of Interbreeding between
21 Neandertals and Modern Humans. *PLoS Genetics* **8**: e1002947. doi:
22 10.1371/journal.pgen.1002947.
- 23 Sankararaman S, Mallick S, Dannemann M, Prüfer K, Kelso J, Pääbo S, Patterson N, Reich D. 2014.
24 The genomic landscape of Neanderthal ancestry in present-day humans. *Nature* **507**: 354–357.

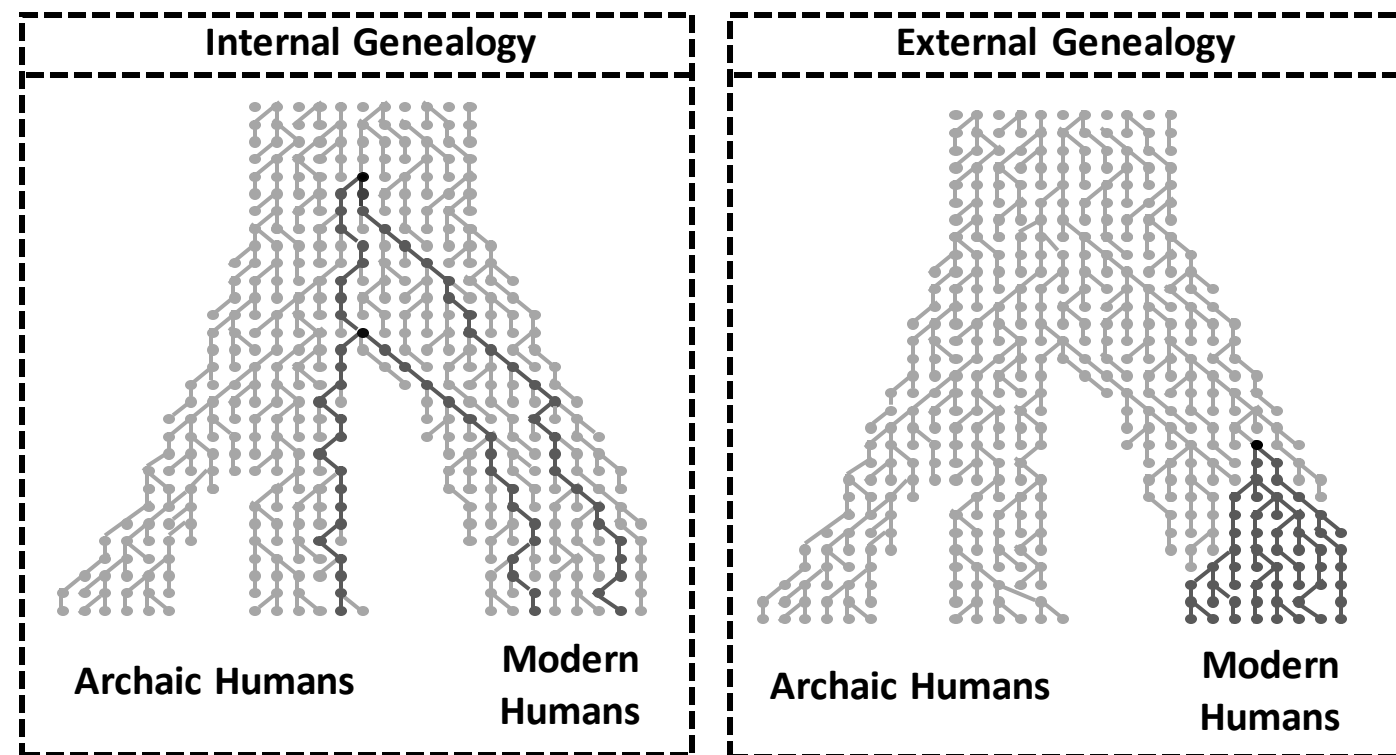
- 1 Sankararaman S, Mallick S, Patterson N, Reich D. 2016. The Combined Landscape of Denisovan and
2 Neanderthal Ancestry in Present-Day Humans. *Current Biology* **26**: 1241–1247.
- 3 Scally A, Dutheil JY, Hillier LW, Jordan GE, Goodhead I, Herrero J, Hobolth A, Lappalainen T,
4 Mailund T, Marques-Bonet T et al. 2012. Insights into hominid evolution from the gorilla
5 genome sequence. *Nature* **483**: 169–175.
- 6 Staab PR, Zhu S, Metzler D, Lunter G. 2014. Scrm: Efficiently Simulating Long Sequences Using the
7 Approximated Coalescent with Recombination. *Bioinformatics* **31**: 1680–1682.
- 8 St Pourcain B, Cents RA, Whitehouse AJ, Haworth CM, Davis OS, O'Reilly PF, Roulstone S, Wren
9 Y, Ang QW, Velders FP, et al. 2014. Common Variation near ROBO2 Is Associated with
10 Expressive Vocabulary in Infancy. *Nature communications* **5**: 4831. doi: 10.1038/ncomms5831.
- 11 Suda S, Iwata K, Shimmura C, Kameno Y, Anitha A, Thanseem I, Nakamura K, Matsuzaki H,
12 Tsuchiya KJ, Sugihara G, et al. 2011. Decreased Expression of Axon-Guidance Receptors in the
13 Anterior Cingulate Cortex in Autism. *Molecular autism* **2**: 14. doi: 10.1186/2040-2392-2-14.
- 14 Tajima F. 1989. Statistical Method for Testing the Neutral Mutation Hypothesis by DNA
15 Polymorphism. *Genetics* **123**: 585–595.
- 16 Varki A, Geschwind DH, Eichler EE. 2008. Explaining Human Uniqueness: Genome Interactions with
17 Environment, Behaviour and Culture. *Nature reviews. Genetics* **9**: 749–763.
- 18 Vernot B, Tucci S, Kelso J, Schraiber JG, Wolf AB, Gittelman RM, Dannemann M, Grote S, McCoy
19 RC, Norton H, et al. 2016. Excavating Neandertal and Denisovan DNA from the Genomes of
20 Melanesian Individuals. *Science* **352**: 235–239.
- 21 Vernot B, Akey JM. 2014. Resurrecting Surviving Neandertal Lineages from Modern Human
22 Genomes. *Science* **343**: 1017–1021.
- 23 Voight BF, Kudaravalli S, Wen X, Pritchard JK. 2006. A Map of Recent Positive Selection in the
24 Human Genome. *PLoS Biology* **4**: 446–458.

- 1 Wang R, Chen CC, Hara E, Rivas MV, Roulhac PL, Howard JT, Chakraborty M, Audet JN, Jarvis
2 ED. 2015. Convergent Differential Regulation of SLIT-ROBO Axon Guidance Genes in the
3 Brains of Vocal Learners. *Journal of Comparative Neurology* **523**: 892–906.
- 4 Weaver TD. 2009. The Meaning of Neandertal Skeletal Morphology. *Proceedings of the National*
5 *Academy of Sciences* **106**: 16028–16033.
- 6 Wright S. 1951. The Genetical Structure of Populations. *Annals of Eugenics* **15**: 322–354.
- 7 Yang MA, Harris K, Slatkin M. 2014. The Projection of a Test Genome onto a Reference Population
8 and Applications to Humans and Archaic Hominins. *Genetics* **198**: 1655–1670.
- 9 Zacher B, Michel M, Schwalb B, Cramer P, Tresch A, Gagneur J. 2016. Accurate Promoter and
10 Enhancer Identification in 127 ENCODE and Roadmap Epigenomics Cell Types and Tissues by
11 GenoSTAN. *PLoS One* **12**: e0169249. doi: 10.1371/journal.pone.0169249.
- 12

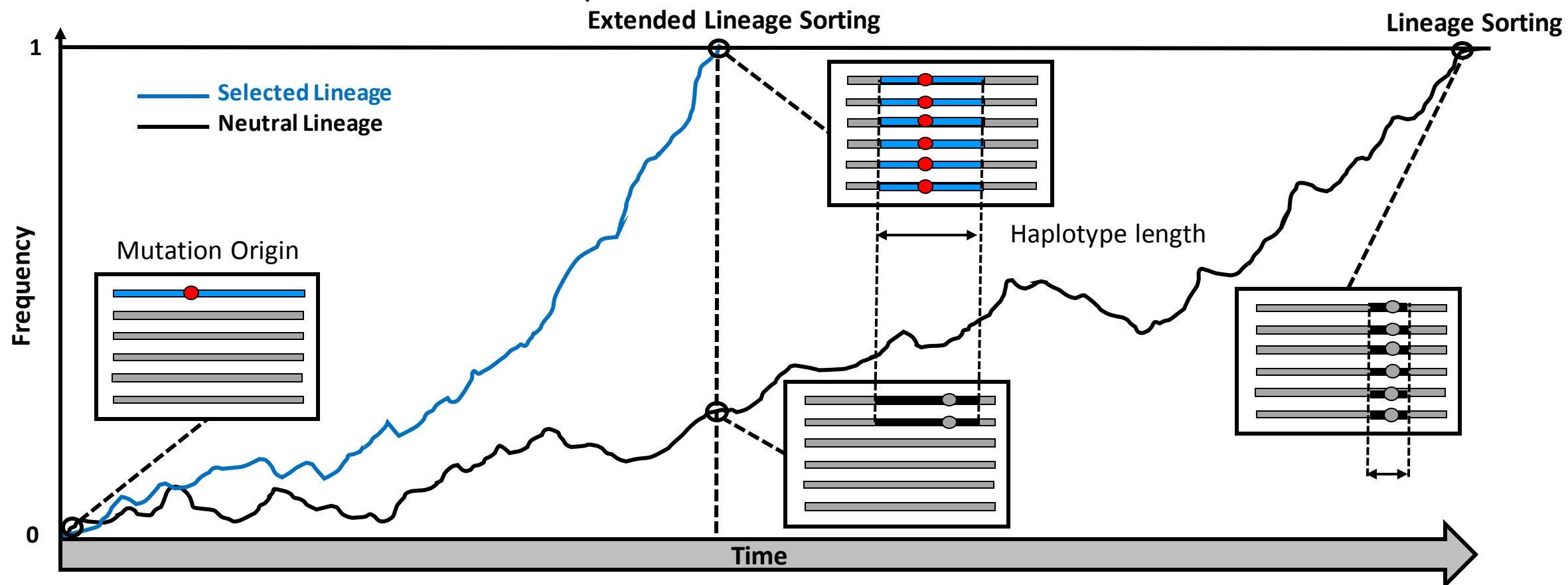
A

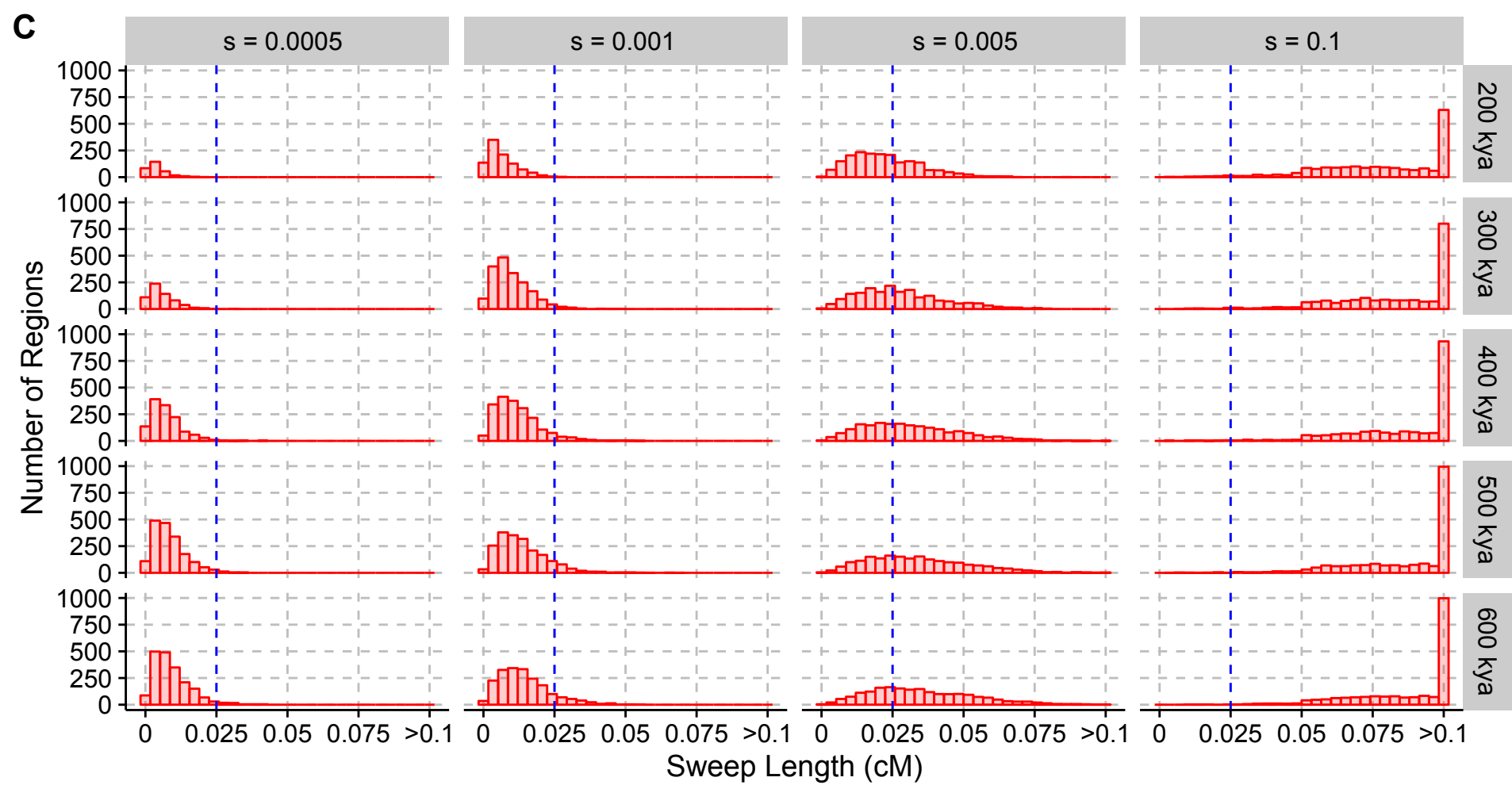
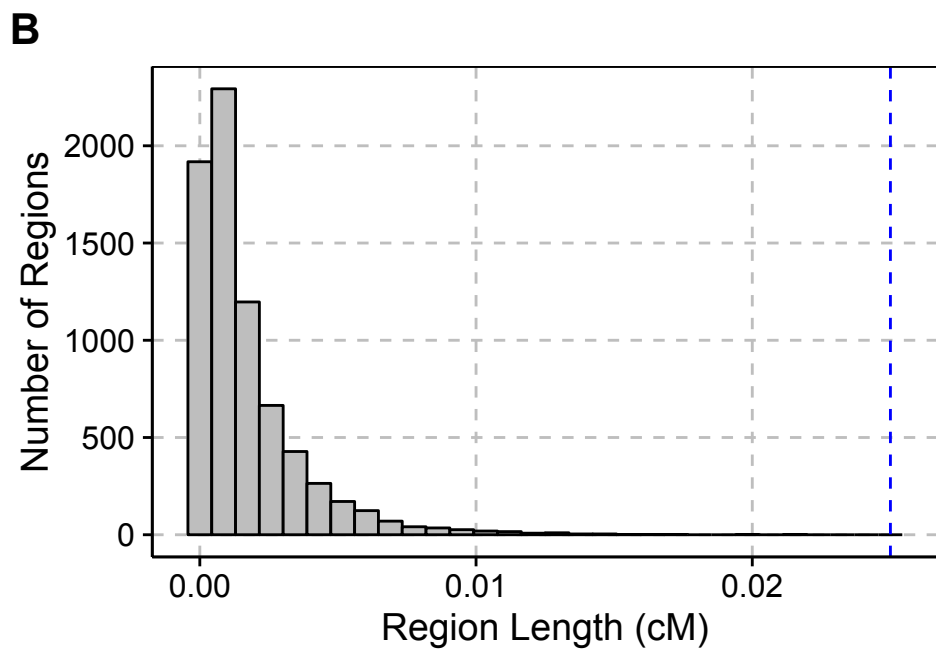
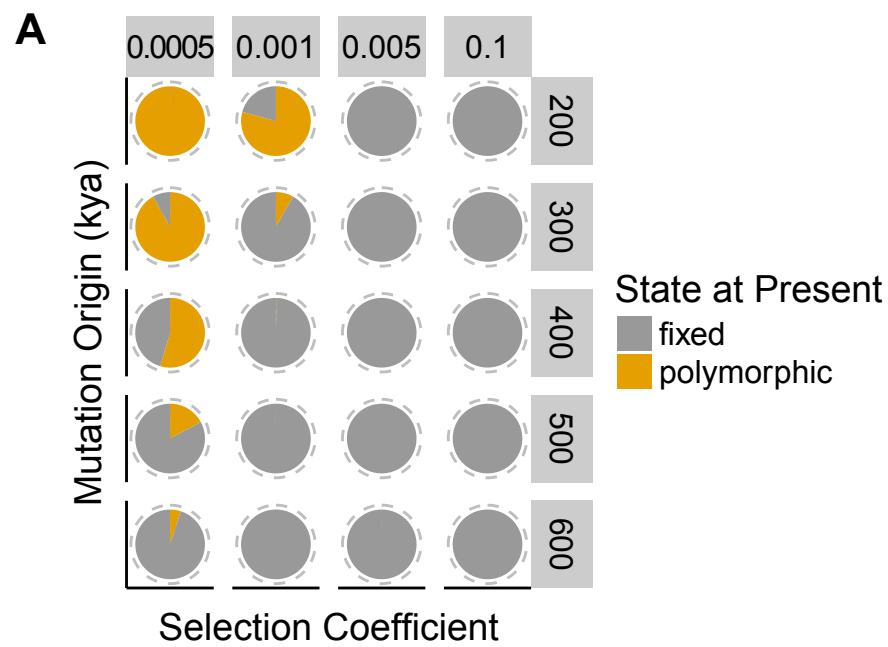


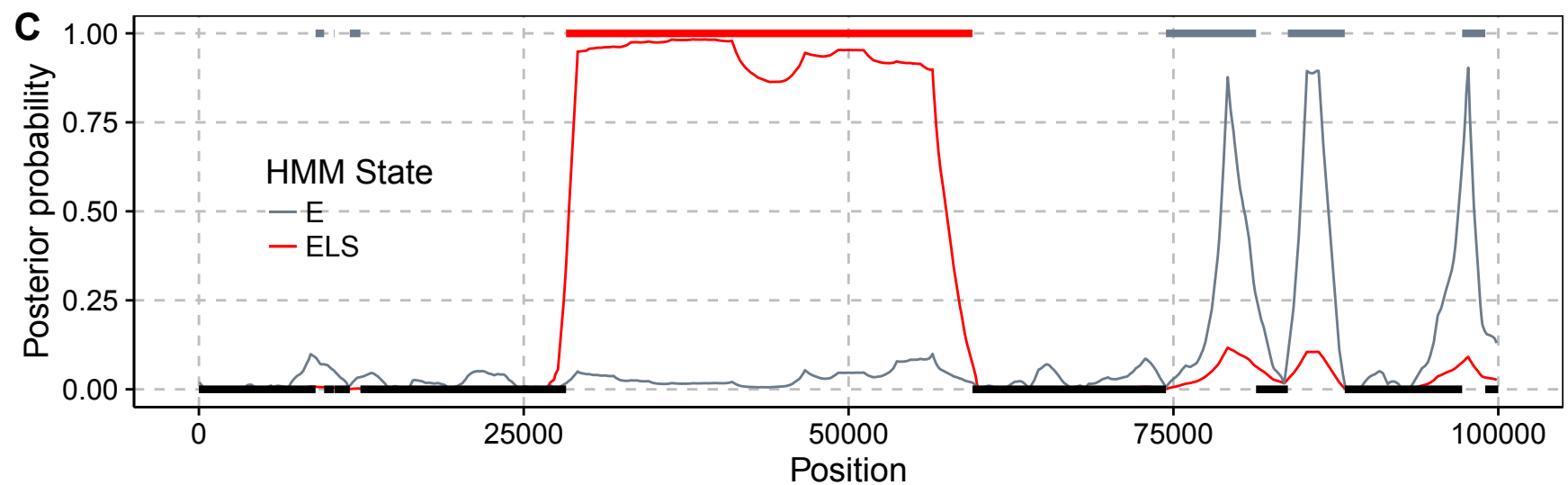
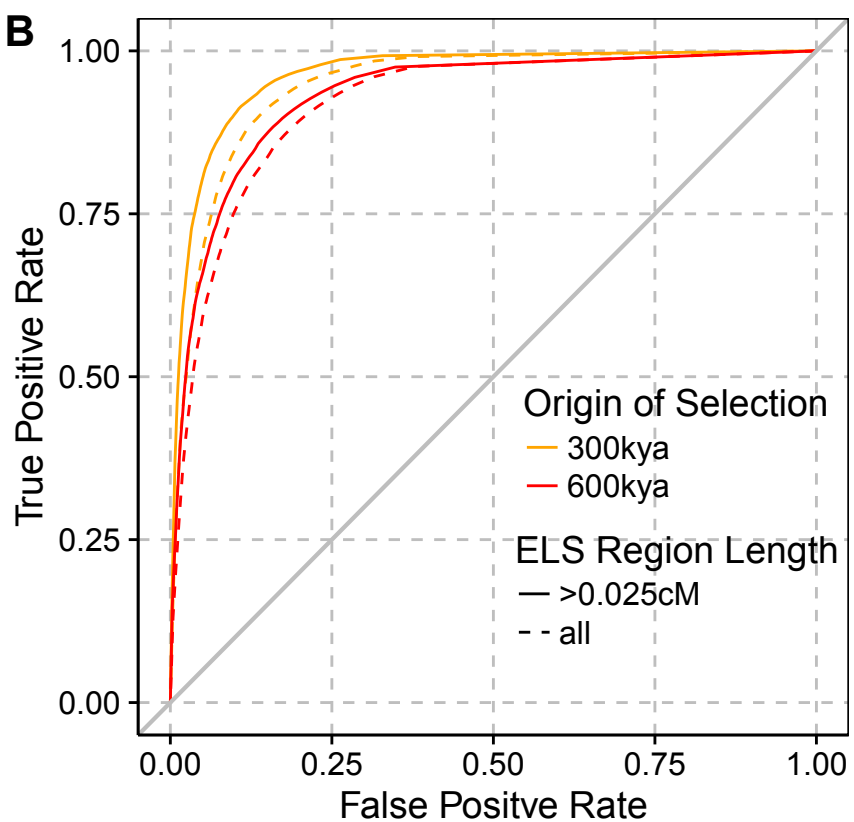
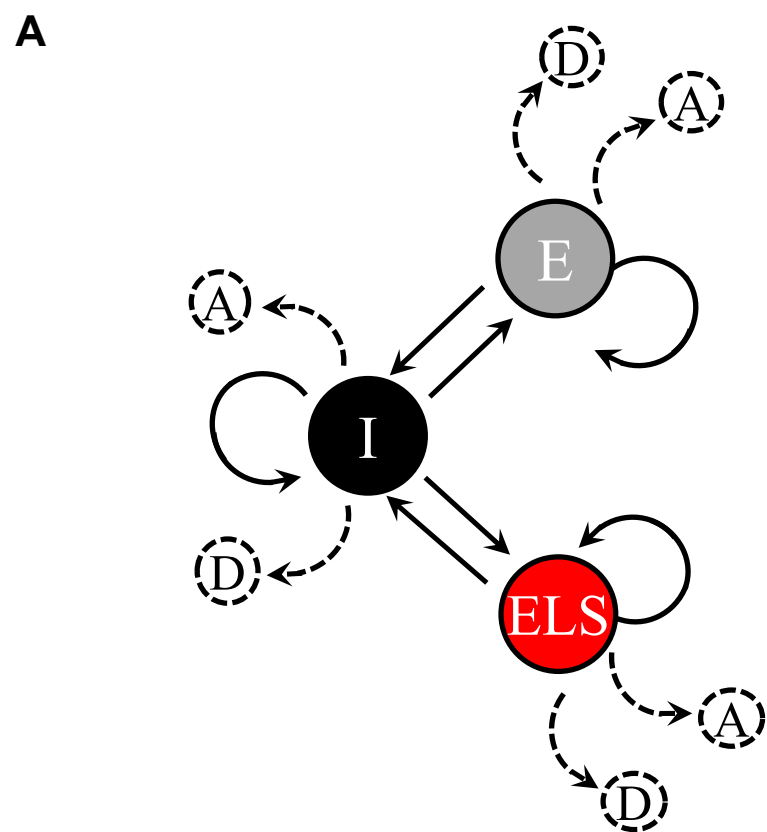
B

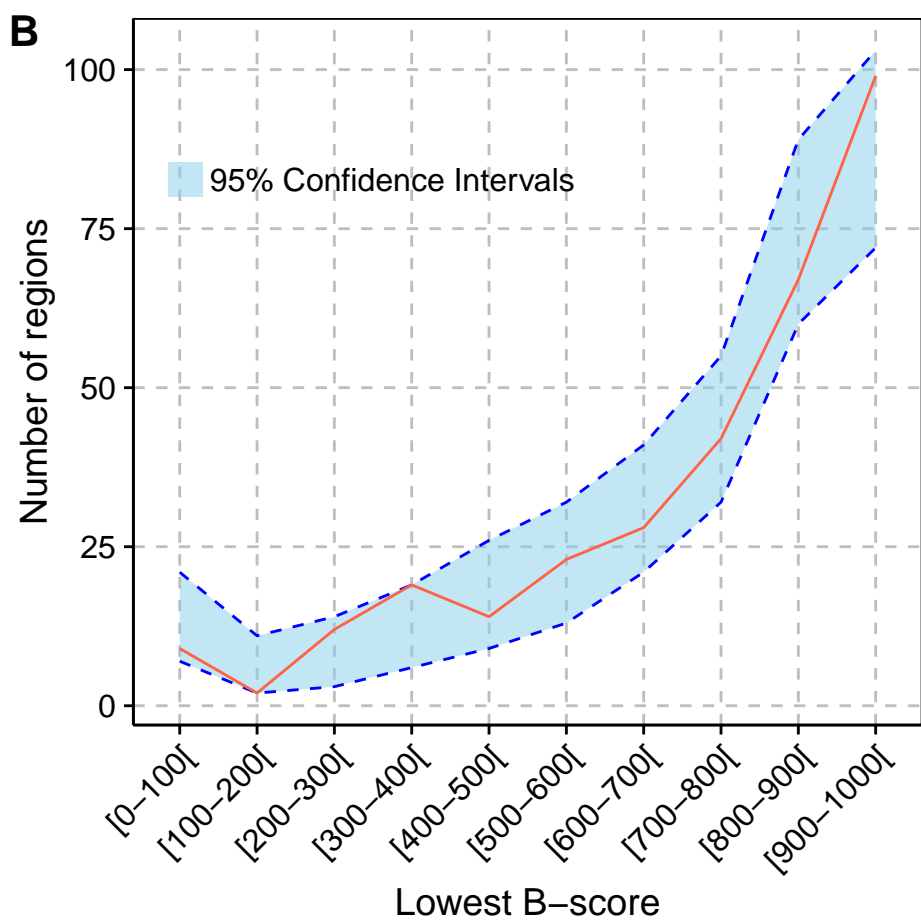
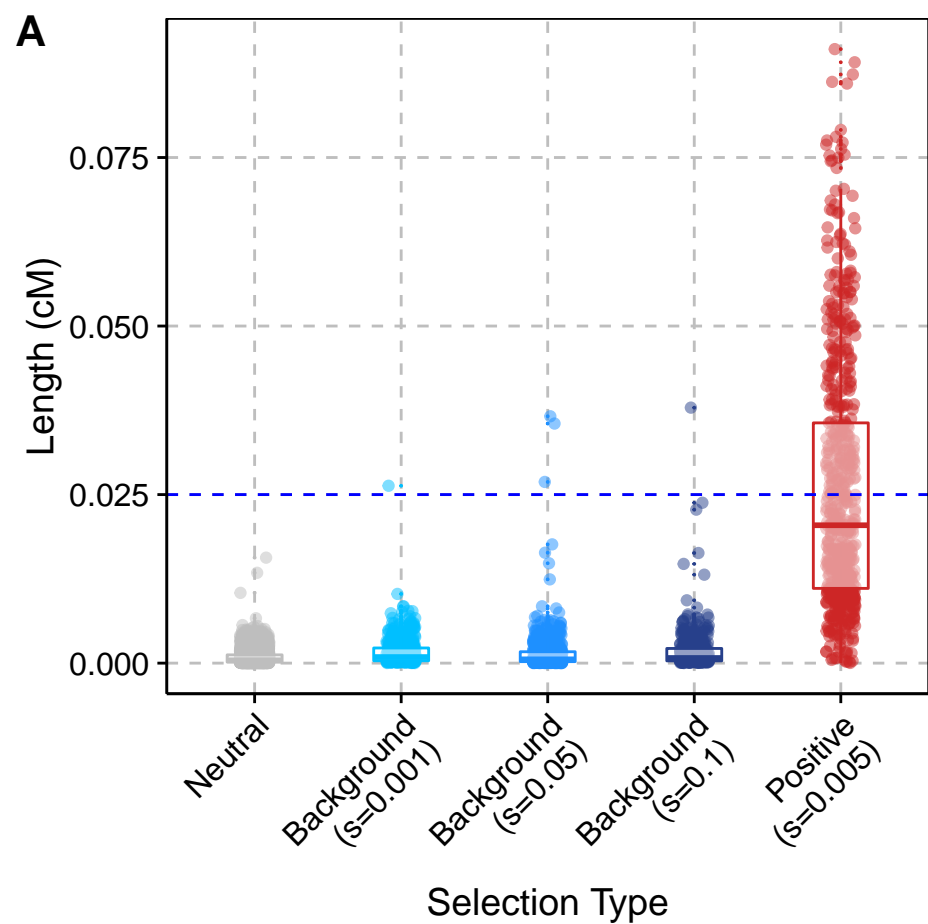


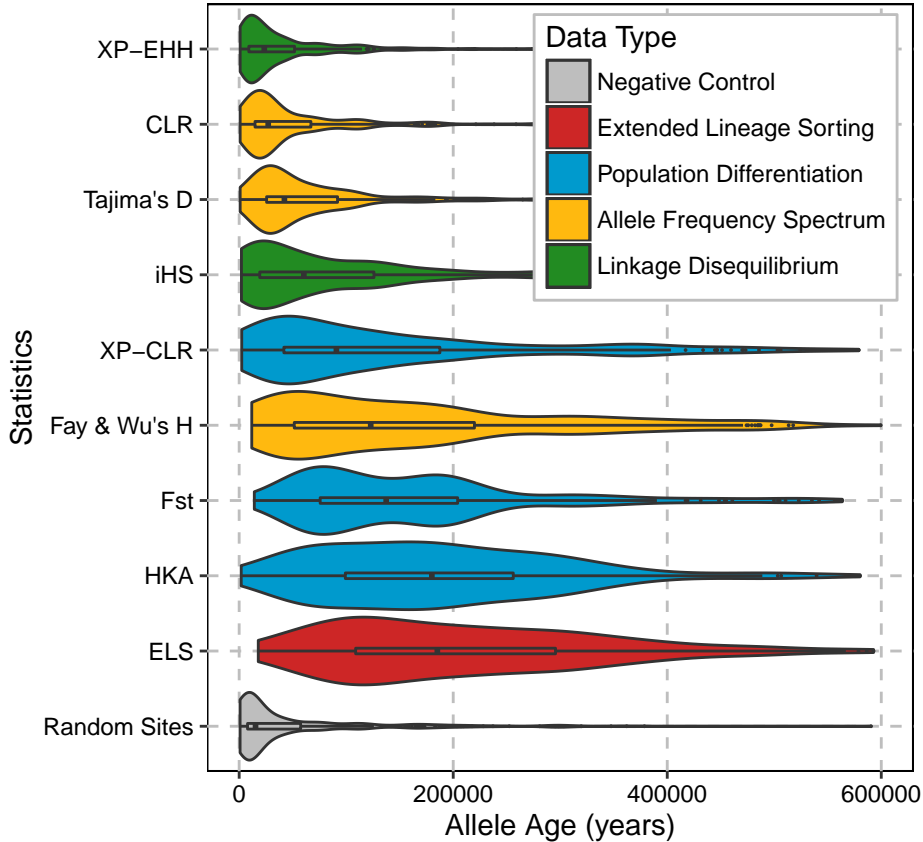
C











Number of Regions Overlapping
at Least One Element

enhancers

genes

promoters

300

200

100

0

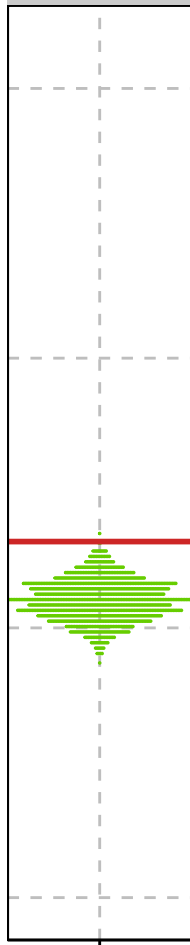
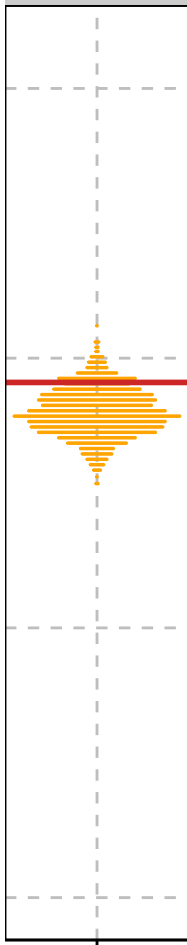
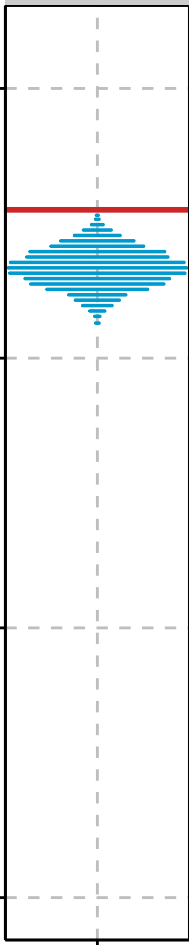


Figure 1: Illustration of the lineage sorting process. (A) Effects on the genealogy. The process starts with a random distribution of lineages when the ancestral population splits. The lineage in black is an outgroup to lineages in blue, so that the blue lineages show a closer relationship between populations than to the black lineage (incomplete lineage sorting). When the blue lineages in the top population reach fixation (through a selective sweep for instance), any lineage from the other populations will constitute an outgroup, thereby completing the sorting of lineages. (B) Two types of genealogies illustrating the possible relationships between an archaic lineage and modern human lineages. (C) Local effects in the genome at different time points. The curves represent the progression of lineage sorting for two independent regions, evolving under neutrality (black curve) and positive selection (blue curve), respectively. Longer fixation times are associated with more recombination so that neutrality produces smaller external regions.

Figure 2: (A) Fraction of selected alleles reaching fixation (grey) or segregating (orange) at present, depending on the strength of selection (columns) and the age of the mutation (rows, in kya) in our simulations. Events for which the selected variant was lost are not shown. (B) Distribution of the genetic length of external regions simulated under neutrality. (C) Distributions of the genetic length of external regions depending on the strength of selection (columns) and age of mutations in kya (rows). The blue line corresponds to the upper limit for the length of external regions produced under neutrality from (B).

Figure 3: (A) Graphical representation of the Extended Lineage Sorting Hidden Markov Model. States are depicted by nodes and transitions by edges. Each state emits an archaic allele as either derived, D, or ancestral, A, depending on the type of site in the modern human population (fixed or segregating at a given frequency). States are labelled I for Internal, E for External and ELS for Extended Lineage Sorting. (B) Receiver Operator Curves for varying cutoffs on the posterior probability of the ELS state and counting the number of sites in ELS regions that were correctly labeled.

All bases labelled ELS outside of simulated ELS regions are considered false positives. Sites in ELS regions with a posterior probability below the cutoff are considered false negatives. (C) Example of the labelling of a simulated ELS region. Horizontal bars indicate true external (top) and internal (bottom) regions. The posterior probability is shown in red for ELS regions and in grey for E regions. The region overlapping position 50,000 (red bar) is caused by a simulated selective sweep.

Figure 4: Effects of background selection. (A) Comparison of the length of ELS regions in simulations of different scenarios. For the distribution under background selection, the s parameter corresponds to the average selection coefficient from the gamma distribution (shape parameter of 0.2). We assumed that the deleterious mutations are recessive with dominance coefficient $h=0.1$. The horizontal blue line corresponds to the length cutoff applied to the real data. (B) Distribution of B-scores in the candidate sweep regions (red curve) compared to sets of random regions with matching physical lengths (blue area with dotted blue lines indicating the 95% confidence intervals over 1000 random sets of regions). The lowest B-score (i.e. stronger background selection) was chosen when a region overlapped several B-score annotations.

Figure 5: Distributions of estimated ages of the modern human segregating derived variants with the highest frequency in putatively selected regions or the age of the derived variants at sites identified by various genome-wide scans. Our candidate regions are labelled as ELS, for Extended Lineage Sorting, other candidate regions are from (Cagan et al. 2016; Pybus et al. 2014). The color coding indicates the type of signal detected by each method. Ages were estimated by ARGweaver (Rasmussen et al. 2014). We only report events between 0 and 600kya.

Figure 6: Enrichment for regulatory elements (enhancers, P -value <0.001 , protein-coding genes, P -value $=0.124$, and promoters, P -value $=0.002$) in the extended set of 314 candidate sweep regions. The distributions were obtained by randomly placing candidate regions in the genome to obtain lists of regions with similar physical length. The red lines represent the value observed in the real extended set.



Detecting ancient positive selection in humans using extended lineage sorting

Séphane Peyégne, Michael James Boyle, Michael Dannemann, et al.

Genome Res. published online July 18, 2017

Access the most recent version at doi:[10.1101/gr.219493.116](https://doi.org/10.1101/gr.219493.116)

P<P	Published online July 18, 2017 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution 4.0 International license), as described at http://creativecommons.org/licenses/by/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<http://genome.cshlp.org/subscriptions>
