

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

3 Stéphane Peyrégne*, 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 in the brain and present a list of genomic regions that are predicted to underlie positively
20 selected human 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).

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

1 signal of selection also fades over time and events of positive selection older than 300kya become
2 undetectable.

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

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

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

1 RESULTS

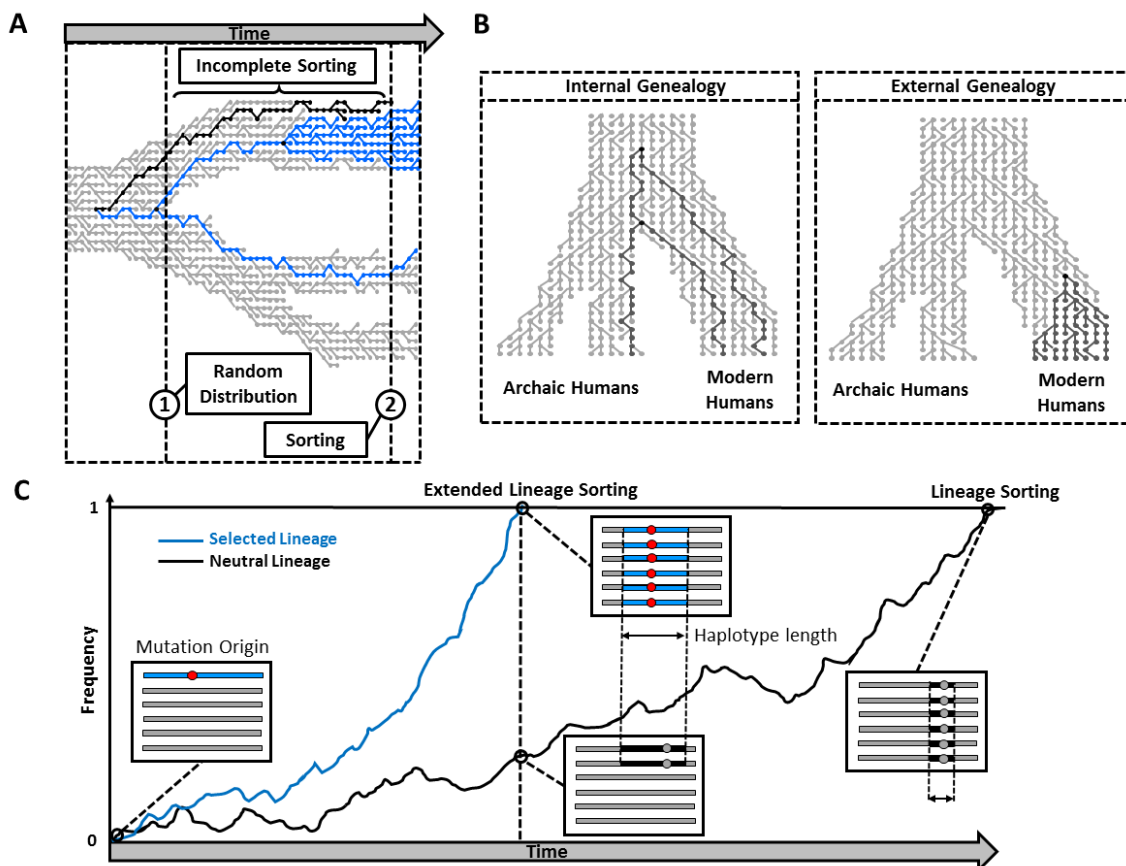
2 Selection causes extended lineage sorting between closely related populations

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

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



12

1 Expected Incomplete Lineage Sorting among Humans to Archaics

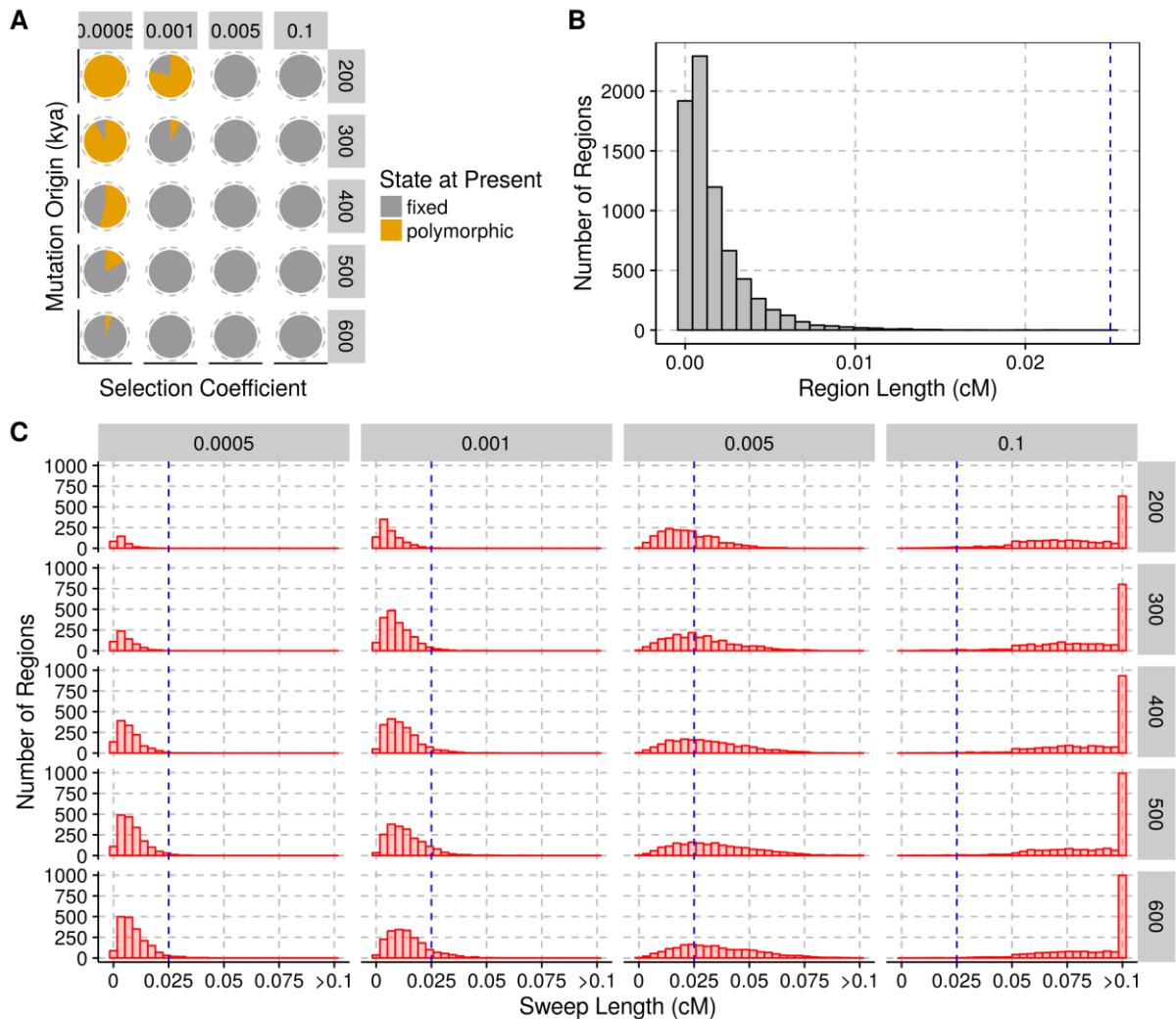
2 We used coalescent simulations to determine the incidence and expected length of regions resulting
3 from incomplete lineage sorting in modern humans. Using a model of human demographic history
4 (Yang et al. 2014), we estimated the fraction of lineage sorting in modern humans in regards to
5 Neandertals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

10 Minimum Strength of Selection to Produce Detectable Sweep Signals

11 We investigated the range of selection coefficients that could have led to the fixation of a lineage after
12 the split with the Archaic hominins, but before the differentiation of genetically modern humans about
13 100–120 kyr ago (Li and Durbin 2011) by simulating mutations occurring at different times and
14 evolving with different selection coefficients. While the simulations show that completed selective
15 sweeps could have occurred with selection coefficients as low as 0.0005 (Figure 2A), the length
16 distribution of haplotypes reaching fixation is indistinguishable from neutrality for selection
17 coefficients under 0.001 (Figure 2, B and C). Under neutrality, the average length of external regions
18 was 0.02 cM and remained below 0.03cM for most simulations with a selection coefficient of 0.001.
19 In contrast, external regions longer than 0.1cM were observed for selection coefficients above 0.05.
20 Therefore, detectable signals are expected to be biased towards strong events with a selection
21 coefficient larger than 0.001.

22 **Figure 2:** (A) Fraction of selected alleles reaching fixation (grey) or segregating (orange) at present,
23 depending on the strength of selection (columns) and the age of the mutation (rows, in kya) in our
24 simulations. Events for which the selected variant was lost are not shown. (B) Distribution of the
25 genetic length of external regions simulated under neutrality. (C) Distributions of the genetic length of

- 1 external regions depending on the strength of selection (columns) and age of mutations in kya (rows).
- 2 The blue line corresponds to the upper limit for the length of external regions produced under
- 3 neutrality from (B).



4

5 Hidden Markov Model to Detect Extended Lineage Sorting

- 6 To detect regions of Extended Lineage Sorting, we modeled the changes of local genealogies along the
- 7 genome with a hidden Markov model. We distinguish two types of genealogies, internal or external,
- 8 depending on whether the archaic lineage falls inside or outside of the human variation respectively
- 9 (Figure 3A). The model includes a third state corresponding to extended lineage sorting, and external
- 10 regions produced by this state are required to be longer, on average, than those produced by the
- 11 external state. The three states are inferred from the state of the archaic allele (ancestral or derived)

1 either at a polymorphic position in modern humans or at a position where modern humans carry a
2 fixed derived variant. In the following, we describe the different statistical properties expected for
3 each type of genealogy.

4 We first consider external regions. At modern human polymorphic sites, the archaic genome is
5 expected to carry the ancestral variant since the derived variant would indicate incomplete lineage
6 sorting. To account for sequencing errors or misassignment of the ancestral state, we allow a
7 probability of 0.01 for carrying the derived allele (see Material and Methods). At sites where the
8 derived allele is fixed, the archaic genome could carry either the derived or ancestral state depending
9 on whether the fixation event occurred before or after the split of the archaic from the modern human
10 lineage.

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

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

1 Accuracy of Parameter Estimates and Inferred Genealogies

2 We first investigated the performance of the parameter inference on simulated data under neutral
3 evolution. We found that the estimated probabilities for encountering ancestral/derived alleles in
4 external and internal regions fit the simulated parameters well (on average less than ± 0.08 from
5 simulated under all tested conditions) (Supplemental Figures S1 and S2), while the estimated length of
6 internal and external regions deviate more from the simulated lengths (around 15% overestimate of the
7 mean length, Supplemental Figure S3). However, we found that the model exhibits better accuracy in
8 labelling the correct genealogies with the estimated length parameters compared to the simulated true
9 values (Supplemental Figure S4). This difference seems to originate from the difficulty in accurately
10 detecting very short external regions or internal regions with very few informative sites. We note that
11 detecting selection is not affected by this problem since we are primarily interested in detecting long
12 external regions.

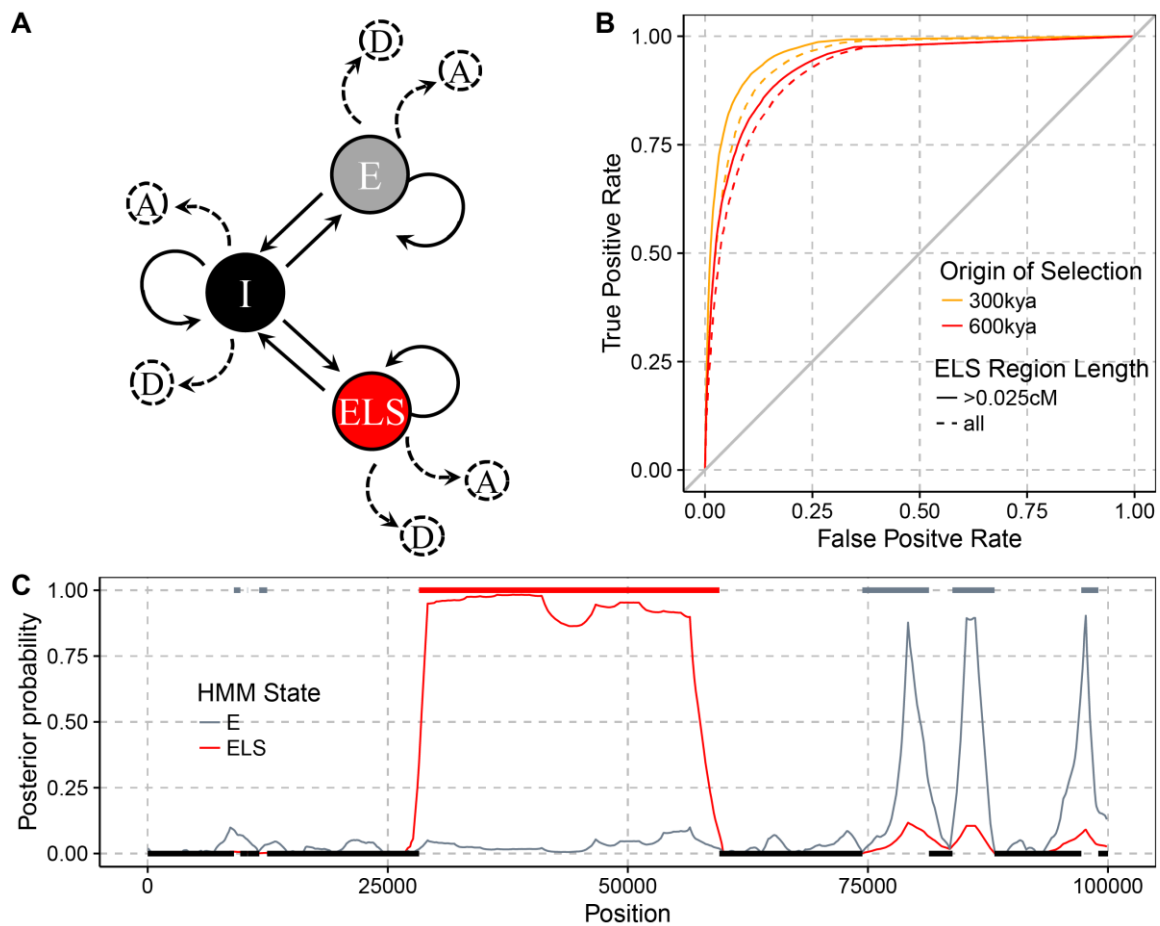
13 We do not expect ELS regions to be detected in our neutral simulations and indeed we found that
14 either the estimated proportion of ELS converged to zero or the maximum likelihood estimate for the
15 length of ELS and external regions converge to the same value (49% and 51% of simulations
16 respectively). A likelihood ratio test comparing a model without the ELS state to the full model with
17 the ELS state also showed no significant improvement with the additional state in almost all neutral
18 simulations (only one likelihood ratio test out of 100 simulations showed a significant improvement
19 after Bonferroni correction for multiple testing).

20 We then evaluated the accuracy of the ELS method to assign the correct genealogy to regions based on
21 sequences obtained through coalescent simulations with selection (Figure 3, B and C). In these
22 simulations, the underlying genealogy at each site along the sequences is known and can be compared
23 to the estimates. To be conservative, we only focus on results with the smallest selection coefficient
24 ($s=0.005$) that produces regions long enough to be detectable. In Figure 3B we show the accuracy for
25 labelling the extended lineage sorting regions dependent on the posterior probability cutoff for the

1 ELS state. The results demonstrate that the model has sufficient power to accurately label sites that
2 experienced selection with a coefficient $s \geq 0.005$ and an occurrence of the beneficial mutation as long
3 as 600,000 years ago.

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

14 **Figure 3:** (A) Graphical representation of the Extended Lineage Sorting Hidden Markov
15 Model. States are depicted by nodes and transitions by edges. Each state emits an archaic
16 allele as either derived, D, or ancestral, A, depending on the type of site in the modern human
17 population (fixed or segregating at a given frequency). States are labelled I for Internal, E for
18 External and ELS for Extended Lineage Sorting. (B) Receiver Operator Curves for varying
19 cutoffs on the posterior probability of the ELS state and counting the number of sites in ELS regions
20 that were correctly labeled. All bases labelled ELS outside of simulated ELS regions are considered
21 false positives. Sites in ELS regions with a posterior probability below the cutoff are considered false
22 negatives. (C) Example of the labelling of a simulated ELS region. Horizontal bars indicate true
23 external (top) and internal (bottom) regions. The posterior probability is shown in red for ELS regions
24 and in grey for E regions. The region overlapping position 50,000 (red bar) is caused by a simulated
25 selective sweep.



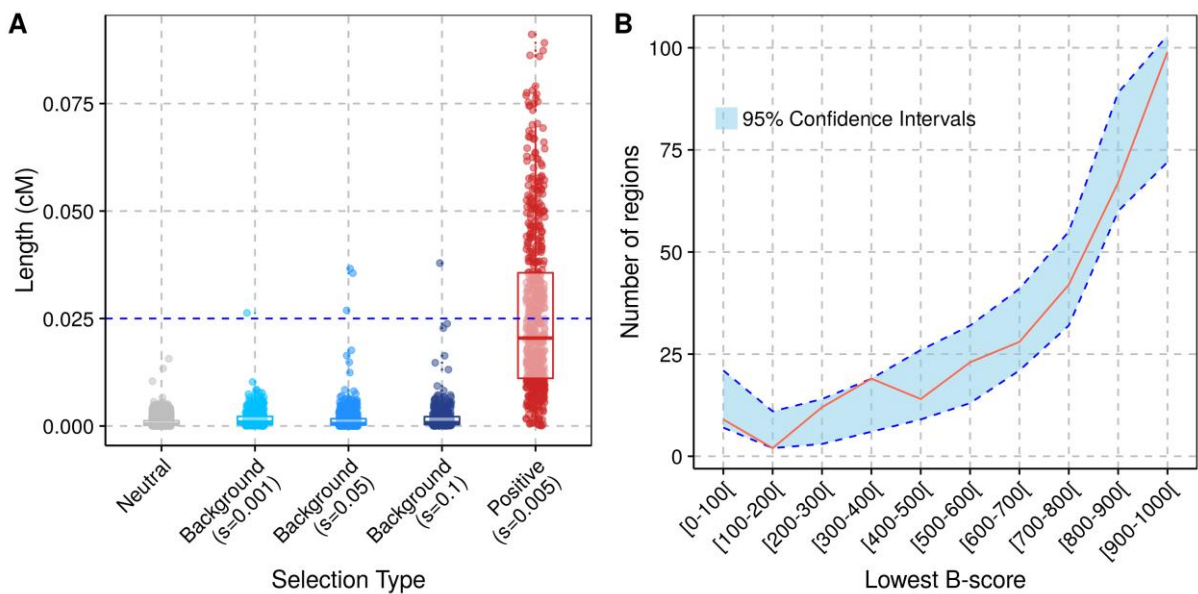
1

2 Role of Background selection

3 Background selection is defined as the constant removal of neutral alleles due to linked deleterious
4 mutations (Charlesworth et al. 1993). In regions of the genome that undergo background selection, a
5 fraction of the population will not contribute to subsequent generations, causing a reduced effective
6 population size. As a consequence, remaining neutral alleles can reach fixation faster than under
7 neutrality, potentially producing unusually long external regions that could be mistaken as signals of
8 positive selection. We investigated the effects of background selection by running forward simulations
9 with parameters that mimic the strength and extent of background selection estimated for the human
10 genome (Messer 2013). While background selection simulations did produce some long outlier
11 regions that fall outside the distribution observed in neutral simulations, most regions are still smaller
12 than regions simulated with positive selection at a conservative selection coefficient of 0.005 (Figure
13 4A). Indeed, among the 1160 external regions detected in our simulations of background selection

1 ($s=0.05$, Figure 4A) only six were labeled as ELS and only three passed the minimal length filter of
2 0.025 cM.

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



13

14 Candidate Regions of Positive Selection on the Human Lineage

15 To identify ancient events of positive selection on the human lineage, we applied the ELS method to
16 African genomes from the 1000 genomes project (Abecasis et al. 2012). We disregarded non-African

1 populations since Neandertal introgression in these populations could mask selective sweeps and lead
2 to false negatives. A model with ELS fits the data significantly better than a model without the ELS
3 state for all chromosomes and for both tested recombination maps (p -value $< 1e-8$, Supplemental
4 Table S2).

5 We identified 81 regions of human extended lineage sorting for which both recombination maps
6 support a genetic length greater than 0.025cM (average length: 0.05 cM). Depending on the
7 recombination map, the longest overlap between the maps is 0.12 (African-American map) or 0.17
8 (deCode map) cM long, which is three to four times longer than the longest regions produced under
9 background selection in our simulations. An additional 233 regions are longer than 0.025cM according
10 to only one recombination map, with 71% of those additional regions showing support for the ELS
11 state using both recombination maps. This suggests that the variation in the candidate set mostly stems
12 from uncertainty about recombination rates. We will refer to the set of 81 regions as the core set
13 (Supplemental File S1) and the set including the 233 putatively selected regions found with just one
14 recombination map as the extended set (314 regions, Supplemental File S2).

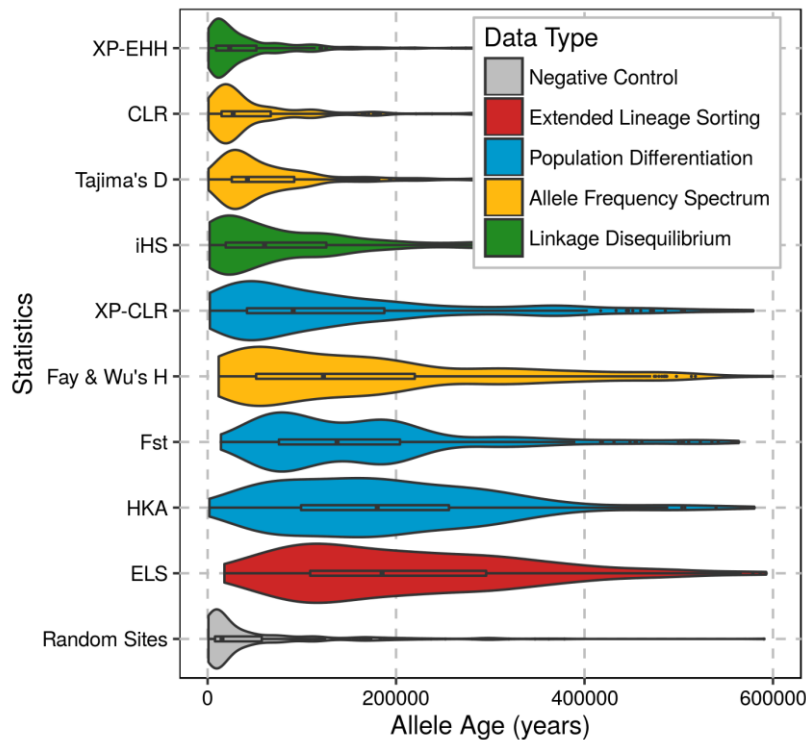
15 For completeness, we also ran our model on the X chromosome and identified 12 additional
16 candidates (43 if we consider candidates found with at least one recombination map), applying a more
17 stringent length cutoff of 0.035 cM to account for the stronger effects of random drift on this
18 chromosome (cf. Material and Methods). Interestingly, we also found a significant increase of
19 posterior probabilities for selection within previously reported regions under potential recurrent
20 selective sweeps in apes (Dutheil et al. 2015; Nam et al. 2015) (Mann-Whitney one-sided test, p -value
21 $< 2.2e-16$, Supplemental Table S3).

22 The detected selection candidate regions on the autosomes do not show a decrease in B scores
23 (McVicker et al. 2009), a local measure of background selection strength, compared with random
24 regions (Figure 4B; Wilcoxon rank sum test comparing the average B-scores with permuted regions,

1 p-value=0.565, or comparing the lowest B-scores in our regions to permuted regions, p-value=0.504).
2 This suggests that candidate regions are not primarily generated by strong background selection.

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

13 **Figure 5:** Distributions of estimated ages of the modern human segregating derived variants
14 with the highest frequency in putatively selected regions or the age of the derived variants at
15 sites identified by various genome-wide scans. Our candidate regions are labelled as ELS, for
16 Extended Lineage Sorting, other candidate regions are from (Cagan et al. 2016; Pybus et al.
17 2014). The color coding indicates the type of signal detected by each method. Ages were
18 estimated by ARGweaver (Rasmussen et al. 2014). We only report events between 0 and
19 600kya.



1

2 Overlap with Genes, Enhancers and Promoters

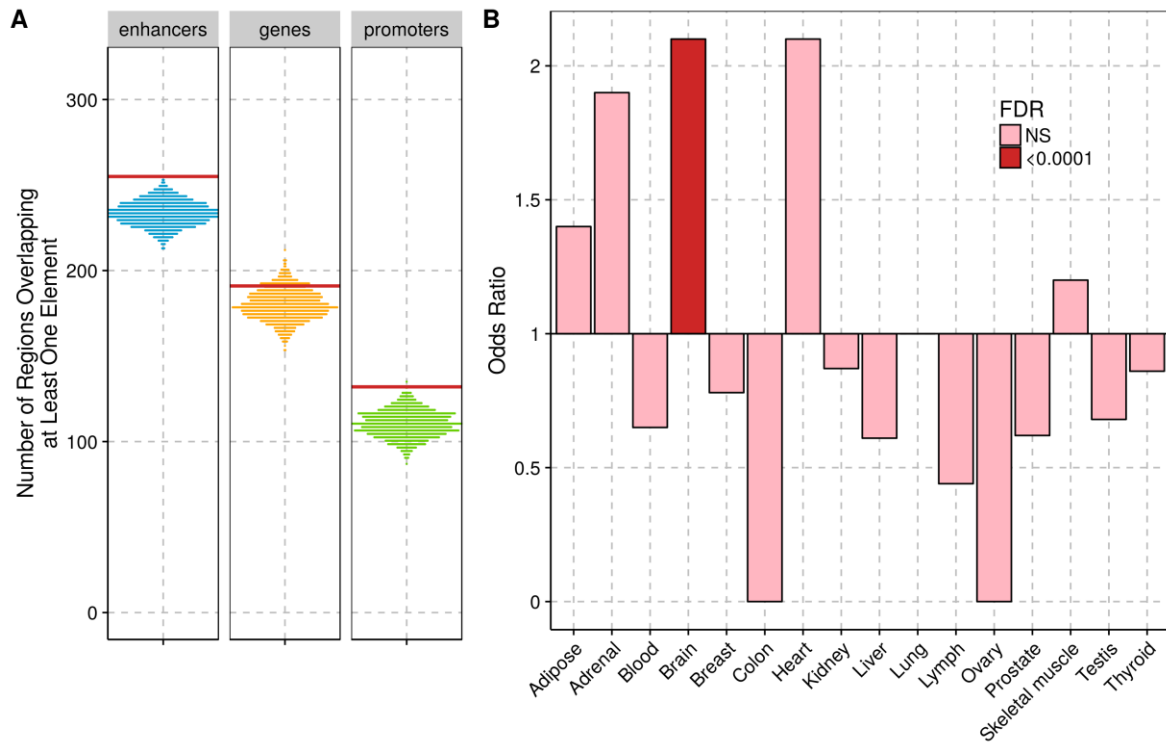
3 Since positive selection acts on advantageous phenotypes that are caused by changes to functional
4 elements in the genome, we would expect that our candidate regions overlap functional elements in the
5 genome more often than expected.

6 We first tested this hypothesis by counting the overlap between sweep candidate regions and protein
7 coding genes (ENSEMBL release 82). We find no statistically significant overlap of ELS regions with
8 protein coding genes compared to randomly placed regions of the same size (p-value = 0.671 and
9 0.124, for core and extended set, respectively; Figure 6A). Previous work has identified 96 proteins
10 that carry human fixed derived non-synonymous changes compared to Neandertal and Denisova,
11 which constitute a particularly interesting subset of potentially functional changes to genes that may
12 have been caused by selective sweeps (Prüfer et al. 2014). We found no overlap between these genes
13 and the core set of sweep candidate regions that were identified by both recombination maps.
14 However, when considering the extended set of sweep candidate regions, 11 regions overlapped such
15 genes: *ADSL*, *BBIP1*, *ENTHD1*, *HERC5*, *KATNA1*, *KIF18A*, *NCOA6*, *PRDM10*, *SCAP*, *SLITRK1* and

1 *ZNHIT2*. This overlap is significantly larger than expected by chance (only 2 genes are expected on
2 average; p -value $< 10^{-3}$). In all instances, the candidate regions contained at least one fixed amino acid
3 change. Since fixed changes are part of the information used to infer external regions, it stands to
4 reason that the presence of such a change may bias towards observing an overlap with candidate
5 regions. However, we note that the overlap with fixed amino acid changes is also significantly larger
6 than the overlap with other fixed changes (963 of 20347 fixed changes fall within candidate regions
7 from the extended set; binomial p -value=0.006).

8 Phenotype may also be influenced by regulatory changes that affect gene expressions. Interestingly,
9 we found a significant enrichment for regions overlapping enhancers and promoters (p -value <0.001
10 and p -value=0.002, respectively; see Figure 6A) when considering the extended set of 314 candidate
11 regions. However, this enrichment was not significant for the smaller core set of candidates.

12 **Figure 6:** Gene expression enrichment (A) Enrichment for regulatory elements (enhancers, p -
13 value <0.001 , protein-coding genes, p -value=0.124, and promoters, p -value=0.002) in the extended set
14 of 314 candidate sweep regions. The distributions were obtained by randomly placing candidate
15 regions in the genome to obtain lists of regions with similar physical length. The red lines represent
16 the value observed in the real extended set. (B) Enrichment for genes with tissue-specific expression
17 (significantly highest expression in one tissue compared to others). NS stands for Not Significant at a
18 FDR (False Discovery Rate) cutoff of 0.05.



1

2 Gene Ontology and Tissue-Specific Expression

3 To further investigate the biological function of our regions, we tested for gene ontology enrichment in
 4 genes within the regions defined by the extended set and find enrichment in a number of categories
 5 (Table 1). Two categories were related to amylase-activity. However, upon further investigation, we
 6 found that this enrichment was driven by one region in a cluster of amylase genes on chromosome 1.
 7 Such a clustering for genes of similar function did not affect the other categories, which showed an
 8 enrichment for genes related to membrane and synapse, suggesting that our regions may have been
 9 selected for brain-related phenotypes.

10 **Table 1:** Gene ontology enrichment for the extended set of candidate genes.

Root	Node	ID	FWER
Molecular Function	alpha-amylase activity	GO:0004556	<0.001
Cellular Component	membrane region	GO:0098589	0.004
Cellular Component	plasma membrane region	GO:0098590	0.004
Cellular Component	synapse	GO:0045202	0.004
Cellular Component	synapse part	GO:0044456	0.005
Molecular Function	amylase activity	GO:0016160	0.005

1 To further test this enrichment, we assigned genes that overlap our extended dataset to tissues for
2 which they show the significantly highest expression (see Material and Methods). In agreement with
3 the Gene Ontology enrichment, we find a marginally significant signal for genes expressed highest in
4 brain (OR=1.60, p-value=0.020, FDR=0.16). In an attempt to include potential regulatory changes in
5 the enrichment test, we further assigned genes that may have been affected by a candidate region
6 upstream or downstream (see Material and Methods). We then found a significant enrichment for
7 genes expressed in the brain (core set: OR=3.10, p-value<10⁻⁴, FDR<10⁻³; extended set: OR=2.10, p-
8 value<10⁻⁶, FDR<10⁻⁴, Figure 6B, Supplemental Table S6 and S7) as well as a marginal signal for the
9 heart in the extended dataset (OR=2.10, p-value=0.01, FDR=0.083, Supplemental Table S6).

10 We then investigated whether the enrichment for genes expressed in the brain was specific to any
11 brain tissue across different developmental stages using the Allan Brain Atlas (Hawrylycz et al. 2012;
12 Miller et al. 2014). We found expression enrichment in several tissues with the strongest signals in the
13 cerebral cortex in fetus (hypergeometric test, FWER=0.029), the hippocampus in children
14 (FWER=0.004; marginal in infant, FWER=0.07) as well as the pineal gland in adults (FWER=0.017,
15 Methods).

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
26 of out-of-Africa ancestral alleles at African fixed derived sites (~26%; bootstrap p-value=0.583). We

1 also find no significant reduction in frequency of Neandertal ancestry around candidate substitutions
2 in sweep regions, when comparing one randomly sampled fixed African substitution per region against
3 random regions matched for size and distance to genes (Supplemental Figure S5 and S6).

4 If selection against the re-introduction of an ancestral variant were very strong, selection may have
5 depleted Neandertal ancestry in a large region surrounding the selected allele. Interestingly we find
6 some of our sweep candidate regions that fall within the longest deserts of both Neandertal and
7 Denisova ancestry (Table 2) (Vernot et al. 2016). A significantly high number of the core set of
8 regions fall in these deserts (5/81 regions, p-value=0.024), while the extended set shows no significant
9 enrichment (9/314 regions, p-value=0.205).

10 **Table 2:** Genes from the core set of candidate regions overlapping with long deserts of Neandertal and
11 Denisovan ancestry.

Chromosome	Start	End	Overlapping Genes	Overlapping Regulatory Domains
chr1	104000000	104154236	<i>AMY2B, RNPC3</i>	<i>COL11A1</i>
chr1	113429666	113560554	<i>SLC16A1</i>	<i>FAM19A3, LRIG2</i>
chr3	77027850	77033270	<i>ROBO2</i>	-
chr7	122320038	122379695	<i>RNF133, RNF148, CADPS2</i>	<i>TAS2R16</i>
chr10	107809941	107866217	-	<i>SORCS1, SORCS3</i>

12 DISCUSSION

13 Many genetic changes set modern humans apart from Neandertals and Denisovans but their functions
14 remain elusive. Most of these changes probably resulted in either no change to the phenotype or to a
15 selectively neutral change. However, in rare instances selection may have favored changes modifying
16 the appearance, behavior and abilities of present-day humans. Unfortunately, current methods to
17 identify selection have limited power to detect such old events of positive selection (Sabeti et al.
18 2006).

19 Here, we introduce a hidden Markov model to detect ancient selective sweeps based on a signal of
20 extended lineage sorting. Using simulations we were able to show that the method can detect older

1 events of selection as long as the selected variant was sufficiently advantageous. The power to detect
2 older events is due to the fact that the method increases in power with the number of mutations that
3 accumulated after the sweep finished. We also showed that background selection can cause false
4 signals and have chosen a minimum length cutoff on candidate regions. While this cutoff reduces the
5 number of false positives due to background selection, we note that this cutoff is expected to exclude
6 *bona fide* events of positive selection, too.

7 We applied the ELS method to 185 African genomes, the Altai Neandertal genome and the Denisovan
8 genome, and detected 81 candidate regions of selection when requiring a minimum genetic length
9 supported by two independent recombination maps. The uncertainty in the recombination maps has a
10 large effect on our results, as shown by the much larger number of 314 regions identified by either
11 recombination map. Recombination rates over the genome are known to evolve rapidly (Lesecque et
12 al. 2014) and of particular concern are recent changes in recombination rates that make some regions
13 appear larger in genetic length than they were in the past. By comparing the current recombination
14 rates in our regions to recombination rates in the ancestral population of both chimpanzee and humans
15 (Munch et al. 2014), we identified some candidate regions that may have increased in recombination
16 rates (Supplemental Table S8). However, it is currently impossible to date the change in
17 recombination rates confidently and these candidate sweeps may post-date the change.

18 A particular strength of our screen for selective sweeps is the ability to detect older events, as
19 indicated by the estimated power to detect simulated events of positive selection of old age and
20 moderate strength. This sets the ELS method apart from previous approaches that made use of archaic
21 genomes, which were geared towards detecting younger events with an age of less than 300,000 years
22 ago (Racimo 2016; Racimo et al. 2014). Despite this difference, we found significant overlap between
23 the ELS candidates and the candidates identified by the other approaches, while the overlap with other
24 types of positive selection scans is smaller. Among our candidates, 71 are novel candidates (283 if
25 considering the extended set) that were not detected in any of the previous screens.

1 While we find no difference in the fraction of genes in selected regions compared to randomly placed
2 regions, we detect an enrichment for enhancers and promoter regions. This result is in agreement with
3 the hypothesis that regulatory changes may play an important role in human-specific phenotypes
4 (Carroll 2003; Enard et al. 2014; King and Wilson 1975). Interestingly, we find an enrichment for
5 genes expressed highest in brain among our candidate regions, leading us to speculate that among the
6 positively selected changes that set us apart from Neandertals are some that are related to behavior or
7 mental capability. While we are unable to substantiate this hypothesis based on our computational
8 analysis alone, we note that several gene candidates falling within sweep regions play a role in the
9 function and development of the brain. A particularly interesting observation is the potential selection
10 on both the ligand *SLIT2* and its receptor *ROBO2*, which reside on chromosome 4 and 3 respectively.
11 Members of the Roundabout (ROBO) gene family play an important role in guiding developing axons
12 in the nervous system through interactions with the ligands SLITs. SLITs proteins act as attractive or
13 repulsive signals for axons expressing different ROBO receptors. *ROBO2* has been further associated
14 with vocabulary growth (St Pourcain et al. 2014), autism (Suda et al. 2011), and dyslexia (Fisher and
15 DeFries 2002) and is involved in the development of neural circuits related to vocal learning in birds
16 (Wang et al. 2015). Interestingly, *ROBO2* is also in a long desert of both Denisovan and Neandertal
17 ancestry in non-Africans.

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

23 We have presented a new approach to detect ancient selective sweeps based on a signal of extended
24 lineage sorting. Applying this approach to modern human data revealed that selection may have acted
25 primarily on regulatory changes and selection may have favored these changes in the brain. With

1 population level sequencing of non-human species becoming more readily available we anticipate that
2 this approach will help to reveal the targets of ancient selection in other species.

3 MATERIALS AND METHODS

4 Data

5 We used 185 unrelated Luhya and Yoruba individuals from the 1000 Genomes Project phase I
6 (Abecasis et al. 2012), corresponding to 370 sets of autosomes and 279 X chromosomes. From this
7 dataset, we extracted allele counts at single nucleotide polymorphism (SNP) sites using vcftools
8 (Danecek et al. 2011). In order to add sites where all Africans differ from the common ancestor with
9 chimpanzee, we first compiled a list of all sites where six high-coverage African genomes (Mbuti, San
10 and Yoriban A and B-panel individuals from Prüfer et al. 2014) are identical. A site was regarded
11 fixed different when the whole genome alignments of at least three out of four ape reference genome
12 assemblies (chimpanzee (panTro3), bonobo (panPan1.1), gorilla (gorGor3) and orangutan (ponAbe2);
13 lastz alignments to the human genome GRCh37/hg19 prepared in-house and by the UCSC genome
14 browser (Speir et al. 2016)) had coverage and were different from the African allele, and when the site
15 was not marked as polymorphic among the 1000 Genomes Luhya and Yoruba individuals.

16 Neandertal and Denisova alleles at polymorphic and fixed positions were extracted from published
17 VCFs and positions were further filtered to sites passing the published map35_100 filter for both the
18 Denisova and Neandertal genotypes (Prüfer et al. 2014). Sites where either Neandertal or Denisova
19 carried a third allele were disregarded.

20 Over all autosomes, 11 million SNPs passed the filters in addition to 6.6 million African fixed
21 variants. For the X chromosome, pseudoautosomal regions, defined as chrX: 60,001-2,699,520, chrX:
22 154,931,044-155,260,560 in hg19 coordinates (<http://www.ncbi.nlm.nih.gov/assembly/2758/>), were
23 filtered out and around 315,000 SNPs as well as 248,000 African fixed variants remained for analysis.

24 Genetic distances between those positions were calculated using the African-American (Hinch et al.
25 2011) and the DeCode (Kong et al. 2010) recombination maps (available in Build 37 from

1 <http://www.well.ox.ac.uk/~anjali/>). Both maps were chosen since they estimate recombination rates
2 from events that occurred within a few generations before present. Recombination maps based on
3 older events (i.e. LD based map) can underestimate recombination rates in regions that underwent
4 recent selective sweeps, potentially masking true signals.

5 Changes of recombination rates along the human lineage could also 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 Hidden Markov model

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

21 The inference further requires the probability for observing an ancestral or derived allele in the archaic
22 at a site i with a derived allele frequency $f_i > 0$ in modern humans (O_i) given that the true genealogy
23 is $k \in \{I, E, ELS\}$: $e_k(O_i) = P(x_i = O_i | \pi_i = k)$. We assume that $\forall o: e_{ELS}(o) = e_E(o)$, i.e. that both

1 external states give rise to ancestral and derived alleles in the archaic with equal probabilities given the
2 same observation. Since external regions are not expected to give rise to derived sites when the
3 derived allele is segregating in modern humans, the only sources for such an observation can be errors
4 or independent coinciding identical mutations and we define an error rate for external regions: $\epsilon_E =$
5 $e_E(O_i = \text{derived}, f_i < 1)$. Similarly fixed derived sites are expected to show the derived allele in the
6 archaics if the local genealogy is internal and we define an error rate for internal regions: $\epsilon_I = e_I(O_i =$
7 $\text{derived}, f_i = 1)$.

8 We compute the posterior probability $P(\pi_i = k | x)$ that an observation O_i came from state k given
9 the 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 $f_k(i) =$
10 $P(O_1 \dots O_i, \pi_i = k)$ and $b_k(i) = P(O_{i+1} \dots O_L | \pi_i = k)$ are the output of the Forward and Backward
11 algorithms respectively (Durbin et al. 1998; Rabiner 1989). $P(x)$ corresponds to the likelihood of the
12 data given our model and was also calculated from the Forward algorithm.

13 Parameter estimate

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

20 We estimated the remaining parameters (average lengths of regions and the proportion of transitions to
21 the ELS state) using the derivative free optimization method COBYLA (Powell 1994) as implemented
22 in the nlopt library (Steven G. Johnson, The NLOpt nonlinear-optimization package) to maximize the
23 log-likelihood values calculated by the Forward algorithm. Convergence was attained in a maximum
24 of 1000 evaluations and the log-likelihood maximization accuracy was set to 10^{-4} . To test for
25 convergence to local maxima, we ran the algorithm twice with different starting points and used the

1 parameters of the run with the highest likelihood to run the re-estimation algorithm a third time
2 starting with those parameters. All three runs gave similar results on all chromosomes.

3 Post-processing

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

14 Simulations

15 We simulated sequences using a model of recent human demography to test the performance of our
16 HMM under different scenarios of neutral evolution, positive selection or background selection. Each
17 simulation consisted of one chimpanzee chromosome, one chromosome from each archaic hominin
18 and 370 human chromosomes, matching the 185 Luhya and Yoruba individuals used in our analysis.
19 For all simulations in this study, a constant mutation rate of 1.45×10^{-8} bp⁻¹.generation⁻¹, a constant
20 recombination rate of 1cM.Mb⁻¹.generation⁻¹ and a generation time of 29 years were assumed. We
21 used estimates of population sizes from (Yang et al. 2014) and population split estimates from (Prüfer
22 et al. 2014) as parameters for the simulated demography (Supplemental Information 1 and 2). Neutral
23 simulations with these parameters using the coalescent simulator scrm (Staab et al. 2014) give a good
24 match to our observed data when plotting derived allele frequency in modern humans against the
25 proportion of derived alleles in the outgroup (Supplemental Figure S7).

1 We generated a total of 100 loci of 1Mb-long sequences under neutrality to investigate the accuracy of
2 labeling external and internal regions using our HMM. To evaluate the length of external regions
3 expected under neutrality for the chromosome X, we simulated 100 loci of 1Mb-long sequences under
4 the demographic model shown in Supplemental Information 1 with the exception that all effective
5 population sizes were reduced to 75% of the original value. To evaluate the accuracy of parameter
6 estimation, we additionally simulated splits of two populations (including an out-group individual)
7 with a constant population size and different split times ranging from 400ky to 1My (step-size of
8 50ky). For each condition, we generated 25 sets of 10 Mb each. In an additional set of 100 loci of
9 1Mb, we introduced random errors by changing the state of the archaic allele with different rates in
10 order to assess our error estimates.

11 To assess our power to detect events of positive selection, we explored selection coefficients ranging
12 from 0.0005 to 0.1 and different times for the occurrence of the selected allele (every 100ky from
13 200kya to 600kya) using the coalescent simulator msms (Ewing and Hermisson 2010). The selected
14 mutation was introduced in the middle of the sequence and we assumed an additive effect of the
15 selected mutation (i.e. the homozygous genotype has twice the advantage stated by the selection
16 coefficient). We performed 2000 simulations of 100kb-long loci for which all demographic parameters
17 match our neutral simulations as described above. We used the `-SForceKeep` switch to drop the
18 simulation if the selected mutation was lost. As 100kb loci are too short to make reliable parameter
19 inferences, we concatenated our simulated sequences, intermittently combining them with 1Mb-long
20 neutral loci from the previous simulations to limit the extent of the sequence affected by positive
21 selection.

22 We investigated how background selection affects lineage sorting in and around a conserved region by
23 performing forward in time simulations using SLiM (Messer 2013). The simulated locus of 500kb
24 length contained a conserved region resembling an ‘average’ human gene (see pg. 19 of the
25 documentary accompanying SLiM (Messer 2013)) and covered 100kb (20%) of the simulated locus.
26 Mutations in the conserved region were assumed to be neutral (25%) or deleterious (75%), with the

1 selection coefficients of the deleterious mutations drawn from a gamma distribution with mean $s =$
2 -0.05 and shape parameter $\alpha = 0.2$. The deleterious mutations were assumed to be partially recessive
3 with dominance coefficient $h = 0.1$ for a set of 100 simulations. To explore the effect of the strength of
4 selection on the results, we produced 2 other sets of 40 simulations each by varying the mean of the
5 gamma distribution ($s = -0.001$ and -0.1).

6 Age Comparison with other Selective Scans

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

11 F_{st} , iHS and XP-EHH are site-based statistics which localise sites that may have been selected,
12 whereas selective scans such as CLR, XP-CLR, Tajima's D , Fay & Wu's H and HKA identify
13 candidate regions. In order to compare the age of the selection events, we assumed that the selected
14 variant in candidate regions was the site with the highest frequency. We note that this procedure will
15 underestimate the age of events if the true selected site reached fixation, as often expected for our
16 method; the comparison is thus conservative.

17 Annotations

18 We used the latest Ensembl gene annotation for hg19 (release 82) to identify protein-coding genes
19 overlapping with our candidate regions. Based on this annotation, a regulatory region was defined as
20 at least 5kb upstream and 1kb downstream of each gene. The regulatory region was extended until it
21 reached a size of 1Mb or came within 5kb upstream or 1kb downstream of a neighboring gene. We
22 additionally used a set of promoters and enhancers mapped by GenoSTAN in 127 cell types and
23 tissues from the ENCODE and Roadmap Epigenomics projects (Zacher et al. 2016).

24 We used B-scores (McVicker et al. 2009) in hg19 coordinates constructed with UCSC's liftover tool
25 to evaluate the extent of background selection in our candidate regions. We also compared our

1 candidate regions on the X chromosome with regions previously suggested to have experienced
2 recurrent selective sweeps in apes (Dutheil et al. 2015; Nam et al. 2015). And, finally, we examined
3 patterns of introgression in our candidate regions with two maps of Neandertal ancestry
4 (Sankararaman et al. 2014; Vernot et al. 2014) and overlapped our regions with long deserts of
5 Neandertal and Denisova ancestry from another recent study (Vernot et al. 2016).

6 To statistically test the overlap of our regions with these annotations, we permuted regions of similar
7 physical sizes in the regions of the genome that passed our quality filters. Quality filtered regions that
8 were smaller than the longest gap present in our candidate ELS regions were regarded as sufficiently
9 short to not prohibit the placement of regions.

10 Gene expression enrichments

11 We defined genes that show tissue-specific expression levels using the Illumina BodyMap 2.0 RNA-
12 seq data (Derrien et al. 2012), which contains expression data from 16 human tissues. We computed
13 differential expression for all genes between a given tissue and all other tissues pooled using the
14 DESeq package (Anders and Huber 2010) and genes were defined to be expressed in a tissue-specific
15 manner when their expression levels were significantly higher (P -value < 0.05) in a given tissue
16 compared to all other tissues. We tested for enrichment of candidate genes in the 16 sets of tissue-
17 specifically expressed genes comparing to genes that were located outside of candidate regions using
18 Fisher's exact test. P -values were corrected for multiple testing using the Benjamini-Hochberg
19 procedure (Benjamini and Hochberg 1995).

20 We used the ABAenrichment package (Grote et al. 2016) to pinpoint regions in the brain where
21 expression of our candidate genes is enriched. As background genes, we only used genes that could
22 have been potentially identified by the sweep screen according to our genomic filters. We used the
23 default parameters and the hypergeometric test.

1 DATA ACCESS

2 The software and input files used in this study have been made available through the website
3 <http://bioinf.eva.mpg.de/ELS/>.

4 ACKNOWLEDGMENTS

5 We would like to thank Michael Lachmann for early discussions about the design of the study, Janet
6 Kelso and Mark Stoneking for many useful comments on the manuscript, Svante Pääbo, Udo Stenzel,
7 Fernando Racimo and Adam Siepel for helpful discussions, and Matthias Ongyerth and Christoph
8 Theunert for help with earlier analysis. This research was funded by the Max Planck Society.

9 DISCLOSURE DECLARATION

10 The authors declare no competing financial interests.

11 AUTHOR CONTRIBUTIONS

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

14 REFERENCES

15 The 1000 Genomes Project Consortium. 2012. An Integrated Map of Genetic Variation from 1,092
16 Human Genomes. *Nature* **491**: 56–65.

17 Anders S, Huber W. 2010. “DESeq: Differential Expression Analysis for Sequence Count Data.”
18 *Genome biology* **11**:R106. doi: 10.1186/gb-2010-11-10-r106

19 Benjamini Y, Hochberg Y. 1995. Controlling the False Discovery Rate: A Practical and Powerful
20 Approach to Multiple Testing. *Journal of the Royal Statistical Society B* **57**: 289–300.

21 Cagan A, Theunert C, Laayouni H, Santpere G, Pybus M, Casals F, Prüfer K, Navarro A, Marques-
22 Bonet T, Bertranpetit J, et al. 2016. Natural Selection in the Great Apes. *Molecular Biology and*

- 1 *Evolution* **33**: 3268-3283.
- 2 Carroll SB. 2003. Genetics and the Making of Homo Sapiens. *Nature* **422**: 849–857.
- 3 Charlesworth B, Morgan MT, Charlesworth D. 1993. The Effect of Deleterious Mutations on Neutral
4 Molecular Variation. *Genetics* **134**: 1289–1303.
- 5 Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE, Lunter G,
6 Marth GT, Sherry ST, et al. 2011. The Variant Call Format and VCFtools. *Bioinformatics* **27**:
7 2156–2158.
- 8 Derrien T, Johnson R, Bussotti G, Tanzer A, Djebali S, Tilgner H, Guernec G, Martin D, Merkel A,
9 Knowles DG, et al. 2012. The GENCODE v7 Catalog of Human Long Noncoding RNAs:
10 Analysis of Their Gene Structure, Evolution, and Expression. *Genome Research* **22**: 1775–1789.
- 11 Durbin R, Eddy S, Krogh A, Mitchison G. 1998. Biological Sequence Analysis: Probabilistic Models
12 of Proteins and Nucleic Acids. Cambridge University Press, Cambridge.
- 13 Dutheil JY, Munch K, Nam K, Mailund T, Schierup MH. 2015. Strong Selective Sweeps on the X
14 Chromosome in the Human-Chimpanzee Ancestor Explain Its Low Divergence. *PLoS Genetics*
15 **11**: e1005451. doi:10.1371/journal.pgen.1005451
- 16 Enard D, Messer PW, Petrov DA. 2014. Genome-Wide Signals of Positive Selection in Human
17 Evolution. *Genome Research* **24**: 885–895.
- 18 Ewing G, Hermisson J. 2010. MSMS: A Coalescent Simulation Program Including Recombination,
19 Demographic Structure and Selection at a Single Locus. *Bioinformatics* **26**: 2064–2065.
- 20 Fisher SE, DeFries JC. 2002. Developmental Dyslexia: Genetic Dissection of a Complex Cognitive
21 Trait. *Nature Reviews Neuroscience* **3**: 767–780.
- 22 Fu Q, Li H, Moorjani P, Jay F, Slepchenko SM, Bondarev AA, Johnson PL, Aximu-Petri A, Prüfer K,
23 de Filippo C, et al. 2014. Genome sequence of a 45,000-year-old modern human from western

- 1 Siberia. *Nature* 514: 445–449.
- 2 Fu Q, Hajdinjak M, Moldovan OT, Constantin S, Mallick S, Skoglund P, Patterson N, Rohland N,
3 Lazaridis I, Nickel B, et al. 2015. An early modern human from Romania with a recent
4 Neanderthal ancestor. *Nature* **524**: 216–219
- 5 Gleeson JG, Allen KM, Fox JW, Lamperti ED, Berkovic S, Scheffer I, Cooper EC, Dobyns WB,
6 Minnerath SR, Ross ME, et al. 1998. Doublecortin, a Brain-Specific Gene Mutated in Human X-
7 Linked Lissencephaly and Double Cortex Syndrome, Encodes a Putative Signaling Protein. *Cell*
8 **92**: 63–72.
- 9 Gleeson JG, Lin PT, Flanagan LA, Walsh CA. 1999. Doublecortin Is a Microtubule-Associated
10 Protein and Is Expressed Widely by Migrating Neurons. *Neuron* **23**: 257–271.
- 11 Green RE, Krause J, Briggs AW, Maricic T, Stenzel U, Kircher M, Patterson N, Li H, Zhai W, Fritz
12 MH, et al. 2010. A draft sequence of the Neandertal genome. *Science* **328**: 710-722.
- 13 Grote S, Prüfer K, Kelso J, Dannemann M. 2016. ABAEnrichment: An R Package to Test for Gene
14 Set Expression Enrichment in the Adult and Developing Human Brain. *Bioinformatics* **32**: 3201-
15 3203.
- 16 Hawrylycz MJ, Lein ES, Guillozet-Bongaarts AL, Shen EH, Ng L, Miller JA, van de Lagemaat LN,
17 Smith KA, Ebbert A, Riley ZL, et al. 2012. An Anatomically Comprehensive Atlas of the Adult
18 Human Brain Transcriptome. *Nature* **489**: 391–399.
- 19 Hinch AG, Tandon A, Patterson N, Song Y, Rohland N, Palmer CD, Chen GK, Wang K, Buxbaum
20 SG, Akylbekova EL, et al. 2011. The Landscape of Recombination in African Americans. *Nature*
21 **476**: 170–175.
- 22 King MC, Wilson AC. 1975. Evolution at Two Levels in Humans and Chimpanzees. *Science* **188**:
23 107–116.
- 24 Kong A, Thorleifsson G, Gudbjartsson DF, Masson G, Sigurdsson A, Jonasdottir A, Walters GB,

- 1 Jonasdottir A, Gylfason A, Kristinsson KT, et al. 2010. Fine-Scale Recombination Rate
2 Differences between Sexes, Populations and Individuals. *Nature* **467**: 1099–1103.
- 3 Kryazhimskiy S, Plotkin JB. 2008. The Population Genetics of dN / dS. *PLoS Genetics* **4**: e1000304.
4 doi:10.1371/journal.pgen.1000304.
- 5 Laland KN, Odling-Smee J, Myles S. 2010. How Culture Shaped the Human Genome: Bringing
6 Genetics and the Human Sciences Together. *Nature Reviews Genetics* **11**: 137–148.
- 7 Pybus OG, Shapiro B. 2009. Natural Selection and Adaptation of Molecular Sequences. The
8 Phylogenetic Handbook. Lemey P, Salemi M, Vandamme AM. pp 415-417. Cambridge
9 University Press, Cambridge.
- 10 Lesecque Y, Glémin S, Lartillot N, Mouchiroud D, Duret L. 2014. The Red Queen Model of
11 Recombination Hotspots Evolution in the Light of Archaic and Modern Human Genomes. *PLoS*
12 *Genetics* **10**: e1004790. doi:10.1371/journal.pgen.1004790.
- 13 Li H, Durbin R. 2011. Inference of Human Population History from Individual Whole-Genome
14 Sequences. *Nature* **475**: 493–496.
- 15 McVicker G, Gordon D, Davis C, Green P. 2009. Widespread Genomic Signatures of Natural
16 Selection in Hominid Evolution. *PLoS Genetics* **5**: e1000471. doi:10.1371/journal.pgen.1000471.
- 17 Messer PW. 2013. SLiM: Simulating Evolution with Selection and Linkage. *Genetics* **194**: 1037–
18 1039.
- 19 Meyer M, Kircher M, Gansauge MT, Li H, Racimo F, Mallick S, Schraiber JG, Jay F, Prüfer K, de
20 Filippo C, et al. 2012. A High-Coverage Genome Sequence from an Archaic Denisovan
21 Individual. *Science* **338**: 222–226.
- 22 Miller JA, Ding SL, Sunkin SM, Smith KA, Ng L, Szafer A, Ebbert A, Riley ZL, Royall JJ, Aiona K,
23 et al. 2014. Transcriptional Landscape of the Prenatal Human Brain. *Nature* **508**: 199–206.

- 1 Munch K, Mailund T, Dutheil JY, Schierup MH. 2014. A Fine-Scale Recombination Map of the
2 Human-Chimpanzee Ancestor Reveals Faster Change in Humans than in Chimpanzees and a
3 Strong Impact of GC-Biased Gene Conversion. *Genome Research* **24**: 467–474.
- 4 Nam K, Munch K, Hobolth A, Dutheil JY, Veeramah KR, Woerner AE, Hammer MF, Great Ape
5 Genome Diversity Project, Mailund T, et al. 2015. Extreme Selective Sweeps Independently
6 Targeted the X Chromosomes of the Great Apes. *Proceedings of the National Academy of
7 Sciences* **112**: 6413–6418.
- 8 Nielsen R, Hellmann I, Hubisz M, Bustamante C, Clark AG. 2007. Recent and Ongoing Selection in
9 the Human Genome. *Nature Reviews Genetics* **8**: 857–868.
- 10 Oleksyk TK, Smith MW, O’Brien SJ. 2010. Genome-Wide Scans for Footprints of Natural Selection.
11 *Philosophical transactions of the Royal Society of London. Series B, Biological sciences* **365**:
12 185–205.
- 13 Pääbo S. 2014. The Human Condition—A Molecular Approach. *Cell* **157**: 216–226.
- 14 Powell MJD. 1994. A Direct Search Optimization Method That Models the Objective and Constraint
15 Functions by Linear Interpolation. *Advances in optimization and numerical analysis* **275**: 51–67.
- 16 Prüfer K, Racimo F, Patterson N, Jay F, Sankararaman S, Sawyer S, Heinze A, Renaud G, Sudmant
17 PH, de Filippo C, et al. 2014. The Complete Genome Sequence of a Neanderthal from the Altai
18 Mountains. *Nature* **505**: 43–49.
- 19 Pybus M, Dall’Olio GM, Luisi P, Uzkudun M, Carreño-Torres A, Pavlidis P, Laayouni H, Bertranpetit
20 J, Engelken J. 2014. 1000 Genomes Selection Browser 1.0: A Genome Browser Dedicated to
21 Signatures of Natural Selection in Modern Humans. *Nucleic Acids Research* **42**: D903-D909.
- 22 Rabiner LR. 1989. A Tutorial on Hidden Markov Models and Selected Applications in Speech
23 Recognition. *Proceedings of the IEEE* **77**: 257–286.
- 24 Racimo F. 2016. Testing for Ancient Selection Using Cross-Population Allele Frequency

- 1 Differentiation. *Genetics* **202**: 733–750.
- 2 Racimo F, Kuhlwilm M, Slatkin M. 2014. A Test for Ancient Selective Sweeps and an Application to
3 Candidate Sites in Modern Humans. *Molecular Biology and Evolution* **31**: 3344–3358.
- 4 Rasmussen MD, Hubisz MJ, Gronau I, Siepel A. 2014. Genome-Wide Inference of Ancestral
5 Recombination Graphs. *PLoS Genetics* **10**: e1004342. doi:10.1371/journal.pgen.1004342.
- 6 Rosenberg NA. 2002. The Probability of Topological Concordance of Gene Trees and Species Trees.
7 *Theoretical Population Biology* **61**: 225–247.
- 8 Sabeti PC, Schaffner SF, Fry B, Lohmueller J, Varilly P, Shamovsky O, Palma A, Mikkelsen TS,
9 Altshuler D, Lander ES. 2006. Positive Natural Selection in the Human Lineage. *Science* **312**:
10 1614–1620.
- 11 Sankararaman S, Patterson N, Li H, Pääbo S, Reich D. 2012. The Date of Interbreeding between
12 Neandertals and Modern Humans. *PLoS Genetics* **8**: e1002947.
13 doi:10.1371/journal.pgen.1002947.
- 14 Sankararaman S, Mallick S, Dannemann M, Prüfer K, Kelso J, Pääbo S, Patterson N, Reich D. 2014.
15 The genomic landscape of Neanderthal ancestry in present-day humans. *Nature* **507**: 354–357.
- 16 Sankararaman S, Mallick S, Patterson N, Reich D. 2016. The Combined Landscape of Denisovan and
17 Neanderthal Ancestry in Present-Day Humans. *Current Biology* **26**: 1241–1247.
- 18 Speir ML, Zweig AS, Rosenbloom KR, Raney BJ, Paten B, Nejad P, Lee BT, Learned K, Karolchik
19 D, Hinrichs AS, et al. 2016. The UCSC Genome Browser Database: 2016 Update. *Nucleic Acids*
20 *Research* **44**: D717–25.
- 21 Staab PR, Zhu S, Metzler D, Lunter G. 2014. Scrm: Efficiently Simulating Long Sequences Using the
22 Approximated Coalescent with Recombination. *Bioinformatics* **31**: 1680–1682.
- 23 St Pourcain B, Cents RA, Whitehouse AJ, Haworth CM, Davis OS, O'Reilly PF, Roulstone S, Wren

- 1 Y, Ang QW, Velders FP, et al. 2014. Common Variation near ROBO2 Is Associated with
2 Expressive Vocabulary in Infancy. *Nature communications* **5**:4831. doi: 10.1038/ncomms5831.
- 3 Suda S, Iwata K, Shimmura C, Kameno Y, Anitha A, Thanseem I, Nakamura K, Matsuzaki H,
4 Tsuchiya KJ, Sugihara G, et al. 2011. Decreased Expression of Axon-Guidance Receptors in the
5 Anterior Cingulate Cortex in Autism. *Molecular autism* **2**: 14. doi: 10.1186/2040-2392-2-14
- 6 Varki A, Geschwind DH, Eichler EE. 2008. Explaining Human Uniqueness: Genome Interactions with
7 Environment, Behaviour and Culture. *Nature reviews Genetics* :749–763.
- 8 Vernot B, Tucci S, Kelso J, Schraiber JG, Wolf AB, Gittelman RM, Dannemann M, Grote S, McCoy
9 RC, Norton H, et al. 2016. Excavating Neandertal and Denisovan DNA from the Genomes of
10 Melanesian Individuals. *Science* **352**:235–239.
- 11 Vernot B, Akey JM. 2014. Resurrecting Surviving Neandertal Lineages from Modern Human
12 Genomes. *Science* **343**: 1017–1021.
- 13 Wang R, Chen CC, Hara E, Rivas MV, Roulhac PL, Howard JT, Chakraborty M, Audet JN, Jarvis
14 ED. 2015. Convergent Differential Regulation of SLIT-ROBO Axon Guidance Genes in the
15 Brains of Vocal Learners. *Journal of Comparative Neurology* **523**: 892–906.
- 16 Weaver TD. 2009. The Meaning of Neandertal Skeletal Morphology. *Proceedings of the National*
17 *Academy of Sciences* **106**: 16028–16033.
- 18 Yang MA, Harris K, Slatkin M. 2014. The Projection of a Test Genome onto a Reference Population
19 and Applications to Humans and Archaic Hominins. *Genetics* **198**: 1655–1670.
- 20 Zacher B, Michel M, Schwalb B, Cramer P, Tresch A, Gagneur J. 2016. Accurate Promoter and
21 Enhancer Identification in 127 ENCODE and Roadmap Epigenomics Cell Types and Tissues by
22 GenoSTAN. bioRxiv doi: <http://dx.doi.org/10.1101/041020>.