

1 **Dietary adaptation of *FADS* genes in Europe varied across time and geography**

2 Kaixiong Ye¹, Feng Gao¹, David Wang¹, Ofer Bar-Yosef², Alon Keinan^{1*}

3 ¹ Department of Biological Statistics and Computational Biology, Cornell University, Ithaca,
4 NY, USA

5 ² Department of Anthropology, Harvard University, Cambridge, MA, USA

6 *corresponding author: alon.keinan@cornell.edu

7 **Abstract:** Fatty acid desaturase (*FADS*) genes encode rate-limiting enzymes for the biosynthesis
8 of omega-6 and omega-3 long chain polyunsaturated fatty acids (LCPUFAs). This biosynthesis
9 is essential for individuals subsisting on LCPUFAs-poor, plant-based diets. Positive selection on
10 *FADS* genes has been reported in multiple populations, but its presence and pattern in Europeans
11 remain elusive. Here, with analyses of ancient and modern DNA, we demonstrated that positive
12 selection acted on the same *FADS* variants both before and after the advent of farming in Europe,
13 but on opposite alleles. Selection in recent farmers also varied geographically, with the strongest
14 signal in Southern Europe. These varying selection patterns concur with anthropological
15 evidence of differences in diets, and with the association of recently-adaptive alleles with higher
16 *FADS1* expression and enhanced LCPUFAs biosynthesis. Genome-wide association studies
17 revealed associations of recently-adaptive alleles with not only LCPUFAs, but also other lipids
18 and decreased risk of several inflammation-related diseases.

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33 Identifying genetic adaptations to local environment, including historical dietary practice, and
34 elucidating their implications on human health and disease are of central interest in human
35 evolutionary genomics¹. The fatty acid desaturase (*FADS*) gene family consists of *FADS1*,
36 *FADS2* and *FADS3*, which evolved by gene duplication². *FADS1* and *FADS2* encode rate-
37 limiting enzymes for the biosynthesis of omega-3 and omega-6 long-chain polyunsaturated fatty
38 acids (LCPUFAs) from plant-sourced shorter-chain precursors (Supplementary Fig. 1).
39 LCPUFAs are indispensable for proper human brain development, cognitive function and
40 immune response^{3,4}. While omega-3 and omega-6 LCPUFAs can be obtained from animal-based
41 diets, their endogenous synthesis is essential to compensate for their absence from plant-based
42 diets. Positive selection on the *FADS* locus, a 100 kilobase (kb) region containing all three genes
43 (Supplementary Fig. 2), has been identified in multiple populations⁵⁻⁹. Our recent study showed
44 that a 22 bp insertion-deletion polymorphism (indel, rs66698963) within *FADS2*, which is
45 associated with *FADS1* expression¹⁰, has been adaptive in Africa, South Asia and parts of East
46 Asia, possibly driven by local historical plant-based diets⁸. We further supported this hypothesis
47 by functional association of the adaptive insertion allele with more efficient biosynthesis⁸. In
48 Greenlandic Inuit, who have traditionally subsisted on a LCPUFAs-rich marine diet, adaptation
49 signals were also observed on the *FADS* locus, with adaptive alleles associated with less efficient
50 biosynthesis⁹.

51 In Europeans, positive selection on the *FADS* locus has only been reported recently in a study
52 based on ancient DNA (aDNA)¹¹. Evidence of positive selection from modern DNA is still
53 lacking even though most above studies also performed similarly-powered tests in Europeans⁵⁻⁸.
54 Moreover, although there are well-established differences in the Neolithization process and in
55 dietary patterns across Europe¹²⁻¹⁴, geographical differences of selection within Europe have not
56 been investigated before. Furthermore, before the advent of farming, pre-Neolithic hunter-
57 gatherers throughout Europe had been subsisting on animal-based diets with significant aquatic
58 contribution¹⁵⁻¹⁷, in contrast to the plant-heavy diets of recent European farmers¹⁸⁻²⁰. We
59 hypothesized that these drastic differences in subsistence strategy and dietary practice before and
60 after the Neolithic revolution within Europe exert different selection pressures on the *FADS*
61 locus. In this study, we combined analyses on ancient and modern DNA to investigate potential
62 positive selection on the *FADS* locus in Europe and to examine whether it exhibits geographical
63 and temporal differences as would be expected from differences in diets. Briefly, we present
64 evidence for positive selection on *opposite alleles* of the same variants before and after the
65 Neolithic revolution, and for varying selection signals between Northern and Southern
66 Europeans in recent history. We interpreted the functional significance of adaptive alleles with
67 analysis of expression quantitative trait loci (eQTLs) and genome-wide association studies
68 (GWAS), both pointing to selection for diminishing LCPUFAs biosynthesis in pre-Neolithic
69 hunter-gatherers but for increasing biosynthesis in recent farmers. Anthropological findings
70 indicate that these selection patterns were likely driven by dietary practice and its changes.

71 **Results**

72 **Evidence of recent positive selection in Europe from both ancient and modern DNA**

73 To systematically evaluate the presence of recent positive selection on the *FADS* locus in
74 Europe, we performed an array of selection tests using both ancient and modern samples. We
75 first generated a uniform set of variants across the locus in a variety of aDNA data sets
76 (Supplementary Table S1) via imputation (Methods). For all these variants, we conducted an
77 aDNA-based test¹¹. This test includes three groups of ancient European samples and four groups
78 of modern samples. The three ancient groups represent the three major ancestry sources of most
79 present-day Europeans: Western and Scandinavian hunter-gatherers (WSHG), early European
80 farmers (EF), and Steppe-Ancestry pastoralists (SA). The four groups of modern samples were
81 drawn from the 1000 Genomes Project (1000GP), representing Tuscans (TSI), Iberians (IBS),
82 British (GBR) and additional northern Europeans (CEU). The test identifies variants with
83 extreme frequency change between ancient and modern samples, suggesting the presence of
84 positive selection during recent European history (not more ancient than 8,500 years ago (ya))¹¹.
85 Our results confirmed the presence of significant selection signals on many variants in the *FADS*
86 locus (Fig. 1), including the previously identified peak SNP rs174546 ($p = 1.04e-21$)¹¹. We
87 observed the most significant signal at an imputed SNP, rs174594 ($p = 1.29e-24$), which was not
88 included in the original study¹¹. SNP rs174570, one of the top adaptive SNPs reported in
89 Greenlandic Inuit⁹, also exhibits a significant signal ($p = 7.64e-18$) while indel rs66698963
90 shows no evidence of positive selection ($p = 3.62e-3$, likely due to data quality, see
91 Supplementary Notes). Overall, the entire peak of selection signals coincides with a linkage
92 disequilibrium (LD) block (referred to as the *FADS1-FADS2* LD block) in Europeans, which
93 extends over a long genomic region of 85 kb, covering the entirety of *FADS1* and most of the
94 much longer *FADS2* (Supplementary Figs. 2 and 3). The dominant haplotype of this LD block
95 (haplotype D; Methods) has a frequency of 63% in modern Europeans and is composed of alleles
96 under positive selection as revealed by the above test. Of note, some alleles on this haplotype are
97 derived (i.e. the new mutation relative to primates) while others are ancestral (Supplementary
98 Fig. 4). Thus, the large number of variants showing genome-wide significant signals could
99 potentially be the result of one or a few variants targeted by strong selection, with extensive
100 hitchhiking of nearby neutral variants.

101 We next performed several selection tests solely based on extant European populations.
102 Considering the five European populations from 1000GP, including samples of Finns (FIN) and
103 the four samples described above, a haplotype-based selection test, nSL²¹, revealed positive
104 selection on many SNPs in the *FADS1-FADS2* LD block. Importantly, this test unraveled the
105 same adaptive alleles as in the above aDNA-based test and a same general trend of stronger
106 signal towards rs174594 (Fig. 2A, Supplementary Fig. 5). For rs174594, the nSL values are
107 significant in all five populations and the signal exhibits a gradient of being stronger in southern
108 Europeans and weaker in northern Europeans (Fig. 2A, Supplementary Fig. 6): TSI ($p =$
109 0.00044), IBS ($p = 0.0020$), CEU ($p = 0.0039$), GBR ($p = 0.0093$), and FIN ($p = 0.017$). Of note,
110 nSL values have been normalized separately in each population to remove demographic
111 effects²¹. The other three variants of interest (rs174546, rs174570, and rs66698963) exhibit no
112 selection signals, except for rs174570 showing borderline significance in the two southernmost
113 populations (TSI: $p = 0.022$; IBS: $p = 0.050$, Fig. 2A). Signals were also observed with nSL in
114 two whole-genome sequencing cohorts of British ancestry from the UK10K project
115 (Supplementary Fig. 7). Another test for positive selection in very recent history (during the past

116 ~2,000-3,000 years), Singleton Density Score (SDS)²², applied in the UK10K data set, also
117 revealed significant signals for multiple SNPs in the *FADS1-FADS2* LD block, with the same
118 adaptive alleles and general trend of localized signals as in the above two tests (Fig. 2B,
119 Supplementary Fig. S8). Significant SDS was observed for rs174594 ($p = 0.045$) and rs174570
120 ($p = 0.045$), but not rs174546. Of note, it is the derived allele for rs174594 that was under
121 selection, while it is the ancestral allele for rs174570. Interestingly, selection on the opposite (or
122 derived) allele of rs174570 has been shown in Greenlandic Inuit⁹. Additional tests of selection
123 consistently highlight the *FADS1-FADS2* LD block as a target of natural selection
124 (Supplementary Figs. S5, S7-S10). Taken together, standard tests on modern DNAs support the
125 aDNA-based results of recent positive selection on the *FADS* locus and, specifically, on the D
126 haplotype of the *FADS1-FADS2* LD block.

127 **Geographical differences of recent positive selection signals across Europe**

128 To rigorously evaluate geographical differences of recent positive selection on the *FADS* locus
129 across Europe, we revisited the aDNA-based selection test¹¹. We started by decomposing the
130 original test for four representative SNPs (Fig. 3A) and then performed the test separately in
131 Northern and Southern Europeans for all variants in the *FADS* locus (Fig. 3B). Our first analysis
132 included four SNPs, three of which (rs174594, rs174546, and rs174570) are top SNPs from this
133 and previous studies^{9,11} and are highlighted in all our analyses, while the fourth (rs4246215) is
134 the one showing the biggest difference in the upcoming South-North comparison analysis. The
135 indel rs66698963 was not highlighted in this and all upcoming analyses because it has no
136 significant selection signals in Europe. The original aDNA-based test evaluates the frequencies
137 of an allele in three ancient samples and four modern 1000GP samples under two hypotheses (H_0
138 and H_1). Under H_1 , maximum likelihood estimates (MLEs) of frequencies in all samples are
139 constrained only by observed allele counts and thus equivalent to the direct observed frequencies
140 (Fig. 3A; blue bars). Among the four modern samples, the observed adaptive allele frequencies
141 for all four SNPs exhibit a South-North gradient with the highest in Tuscans and the lowest in
142 Finns, consistent with the gradient of selection signals observed before based on modern DNA.
143 Among the three ancient samples, the observed allele frequencies, equivalent to the frequencies
144 upon admixture (Fig. 3A, orange bars for ancient groups), are always the lowest and often zero
145 in the WSHG sample.

146 Under H_0 , the MLEs of frequencies are constrained by the observed allele counts and an
147 additional assumption that an allele's frequencies in the four modern samples are each a linear
148 combination of its frequencies in the three ancient samples. Considering the later assumption
149 alone, we can predict the frequencies of adaptive alleles right after admixture for each modern
150 population. The admixture contribution of WSHG, as estimated genome-wide, is higher towards
151 the North, constituting of 0%, 0%, 19.6%, and 36.2% for TSI, IBS, CEU, and GBR,
152 respectively¹¹. Thus, the predicted adaptive allele frequencies upon admixture for these four
153 modern populations are usually lower in the North (Fig. 3A; orange bars in modern populations),
154 suggesting higher starting frequencies in the South at the onset of selection. Further considering
155 observed allele counts, we obtained the MLEs of frequencies under H_0 (Fig. 3A; yellow bars in
156 modern populations). As expected, the predicted allele frequencies are higher in the South. But
157 more importantly, the differences between H_0 and H_1 estimates in modern populations (Fig. 3A;
158 indicated differences between yellow and blue bars) are still higher in the South, suggesting that
159 in addition to population-specific admixture proportions and different starting frequencies, more

160 recent factors, such as stronger selection pressure, earlier onset of selection, or unmodeled recent
161 demographic history, might contribute to the observation of stronger selection signals in the
162 South.

163 To evaluate the potential confounding effects of varying demographic history that is not captured
164 by the model, we evaluate all variants in the 3 Mb region surrounding the *FADS* locus. We
165 applied the aDNA-based selection test separately for the two Southern and the two Northern
166 populations. All variants that were significant in the combined analyses (Fig. 1) were also
167 significant in each of the two separate analyses, but many exhibited much stronger signals in
168 Southern populations (Fig. 3B; Supplementary Fig. S11). The maximum difference was found
169 for SNP rs4246215, whose p value in Southern populations is 12 orders of magnitude stronger
170 than that in Northern populations. SNP rs174594, rs174546 and rs174570 also have signals that
171 are several orders of magnitude stronger in the South. A further decomposition of the selection
172 test and comparison of maximum likelihoods under H_0 and H_1 between South and North revealed
173 that a stronger deviation under H_0 in the South is driving the signal (Supplementary Fig. S12). It
174 is noteworthy that the pattern of stronger signal in the South is observed only for some but not all
175 SNPs, excluding the possibility of systemic bias and pointing at variants-specific properties,
176 likely for variants that were under selection and the nearby variants in LD. Indeed, the candidate
177 adaptive haplotype D also exhibits frequency patterns that are consistent with adaptive alleles of
178 the four representative SNPs (Fig. 3C). Hence, these results suggest that there might be stronger
179 selection pressure or earlier onset of positive selection on the *FADS1-FADS2* LD block in
180 Southern Europeans.

181 **Opposite selection signals in pre-Neolithic European hunter-gatherers**

182 Motivated by the very different diet of pre-Neolithic European hunter-gatherers, we set to test
183 the action of natural selection on the *FADS* locus before the Neolithic revolution. We started by
184 examining the frequency trajectory of haplotype D, the candidate adaptive haplotype in recent
185 European history. As noted above, its frequency increased drastically in Europe after the
186 Neolithic revolution (Fig. 3C, the contrast between orange and blue bars). In stark contrast, it
187 shows a clear trajectory of decreasing frequency over time among pre-Neolithic hunter-
188 gatherers²³ (Fig. 4A): starting from 32% in the ~30,000-year-old (yo) “Věstonice cluster”,
189 through 21% in the ~15,000 yo “El Mirón cluster”, to 13% in the ~10,000 yo “Villabruna
190 cluster”, and to being practically absent in the ~7,500 yo WSHG group. We hypothesized that
191 there was positive selection on alleles opposite to the recently adaptive ones on haplotype D.

192 To search for variants with evidence of positive selection during the pre-Neolithic period, we
193 considered the allele frequency time series for all variants around the *FADS* locus. We applied to
194 each variant two rigorous, recently-published Bayesian methods^{24,25} to infer selection
195 coefficients from time series data. Under a simple demographic model of constant population
196 size, both methods consistently highlighted two SNPs (rs174570 and rs2851682) within the
197 *FADS1-FADS2* LD block to be under positive selection during the pre-Neolithic period tested,
198 approximately 30,000-7,500 ya (Supplementary Figs. 13 and 14). The Schraiber *et al.* method is
199 capable of processing more complicated demographic models²⁴. With this method and
200 considering a more realistic European demographic model²⁶, the same two SNPs were
201 highlighted (Supplementary Fig. 15). The derived alleles of these two SNPs have similar
202 frequency trajectories during the examined period, increasing from 36% to 78% (Fig. 4B).

203 Estimated selection coefficients for homozygotes of adaptive allele (s) for these two SNPs are
204 similar across methods and demographic models. With the Schraiber *et al.* method and the
205 realistic demographic model, the marginal maximum *a posteriori* estimate of s for rs174570 is
206 0.38% (95% credible interval (CI): 0.038% – 0.92%) while the estimated derived allele age is
207 57,380 years (95% CI: 157,690 – 41,930 years) (Fig. 4C, Supplementary Fig. 16). For
208 rs2851682, the estimated s is 0.40% (95% CI: 0.028% – 1.12%) while its derived allele age is
209 53,440 years (95% CI: 139,620 – 39,320 years) (Fig. 4D, Supplementary Fig. 17). In addition to
210 these two SNPs, ApproxWF²⁵ revealed significant signals for 44 SNPs in the *FADS1-FADS2* LD
211 block (Supplementary Fig. 14), including rs174546 and rs174594, whose ancestral allele
212 frequencies increased from about 65% to almost fixation (Fig. 4B). Importantly, these SNPs have
213 similar estimated s (0.28% - 0.62%) and their adaptive alleles are alternative (or opposite) to the
214 ones under selection in recent history.

215 Considering the haplotype structure of the *FADS1-FADS2* LD block (Fig. 5A), we identified a
216 haplotype (referred to as M2), which is comprised of alleles that are mostly alternative to those
217 on haplotype D (Supplementary Fig. S4). M2 appears in modern Europeans at a frequency of
218 10% but is much more common in Eskimos from Eastern Siberia, presumably for similar reasons
219 that the derived allele of rs174570 is prevalent in Greenlandic Inuit. M2 exhibits increasing
220 frequency over time in pre-Neolithic hunter-gatherers (Supplementary Table S2), suggesting that
221 the allele(s) targeted by selection during that period are likely on M2.

222 **The temporal and global evolutionary trajectory of *FADS* haplotypes**

223 To go beyond the dominant D haplotype and study different haplotypes in the *FADS1-FADS2*
224 LD block, their frequency changes over time and their current global distributions, we performed
225 haplotype network and frequency analysis on 450 and 5,052 haplotypes from ancient and modern
226 DNA, respectively (Fig. 5, Supplementary Fig. S18, Supplementary Tables S2-S4). The top five
227 haplotypes in modern Europeans, designated as D, M1, M2, M3 and M4 from the most to the
228 least common (63%, 15%, 10%, 5%, 4%, respectively), were all present in aDNA and modern
229 Africans. M1, M2 and M4 are closer to the consensus ancestral haplotype observed in primates
230 while D and M3 are more distant (Fig. 5A). Among the Out-of-Africa ancestors, the frequencies
231 of D and M2 were probably ~35% and ~27% because those frequencies were observed in both
232 the oldest European hunter-gatherer group, the ~30,000 yo “Věstonice cluster”, and the ~14,500
233 yo Epipalaeolithic Natufian hunter-gatherers in the Levant (Fig. 5B, Supplementary Table S3).
234 Among pre-Neolithic European hunter-gatherers, positive selection on M2 increased its
235 frequency from 29% to 56% from approximately 30,000 to 7,500 ya, while the D haplotype
236 practically disappeared by the advent of farming (Figs. 4A and 5B). With the arrival of farmers
237 and Steppe-Ancestry pastoralists, D was re-introduced into Europe. Since the Neolithic
238 revolution, positive selection on D increased its frequency dramatically to 63% while that of M2
239 has decreased to only 10% among present-day Europeans. Globally, D is also present at high
240 frequency in South Asia (82%) but absent in modern-day Eskimos (Fig. 5C). In contrast, M2 has
241 very low frequency in South Asia (3%) but moderate frequency in Eskimos (27%). Further
242 detailed description of evolutionary trajectories of haplotypes in this region could be found in
243 Supplementary Notes.

244 The geographical frequency patterns of representative variants (rs174570, rs174594, rs174546,
245 rs66698963, and rs2851682; Fig. 6, Supplementary Figs. 19-23) mostly mirror those of key

246 haplotypes, but with discrepancies providing insights into casual variants and allele ages. One
247 major discrepancy was found in Africa. The derived alleles of rs174570 and rs2851682 remain
248 almost absent in Africa, consistent with their allele age estimates of ~55,000 years (Figs. 4C and
249 4D) and ruling out their involvement in the positive selection on *FADS* genes in Africa^{5,6,8}.
250 Considering the much weaker LD structure of the *FADS* locus in Africa (Supplementary Fig.
251 24), it is possible that selection in Africa may be on haplotypes and variants that are different
252 from those in Europe.

253 **Functional and medical implications of adaptive variants**

254 Previous studies on adaptive evolution of the *FADS* locus suggested that adaptive alleles are
255 associated with expression levels of *FADS* genes^{5,6,8}. To test this possibility in the context of this
256 large-scale analysis, we considered data from the Genotype-Tissue Expression (GTEx) project²⁷.
257 Our results point to many SNPs on the *FADS1-FADS2* LD block being eQTLs of *FADS* genes.
258 Out of a total of 44 tissues, these eQTLs at genome-wide significance level are associated with
259 the expression of *FADS1*, *FADS2*, and *FADS3* in 12, 23, and 4 tissues, respectively, for a total of
260 27 tissues (Supplementary Figs. 25-27). Considering the peak SNP rs174594 alone, nominally
261 significant associations with these three genes were found in 29, 28 and 4 tissues, respectively.
262 Importantly, out of these tissues with association signals, the adaptive allele in recent European
263 history is associated with higher expression of *FADS1*, lower expression of *FADS2* and higher
264 expression of *FADS3* in 28, 27 and 4 tissues, respectively. The general trend that recently
265 adaptive allele is associated with higher expression of *FADS1* but lower expression of *FADS2*
266 was also observed for other representative SNPs (rs174546, rs174570, and rs2851682).

267 GWAS have revealed 178 association signals with 44 different traits in the *FADS1-FADS2* LD
268 block, as recorded in the GWAS catalog (Supplementary Tables S5-S9)²⁸. All effects reported in
269 the following are based on individuals of European ancestry, while some are also replicated in
270 other ethnic groups. We report the direction of association in terms of recently adaptive alleles,
271 while the direction is opposite for adaptive alleles in pre-Neolithic hunter-gatherers. Dissecting
272 different associations, (1) the most prominent group of associated traits are polyunsaturated fatty
273 acids (PUFAs, Supplementary Fig. S1), including LCPUFAs and their shorter-chain precursors.
274 Alleles on haplotype D are associated with higher levels of arachidonic acid (AA)²⁹⁻³¹, adrenic
275 acid (AdrA)^{29,31-33}, eicosapentaenoic acid (EPA)^{31,34} and docosapentaenoic acid (DPA)^{31,32,34}, but
276 with lower levels of dihomo-gamma-linolenic acid (DGLA)²⁹⁻³², all of which suggest increased
277 activity of delta-5 desaturase encoded by *FADS1*^{31,35}. This is consistent with the association of
278 recently adaptive alleles with higher *FADS1* expression. Surprisingly, these alleles are associated
279 with higher levels of gamma-linolenic acid (GLA)^{29,30,32} and stearidonic acid (SDA)³¹, but with
280 lower levels of linoleic acid (LA)^{29,30,32,36} and alpha-linolenic acid (ALA)^{30,32,34}, suggesting
281 increased activity of delta-6 desaturase encoded by *FADS2*³⁰. However, the above eQTL analysis
282 suggested that recently adaptive alleles tend to be associated with lower *FADS2* expression.
283 Some of these association signals have been replicated across Europeans^{29,31-36}, Africans³⁴, East
284 Asians^{30,34}, and Hispanic/Latino³⁴. (2) Besides PUFAs, recently adaptive alleles are associated
285 with decreased cis/trans-18:2 fatty acids³⁷, which in turn is associated with lower risks for
286 systemic inflammation and cardiac death³⁷. Consistently, these alleles are also associated with
287 decreased resting heart rate^{38,39}, which reduces risks of cardiovascular disease and mortality. (3)

288 With regards to other lipid levels, recently adaptive alleles have been associated with higher
289 levels of high-density lipoprotein cholesterol (HDL)⁴⁰⁻⁴⁵, low-density lipoprotein cholesterol
290 (LDL)^{40-42,46} and total cholesterol⁴⁰⁻⁴², but with lower levels of triglycerides^{40,41,44,45}. (4) In terms
291 of direct association with disease risk, these alleles are associated with lower risk of
292 inflammatory bowel diseases, both Crohn's disease⁴⁷⁻⁴⁹ and ulcerative colitis⁴⁹, and of bipolar
293 disorder⁵⁰.

294 Going beyond known associations from the GWAS catalog, we analyzed data from the two
295 sequencing cohorts of the UK10K study. Focusing on the peak SNP rs174594, we confirmed the
296 association of the recently adaptive allele with higher levels of TC, LDL, and HDL. We further
297 revealed its association with higher levels of additional lipids, Apo A1 and Apo B
298 (Supplementary Fig. 28). Taken together, recently adaptive alleles, beyond their direct
299 association with fatty acid levels, are associated with factors that are mostly protective against
300 inflammatory and cardiovascular diseases, and indeed show direct association with decreased
301 risk of inflammatory bowel diseases.

302 Discussion

303 Recent positive selection on *FADS* genes after the Neolithic revolution in Europe has been
304 previously reported¹¹. Here, we provided a more detailed view of this recent selection and
305 revealed that it varied geographically, between the North and the South (Figs. 1-3). We further
306 discovered a unique phenomenon that before the Neolithic revolution, the same variants were
307 also subject to positive selection, but with the opposite alleles being selected (Fig. 4). We
308 showed that alleles diminishing LCPUFAs biosynthesis were adaptive before the Neolithic
309 revolution, while alleles enhancing LCPUFAs biosynthesis were adaptive after the Neolithic
310 revolution. In Supplementary Notes, we provided detailed discussions of our results, including 1)
311 interpreting results from different selection tests, especially considering the complications of
312 selection on alternative alleles in two historic periods and selection on standing variations in
313 recent history; 2) interpreting results concerning South-North differences, including
314 consideration of potential geographical differences in demographic history; 3) interpreting
315 eQTLs and GWAS results. Here, we focus instead on interpreting the selection patterns in light
316 of anthropological findings.

317 The dispersal of the Neolithic package into Europe about 8,500 ya caused a sharp dietary shift
318 from an animal-based diet with significant aquatic contribution to a terrestrial plant-heavy diet
319 including dairy products¹⁵⁻²⁰. For pre-Neolithic European hunter-gatherers, the significant role of
320 aquatic food, either marine or freshwater, has been established in sites along the Atlantic
321 coast^{17,51-53}, around the Baltic sea¹⁷, and along the Danube river⁵⁴. The content of LCPUFAs are
322 usually the highest in aquatic foods, lower in animal meat and milk, and almost negligible in
323 most plants⁵⁵. Consistent with the subsistence strategy and dietary pattern, positive selection on
324 *FADS* genes in pre-Neolithic hunter-gatherers was on alleles associated with less efficient
325 LCPUFAs biosynthesis, possibly compensating for the high dietary input. In addition to
326 obtaining sufficient amounts of LCPUFAs, maintaining a balanced ratio of omega-6 to omega-3
327 is critical for human health⁵⁶. Hence, it is also plausible that positive selection in hunter-gatherers
328 was in response to an unbalanced omega-6 to omega-3 ratio (e.g. too much omega-3 LCPUFAs).
329 Positive selection on *FADS* genes has also been observed in modern Greenlandic Inuit, who
330 subsist on a seafood diet⁹. Specifically, the derived allele of rs174570 exhibits positive selection

331 signals in both pre-Neolithic European hunter-gatherers and extant Greenlandic Inuit. More
332 generally, haplotype M2, the candidate adaptive haplotype in pre-Neolithic Europe, is also
333 common in the extant Eskimo samples examined in our study. It is noteworthy that aquatic food
334 was less prevalent among pre-Neolithic hunter-gatherers around the Mediterranean basin,
335 possibly due to the low productivity of the Mediterranean Sea⁵⁷⁻⁵⁹. It would be interesting to
336 examine the geographical differences of selection in pre-Neolithic Europe. However, pre-
337 Neolithic aDNA is still scarce, prohibiting such an analysis at present.

338 The Neolithization of Europe^{12,60,61} started in the Southeast region around 8,500 ya when farming
339 and herding spread into the Aegean and the Balkans. Despite a few temporary stops, it continued
340 spreading into central and northern Europe following the Danube River and its tributaries, and
341 along the Mediterranean coast. It arrived at the Italian Peninsula about 8,000 ya and shortly after
342 reached Iberia by 7,500 ya. While farming rapidly spread across the loess plains of Central
343 Europe and reached the Paris Basin by 7,000 ya, it took another 1,000 or more years before it
344 spread into Britain and Northern Europe around 6,000 ya. From that time on, European farmers
345 relied heavily on their domesticated animals and plants. Compared to pre-Neolithic hunter-
346 gatherers, European farmers consumed much more plants and less aquatic foods^{18-20,62}.
347 Consistent with the lack of LCPUFAs in plant-based diets, positive selection on *FADS* genes
348 during recent European history has been on alleles associated with enhanced LCPUFAs
349 biosynthesis from plant-derived precursors (LA and ALA). Positive selection for enhanced
350 LCPUFAs synthesis has also been observed previously in Africans, South Asians and some East
351 Asians, possibly driven by their traditional plant-based diets^{5,6,8}.

352 Despite the overall trend of relying heavily on domesticated plants, there are geographical
353 differences of dietary patterns among European farmers. In addition to the 2,000-year-late arrival
354 of farming at Northern Europe, animal husbandry and the consumption of animal milk became
355 gradually more prevalent as Neolithic farmers spread to the Northwest^{18,61,63-65}. Moreover,
356 similar to their pre-Neolithic predecessors, Northwestern European farmers close to the Atlantic
357 Ocean or the Baltic Sea still consumed more marine food than their Southern counterparts in the
358 Mediterranean basin^{66,67}. It is noteworthy that historic dairying practice in Northwestern Europe
359 has driven the adaptive evolution of lactase persistence in Europe to reach the highest prevalence
360 in this region⁶⁴. In this study, we observed that recent selection signals for alleles enhancing
361 LCPUFAs biosynthesis are stronger in Southern than in Northern Europeans, even after
362 considering the later arrival of farming and the lower starting allele frequencies in the North. The
363 higher aquatic contribution and stronger reliance on animal meat and milk might be responsible
364 for a weaker selection pressure in the North. However, since GWAS results have unraveled
365 many traits and diseases associated with *FADS* genes, it is possible that other environmental
366 factors beyond diet were involved.

367 **Conclusions**

368 We presented several lines of evidence for positive selection on *FADS* genes in Europe and for
369 its geographically and temporally varying patterns. These patterns concur with mounting
370 anthropological evidence of geographical variability and historical change in dietary patterns.
371 Specifically, in pre-Neolithic hunter-gatherers subsisting on animal-based diets with significant
372 aquatic contribution, LCPUFAs-synthesis-diminishing alleles have been adaptive. In recent
373 European farmers subsisting on plant-heavy diets, LCPUFAs-synthesis-enhancing alleles have
374 been adaptive. Importantly, these are not simply any alleles with opposite functional

375 consequence, but are alternative alleles of the same variants such that when one is under
376 selection and increases in frequency, the other will decrease in frequency. To the best of our
377 knowledge, this is the first example of its kind in humans. Moreover, we reported geographically
378 varying patterns of recent selection that are in line with a stronger dietary reliance on plants in
379 Southern European farmers. These unique, varying patterns of positive selection in different
380 dietary environments, together with the large number of traits and diseases associated with the
381 adaptive region, highlight the importance and potential of matching diet to genome in the future
382 nutritional practice.

383 **Methods**

384 **Data sets.** The ancient DNA (aDNA) data set was compiled from two previous studies^{23,68},
385 which in turn were assembled from many studies, in addition to new sequenced samples. These
386 two data sets were merged by removing overlapping samples. In total, there are 325 ancient
387 samples included in this study. Information about these samples and their original references
388 could be found in Supplementary Table S1. For the aDNA-based test for recent selection, a
389 subset of 178 ancient samples were used and clustered into three groups as in the original
390 study¹¹, representing the three major ancestral sources for most present-day European
391 populations. These three groups are: West and Scandinavian hunter-gatherers (WSHG, N=9),
392 early European farmers (EF, N=76), and individuals of Steppe-pastoralist Ancestry (SA, N=93).
393 Three samples in the EF group in the original study were excluded from our analysis because
394 they are genetic outliers to this group based on additional analysis⁶⁸. For aDNA-based tests for
395 ancient selection in pre-Neolithic European hunter-gatherers, a subset of 42 ancient samples
396 were used and four groups were defined. In addition to the WSHG (N=9), the other three groups
397 were as originally defined in a previous study²³: the “Věstonice cluster”, composed of 14 pre-
398 Last Glacial Maximum individuals from 34,000-26,000 ya; the “El Mirón cluster”, composed of
399 7 post-Last Glacial Maximum individuals from 19,000-14,000 ya; the “Villabruna cluster”,
400 composed of 12 post-Last Glacial Maximum individuals from 14,000-7,000 ya. There were three
401 Western hunter-gatherers that were originally included in the “Villabruna cluster”²³, but we
402 included them in WSHG in the current study because of their similar ages in addition to genetic
403 affinity¹¹. In haplotype network analysis, all aDNAs included in the two aDNA-based selection
404 tests were also included. In addition, we included some well-known ancient samples, such as the
405 Neanderthal, Denisovan, and Ust’-Ishim. In total, there were 225 ancient samples (450
406 haplotypes). For geographical frequency distribution analysis, a total of 300 ancient samples
407 were used and classified into 29 previously defined groups^{11,23,68} based on their genetic affinity,
408 sampling locations and estimated ages.

409 The 1000 Genomes Project (1000GP, phase 3)⁷ has sequencing-based genome-wide SNPs for
410 2,504 individuals from 5 continental regions and 26 global populations. Detailed description of
411 these populations and their sample sizes are in Supplementary Methods. The Human Genome
412 Diversity Project (HGDP)⁶⁹ has genotyping-based genome-wide SNPs for 939 unrelated
413 individuals from 51 populations. The data from the Population Reference Sample (POPRES)⁷⁰
414 were retrieved from dbGaP with permission. Only 3,192 Europeans were included in our
415 analysis. The country of origin of each sample was defined with two approaches. Firstly, a “strict

416 consensus” approach was used: an individual’s country of origin was called if and only if all four
417 of his/her grandparents shared the same country of origin. Secondly, a more inclusive approach
418 was used to further include individuals that had no information about their grandparents. In this
419 case, their countries of birth were used. Both approaches yielded similar results and only results
420 from the inclusive approach are reported. The 22 Eskimo samples were extracted from the
421 Human Origins dataset⁷¹.

422 The two sequencing cohorts of UK10K were obtained from European Genome-phenome
423 Archive with permission⁷². These two cohorts, called ALSPAC and TwinsUK, included low-
424 depth whole-genome sequencing data and a range of quantitative traits for 3,781 British
425 individuals of European ancestry (N=1,927 and 1,854 for ALSPAC and TwinsUK,
426 respectively)⁷².

427 **Imputation for ancient and modern DNA.** Genotype imputation was performed using Beagle
428 4.1⁷³ separately for data sets of aDNA, HGDP and POPRES. The 1000GP phase 3 data were
429 used as the reference panel⁷. Imputation was performed for a 5-Mb region surrounding the *FADS*
430 locus (hg19:chr11: 59,100,000-64,100,000), although most of our analysis was restricted to a 200
431 kb region (hg19:chr11:61,500,000-61,700,000). For most of our analysis (e.g. estimated allele
432 count or frequency for each group), genotype probabilities were taken into account without
433 setting a specific cutoff. For haplotype-based analysis (e.g. estimated haplotype frequency for
434 each group), a cutoff of 0.8 was enforced and haplotypes were defined with missing data and
435 following the phasing information from imputation.

436 Genotype imputation for aDNA has been shown to be desirable and reliable⁷⁴. We also evaluated
437 the imputation quality for aDNA by comparing with the two modern data sets (Supplementary
438 Fig. S29). Overall, the imputation accuracy for ungenotyped SNPs, measured with allelic R^2 and
439 dosage R^2 , is comparable between aDNA and HGDP, but is higher in aDNA when compared
440 with POPRES. Note that sample sizes are much larger for HGDP (N=939) and POPRES
441 (N=3,192), compared to aDNA (N=325). The comparable or even higher imputation quality in
442 aDNA was achieved because of the higher density of genotyped SNPs in the region.

443 **Linkage disequilibrium and haplotype network analysis.** Linkage disequilibrium (LD)
444 analysis was performed with the Haploview software (version 4.2)⁷⁵. Analysis was performed on
445 a 200-kb region (chr11:61,500,000-61,700,000), covering all three *FADS* genes. Variants were
446 included in the analysis if they fulfilled the following criteria: 1) biallelic; 2) minor allele
447 frequency (MAF) in the sample not less than 5%; 3) with rsID; 4) *p* value for Hardy-Weinberg
448 equilibrium test larger than 0.001. Analysis was performed separately for the combined UK10K
449 cohort and each of the five European populations in 1000G.

450 Haplotype network analysis was performed with an R software package, pegas⁷⁶. To reduce the
451 number of SNPs and thus the number of haplotypes included in the analysis, we restricted this
452 analysis to part of the 85 kb *FADS1-FADS2* LD block, starting 5 kb downstream of *FADS1* to
453 the end of the LD block (a 60-kb region). To further reduce the number of SNPs, in the analysis
454 with all 1000GP European samples, we applied an iterative algorithm⁷⁷ to merge haplotypes that
455 have no more than three nucleotide differences by removing the differing SNPs. The algorithm

456 stops when all remaining haplotypes are more than 3 nucleotides away. With this procedure, we
457 were able to reduce the number of total haplotypes from 81 to 12, with the number of SNPs
458 decreased from 88 to 34 (Supplementary Fig. S30). This set of 34 representative SNPs was used
459 in all haplotype-based analysis in aDNA, 1000GP, HGDP and POPRES. Missing data (e.g. from
460 a low imputation genotype probability) were included in the haplotype network analysis.

461 Of note, for the 12 haplotypes identified in 1000GP European samples, only five of them have
462 frequency higher than 1% (Supplementary Table S2). These five haplotypes were designated as
463 D, M1, M2, M3 and M4, from the most common to the least.

464 **Ancient DNA-based test for recent selection in Europe.** The test was performed as described
465 before¹¹. Briefly, most European populations could be modelled as a mixture of three ancient
466 source populations at fixed proportions. The three ancient source populations are West or
467 Scandinavian hunter-gatherers (WSHG), early European farmers (EF), and Steppe-Ancestry
468 pastoralist (SA) (Supplementary Table S1). For modern European populations in 1000G, the
469 proportions of these three ancestral sources estimated at genome-wide level are (0.196, 0.257,
470 0.547) for CEU, (0.362, 0.229, 0.409) for GBR, (0, 0.686, 0.314) for IBS, and (0, 0.645, 0.355)
471 for TSI. FIN was not used because it does not fit this three-population model¹¹. Under neutrality,
472 the frequencies of a SNP (e.g. reference allele) in present-day European populations are expected
473 to be the linear combination of its frequencies in the three ancient source populations. This
474 serves as the null hypothesis: $p_{mod} = Cp_{anc}$, where p_{mod} is the frequencies in A modern
475 populations, p_{anc} is the frequencies in B ancient source populations while C is an AxB matrix
476 with each row representing the estimated ancestral proportions for one modern population. The
477 alternative hypothesis is that p_{mod} is unconstrained by p_{anc} . The frequency in each population is
478 modelled with binomial distribution: $L(p; D) = B(X, 2N, p)$, where X is the number of
479 designated allele observed while N is the sample size. In ancient populations, X is the expected
480 number of designated allele observed, taking into account uncertainty in imputation. We write
481 $\ell(p; D)$ for the log-likelihood. The log-likelihood for SNP frequencies in all three ancient
482 populations and four modern populations are: $\ell(\vec{p}; \vec{D}) = \sum_{i=1}^A \ell(p_i; D_i) + \sum_{j=1}^B \ell(p_j; D_j)$.
483 Under the null hypothesis, there are B parameters in the model, corresponding to the frequencies
484 in B ancient populations. Under the alternative hypothesis, there are A+B parameters,
485 corresponding to the frequencies in A modern populations and B ancestral populations. We
486 numerically maximized the likelihood separately under each hypothesis and evaluate the statistic
487 (twice the difference in log-likelihood) with the null χ_A^2 distribution. Inflation was observed with
488 this statistic in a previous genome-wide analysis and a $\lambda = 1.38$ was used for correction¹¹.
489 Following this, we applied the same factor in correcting the p values in our analysis. For
490 genotyped SNPs previously tested, similar scales of statistical significance were observed as in
491 the previous study (Supplementary Fig. 31). We note that for the purpose of refining the
492 selection signal with imputed variants, only relative significance levels across variants are
493 informative.

494 In addition to combining signals from four present-day European populations, we further
495 performed tests separately in the two South European populations (IBS and TSI) and in the two
496 North European populations (CEU and GBR). In these two cases, A = 2 and the null distribution

497 is χ_2^2 . For comparison between the North and the South, we used three statistics: the final p
498 value, the maximum likelihood under the null hypothesis, and the maximum likelihood under the
499 alternative hypothesis.

500 **Ancient DNA-based test for ancient selection in pre-Neolithic European hunter-gatherers.**
501 Two Bayesian methods, the Schraiber *et al.* method²⁴ and the ApproxWF²⁵, were applied to infer
502 natural selection from allele frequency time series data. The two software were downloaded from
503 <https://github.com/Schraiber/selection> and <https://bitbucket.org/phaentu/approxwf/downloads/>,
504 respectively. The Schraiber *et al.* method models the evolutionary trajectory of an allele under a
505 specified demographic history and estimates selection coefficients (s_1 and s_2) for heterozygotes
506 and homozygotes of the allele under study. This method has two modes, with or without the
507 simultaneous estimation of allele age. Without the estimation of allele age, this method models
508 the frequency trajectory only between the first and last time points provided and its estimates of
509 selection coefficients describe the selection force during this period only. With the simultaneous
510 estimation of allele age, this method models the frequency trajectory starting from the first
511 appearance of the allele to the last time point provided. In this case, the selection coefficients
512 describe the selection force starting from the mutation of the allele, which therefore should be the
513 derived allele. For demographic history, we used two models: a constant population size model
514 with $N_e=10,000$ and a more realistic model with two historic epochs of bottleneck and recent
515 exponential growth²⁶. However, the recent epoch of exponential growth does not have an impact
516 on our analysis because for our analysis the most recent sample, WSHG, has an age estimate of
517 ~ 7500 years, predating the onset of exponential growth (3520 ya, assuming 25 years per
518 generation). ApproxWF can simultaneously estimate selection coefficient and demographic
519 history (only for constant population size model). For our purpose, we set the demographic
520 history as $N_e=10,000$. It estimates selection coefficient for homozygotes, s , and dominance
521 coefficient, h . The selection coefficient estimated is for the time points specified by the input
522 data.

523 Four groups of pre-Neolithic European hunter-gatherers were included in our test: the Věstonice
524 cluster (median sample age: 30,076 yo), the El Mirón cluster (14,959 yo), the Villabruna cluster
525 (10,059 yo) and WSHG (7,769 yo). To identify SNPs with evidence of positive selection during
526 the historic period from Věstonice to WSHG, we applied both methods on most SNPs in the
527 *FADS* locus. The Schraiber *et al.* method was run twice with two demographic models while
528 ApproxWF was run once with the constant size model. For the two candidate SNPs (rs174570
529 and rs2851682), we further ran the Schraiber *et al.* method with the more realistic demographic
530 model to simultaneously estimate their selection coefficients and allele ages. Statistical
531 significance was considered if the 95% CI of selection coefficient does not overlap with 0.
532 Details about running the two software were in Supplementary Methods.

533 **Modern DNA-based selection tests.** We performed two types of selection tests for modern
534 DNAs: site frequency spectrum (SFS)-based and haplotype-based tests. These tests were
535 performed separately in each of the five European populations from 1000G and each of the two
536 cohorts from UK10K. For SFS-based tests, we calculated genetic diversity (π), Tajima's D^{78} , and
537 Fay and Wu's H^{79} , using in-house Perl scripts. We calculated these three statistics with a sliding-

538 window approach (window size = 5 kb and moving step = 1 kb). Statistical significance for these
539 statistics were assessed using the genome-wide empirical distribution. Haplotype-based tests,
540 including iHS⁸⁰ and nSL²¹, were calculated using software selscan (version 1.1.0a)⁸¹. Only
541 common biallelic variants (MAF > 5%) were included in the analysis. Genetic variants without
542 ancestral information were excluded. These two statistics were normalized in frequency bins (1%
543 interval) and the statistical significance of the normalized iHS and nSL were evaluated with the
544 empirical genome-wide distribution. The haplotype bifurcation diagrams and EHH decay plots
545 were drawn using an R package, rehh⁸². Singleton Density Score (SDS) based on UK10K was
546 directly retrieved from a previous study²².

547 **Geographical frequency distribution analysis.** For plots of geographical frequency
548 distribution, the geographical map was plotted with an R software package, maps
549 (<https://CRAN.R-project.org/package=maps>) while the pie charts were added with the mapplots
550 package (<https://cran.r-project.org/web/packages/mapplots/index.html>). Haplotype frequencies
551 were calculated based on haplotype network analysis with pegas⁷⁶, which groups haplotypes
552 while taking into account missing data. SNP frequencies were either the observed frequency, if
553 the SNP was genotyped, or the expected frequency based on genotype probability, if the SNP
554 was imputed.

555 **Targeted association analysis for SNP rs174594 in UK10K.** We performed association
556 analysis for rs174594 in two UK10K datasets – ALSPAC and TwinsUK⁷². For both datasets, we
557 analyzed height, weight, BMI and lipid-related traits including total cholesterol, low density
558 lipoprotein, very low density lipoprotein, high density lipoprotein, Apolipoprotein A-I (APOA1),
559 Apolipoprotein B (APOB) and triglyceride. We performed principal components analysis using
560 smartpca from EIGENSTRAT software⁸³ with genome-wide autosomal SNPs and we added top
561 4 principal components as covariates for all association analysis. We also used age as a covariate
562 for all association analysis. Sex was added as a covariate only for ALSPAC dataset since all
563 individuals in TwinsUK dataset are female. For all lipid-related traits, we also added BMI as a
564 covariate.

565 **Data availability.**

566 Ancient DNA: <https://reich.hms.harvard.edu/datasets>
567 1000 Genomes Project: <ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/>
568 Human Genome Diversity Project (HGDP): <http://www.hagsc.org/hgdp/files.html>
569 Population Reference Sample (POPRES): dbGaP Study Accession: phs000145.v4.p2
570 UK10K: https://www.uk10k.org/data_access.html
571 Singleton Density Score (SDS): <https://github.com/yairf/SDS>

572 **Code availability.** Most analyses were conducted with available software and packages as
573 described in the respective subsections of Methods. Customized Perl and R scripts were used in
574 performing site frequency spectrum-based selection tests, and for general plotting purposes. All
575 these scripts are available upon request.

576 References

- 577
- 578 1 Fan, S., Hansen, M. E. B., Lo, Y. & Tishkoff, S. A. Going global by adapting local: A
579 review of recent human adaptation. *Science* **354**, 54-59 (2016).
 - 580 2 Nakamura, M. T. & Nara, T. Y. Structure, function, and dietary regulation of delta6,
581 delta5, and delta9 desaturases. *Annu Rev Nutr* **24**, 345-376 (2004).
 - 582 3 Raphael, W. & Sordillo, L. M. Dietary polyunsaturated fatty acids and inflammation: the
583 role of phospholipid biosynthesis. *Int J Mol Sci* **14**, 21167-21188 (2013).
 - 584 4 Bazinet, R. P. & Laye, S. Polyunsaturated fatty acids and their metabolites in brain
585 function and disease. *Nat Rev Neurosci* **15**, 771-785 (2014).
 - 586 5 Mathias, R. A. *et al.* Adaptive evolution of the FADS gene cluster within Africa. *PLoS*
587 *One* **7**, e44926 (2012).
 - 588 6 Ameer, A. *et al.* Genetic adaptation of fatty-acid metabolism: a human-specific haplotype
589 increasing the biosynthesis of long-chain omega-3 and omega-6 fatty acids. *Am J Hum*
590 *Genet* **90**, 809-820 (2012).
 - 591 7 The 1000 Genomes Project Consortium *et al.* A global reference for human genetic
592 variation. *Nature* **526**, 68-74 (2015).
 - 593 8 Kothapalli, K. S. *et al.* Positive Selection on a Regulatory Insertion-Deletion
594 Polymorphism in FADS2 Influences Apparent Endogenous Synthesis of Arachidonic
595 Acid. *Mol Biol Evol* **33**, 1726-1739 (2016).
 - 596 9 Fumagalli, M. *et al.* Greenlandic Inuit show genetic signatures of diet and climate
597 adaptation. *Science* **349**, 1343-1347 (2015).
 - 598 10 Reardon, H. T. *et al.* Insertion-deletions in a FADS2 intron 1 conserved regulatory locus
599 control expression of fatty acid desaturases 1 and 2 and modulate response to simvastatin.
600 *Prostaglandins Leukot Essent Fatty Acids* **87**, 25-33 (2012).
 - 601 11 Mathieson, I. *et al.* Genome-wide patterns of selection in 230 ancient Eurasians. *Nature*
602 **528**, 499-503 (2015).
 - 603 12 Bar-Yosef, O. in *On Human Nature: Biology, Psychology, Ethics, Politics, and Religion*
604 (eds M. Tibayrenc & F. J. Ayala) Ch. 19, 297-331 (Academic Press, 2017).
 - 605 13 Coward, F., Shennan, S., Colledge, S., Conolly, J. & Collard, M. The spread of Neolithic
606 plant economies from the Near East to northwest Europe: a phylogenetic analysis.
607 *Journal of Archaeological Science* **35**, 42-56 (2008).
 - 608 14 Bogaard, A. *et al.* Crop manuring and intensive land management by Europe's first
609 farmers. *Proc Natl Acad Sci U S A* **110**, 12589-12594 (2013).
 - 610 15 Richards, M. P. in *The Evolution of Hominin Diets: Integrating Approaches to the Study*
611 *of Palaeolithic Subsistence* (eds J. J. Hublin & M. P. Richards) 251-257 (Springer
612 Science; Business Media, 2009).
 - 613 16 Richards, M. P., Schulting, R. J. & Hedges, R. E. Archaeology: sharp shift in diet at onset
614 of Neolithic. *Nature* **425**, 366 (2003).
 - 615 17 Richards, M. P., Price, T. D. & Koch, E. Mesolithic and Neolithic Subsistence in
616 Denmark: New Stable Isotope Data. *Current Anthropology* **44**, 288-295 (2003).
 - 617 18 Fraser, R. A., Bogaard, A., Schäfer, M., Arbogast, R. & Heaton, T. H. E. Integrating
618 botanical, faunal and human stable carbon and nitrogen isotope values to reconstruct land
619 use and palaeodiet at LBK Vaihingen an der Enz, Baden-Württemberg. *World*
620 *Archaeology* **45**, 492-517 (2013).

- 621 19 Knipper, C. *et al.* What is on the menu in a Celtic town? Iron Age diet reconstructed at
622 Basel-Gasfabrik, Switzerland. *Archaeological and Anthropological Sciences* (2016).
- 623 20 López-Costas, O., Müldner, G. & Martínez Cortizas, A. Diet and lifestyle in Bronze Age
624 Northwest Spain: the collective burial of Cova do Santo. *Journal of Archaeological*
625 *Science* **55**, 209-218 (2015).
- 626 21 Ferrer-Admetlla, A., Liang, M., Korneliussen, T. & Nielsen, R. On detecting incomplete
627 soft or hard selective sweeps using haplotype structure. *Mol Biol Evol* **31**, 1275-1291
628 (2014).
- 629 22 Field, Y. *et al.* Detection of human adaptation during the past 2000 years. *Science* **354**,
630 760-764 (2016).
- 631 23 Fu, Q. *et al.* The genetic history of Ice Age Europe. *Nature* **534**, 200-205 (2016).
- 632 24 Schraiber, J. G., Evans, S. N. & Slatkin, M. Bayesian Inference of Natural Selection from
633 Allele Frequency Time Series. *Genetics* **203**, 493-511 (2016).
- 634 25 Ferrer-Admetlla, A., Leuenberger, C., Jensen, J. D. & Wegmann, D. An Approximate
635 Markov Model for the Wright-Fisher Diffusion and Its Application to Time Series Data.
636 *Genetics* **203**, 831-846 (2016).
- 637 26 Gazave, E. *et al.* Neutral genomic regions refine models of recent rapid human
638 population growth. *Proc Natl Acad Sci U S A* **111**, 757-762 (2014).
- 639 27 GTEx Consortium. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue
640 gene regulation in humans. *Science* **348**, 648-660 (2015).
- 641 28 Welter, D. *et al.* The NHGRI GWAS Catalog, a curated resource of SNP-trait
642 associations. *Nucleic Acids Res* **42**, D1001-1006 (2014).
- 643 29 Guan, W. *et al.* Genome-wide association study of plasma N6 polyunsaturated fatty acids
644 within the cohorts for heart and aging research in genomic epidemiology consortium.
645 *Circ Cardiovasc Genet* **7**, 321-331 (2014).
- 646 30 Dorajoo, R. *et al.* A genome-wide association study of n-3 and n-6 plasma fatty acids in a
647 Singaporean Chinese population. *Genes Nutr* **10**, 53 (2015).
- 648 31 Shin, S. Y. *et al.* An atlas of genetic influences on human blood metabolites. *Nat Genet*
649 **46**, 543-550 (2014).
- 650 32 Tintle, N. L. *et al.* A genome-wide association study of saturated, mono- and
651 polyunsaturated red blood cell fatty acids in the Framingham Heart Offspring Study.
652 *Prostaglandins Leukot Essent Fatty Acids* **94**, 65-72 (2015).
- 653 33 Xie, W. *et al.* Genetic variants associated with glycine metabolism and their role in
654 insulin sensitivity and type 2 diabetes. *Diabetes* **62**, 2141-2150 (2013).
- 655 34 Lemaitre, R. N. *et al.* Genetic loci associated with plasma phospholipid n-3 fatty acids: a
656 meta-analysis of genome-wide association studies from the CHARGE Consortium. *PLoS*
657 *Genet* **7**, e1002193 (2011).
- 658 35 Gieger, C. *et al.* Genetics meets metabolomics: a genome-wide association study of
659 metabolite profiles in human serum. *PLoS Genet* **4**, e1000282 (2008).
- 660 36 Kettunen, J. *et al.* Genome-wide association study identifies multiple loci influencing
661 human serum metabolite levels. *Nat Genet* **44**, 269-276 (2012).
- 662 37 Mozaffarian, D. *et al.* Genetic loci associated with circulating phospholipid trans fatty
663 acids: a meta-analysis of genome-wide association studies from the CHARGE
664 Consortium. *Am J Clin Nutr* **101**, 398-406 (2015).
- 665 38 Eijgelsheim, M. *et al.* Genome-wide association analysis identifies multiple loci related
666 to resting heart rate. *Hum Mol Genet* **19**, 3885-3894 (2010).

- 667 39 den Hoed, M. *et al.* Identification of heart rate-associated loci and their effects on cardiac
668 conduction and rhythm disorders. *Nat Genet* **45**, 621-631 (2013).
- 669 40 Global Lipids Genetics Consortium *et al.* Discovery and refinement of loci associated
670 with lipid levels. *Nat Genet* **45**, 1274-1283 (2013).
- 671 41 Teslovich, T. M. *et al.* Biological, clinical and population relevance of 95 loci for blood
672 lipids. *Nature* **466**, 707-713 (2010).
- 673 42 Aulchenko, Y. S. *et al.* Loci influencing lipid levels and coronary heart disease risk in 16
674 European population cohorts. *Nat Genet* **41**, 47-55 (2009).
- 675 43 Zabaneh, D. & Balding, D. J. A genome-wide association study of the metabolic
676 syndrome in Indian Asian men. *PLoS One* **5**, e11961 (2010).
- 677 44 Kathiresan, S. *et al.* Common variants at 30 loci contribute to polygenic dyslipidemia.
678 *Nat Genet* **41**, 56-65 (2009).
- 679 45 Waterworth, D. M. *et al.* Genetic variants influencing circulating lipid levels and risk of
680 coronary artery disease. *Arterioscler Thromb Vasc Biol* **30**, 2264-2276 (2010).
- 681 46 Sabatti, C. *et al.* Genome-wide association analysis of metabolic traits in a birth cohort
682 from a founder population. *Nat Genet* **41**, 35-46 (2009).
- 683 47 Franke, A. *et al.* Genome-wide meta-analysis increases to 71 the number of confirmed
684 Crohn's disease susceptibility loci. *Nat Genet* **42**, 1118-1125 (2010).
- 685 48 Liu, J. Z. *et al.* Association analyses identify 38 susceptibility loci for inflammatory
686 bowel disease and highlight shared genetic risk across populations. *Nat Genet* **47**, 979-
687 986 (2015).
- 688 49 Jostins, L. *et al.* Host-microbe interactions have shaped the genetic architecture of
689 inflammatory bowel disease. *Nature* **491**, 119-124 (2012).
- 690 50 Ikeda, M. *et al.* A genome-wide association study identifies two novel susceptibility loci
691 and trans population polygenicity associated with bipolar disorder. *Mol Psychiatry*
692 (2017).
- 693 51 Richards, M. P. & Hedges, R. E. M. Stable Isotope Evidence for Similarities in the Types
694 of Marine Foods Used by Late Mesolithic Humans at Sites Along the Atlantic Coast of
695 Europe. *Journal of Archaeological Science* **26**, 717-722 (1999).
- 696 52 Lubell, D., Jackes, M., Schwarcz, H. & Knyf, M. The Mesolithic-Neolithic Transition in
697 Portugal: Isotopic and Dental Evidence of Diet. *Journal of Archaeological Science* **21**,
698 201-216 (1994).
- 699 53 Richards, M. P. & Mellars, P. A. Stable isotopes and the seasonality of the Oronsay
700 middens. *Antiquity* **72**, 178-184 (1998).
- 701 54 Bonsall, C. *et al.* Mesolithic and Early Neolithic in the Iron Gates: A Palaeodietary
702 Perspective. *Journal of European Archaeology* **5**, 50-92 (1997).
- 703 55 Abedi, E. & Sahari, M. A. Long-chain polyunsaturated fatty acid sources and evaluation
704 of their nutritional and functional properties. *Food Sci Nutr* **2**, 443-463 (2014).
- 705 56 Simopoulos, A. P. Evolutionary aspects of diet: the omega-6/omega-3 ratio and the brain.
706 *Mol Neurobiol* **44**, 203-215 (2011).
- 707 57 Mannino, M. A., Thomas, K. D., Leng, M. J., Di Salvo, R. & Richards, M. P. Stuck to the
708 shore? Investigating prehistoric hunter-gatherer subsistence, mobility and territoriality in
709 a Mediterranean coastal landscape through isotope analyses on marine mollusc shell
710 carbonates and human bone collagen. *Quaternary International* **244**, 88-104 (2011).
- 711 58 Mannino, M. A. *et al.* Origin and diet of the prehistoric hunter-gatherers on the
712 mediterranean island of Favignana (Egadi Islands, Sicily). *PLoS One* **7**, e49802 (2012).

- 713 59 Lightfoot, E., Boneva, B., Miracle, P. T., Šlaus, M. & O'Connell, T. C. Exploring the
714 Mesolithic and Neolithic transition in Croatia through isotopic investigations. *Antiquity*
715 **85**, 73-86 (2015).
- 716 60 Bocquet-Appel, J.-P., Naji, S., Vander Linden, M. & Kozłowski, J. Understanding the
717 rates of expansion of the farming system in Europe. *Journal of Archaeological Science*
718 **39**, 531-546 (2012).
- 719 61 Rowley-Conwy, P. Westward Ho! The Spread of Agriculture from Central Europe to the
720 Atlantic. *Current Anthropology* **52**, S431-S451 (2011).
- 721 62 Vigne, J.-D. in *The Neolithic Demographic Transition and its Consequences* (eds J.-P.
722 Bocquet-Appel & O. Bar-Yosef) 179-205 (Springer Science+Business Media B.V.,
723 2008).
- 724 63 Cramp, L. J. *et al.* Immediate replacement of fishing with dairying by the earliest farmers
725 of the Northeast Atlantic archipelagos. *Proc Biol Sci* **281**, 20132372 (2014).
- 726 64 Curry, A. Archaeology: The milk revolution. *Nature* **500**, 20-22 (2013).
- 727 65 Salque, M. *et al.* Earliest evidence for cheese making in the sixth millennium BC in
728 northern Europe. *Nature* **493**, 522-525 (2013).
- 729 66 Lidén, K., Eriksson, G., Nordqvist, B., Götherström, A. & Bendixen, E. “The wet and the
730 wild followed by the dry and the tame” – or did they occur at the same time? Diet in
731 Mesolithic – Neolithic southern Sweden. *Antiquity* **78**, 23-33 (2004).
- 732 67 Rottoli, M. & Castiglioni, E. Prehistory of plant growing and collecting in northern Italy,
733 based on seed remains from the early Neolithic to the Chalcolithic (c. 5600–2100 cal
734 b.c.). *Vegetation History and Archaeobotany* **18**, 91-103 (2008).
- 735 68 Lazaridis, I. *et al.* Genomic insights into the origin of farming in the ancient Near East.
736 *Nature* **536**, 419-424 (2016).
- 737 69 Li, J. Z. *et al.* Worldwide human relationships inferred from genome-wide patterns of
738 variation. *Science* **319**, 1100-1104 (2008).
- 739 70 Nelson, M. R. *et al.* The Population Reference Sample, POPRES: a resource for
740 population, disease, and pharmacological genetics research. *Am J Hum Genet* **83**, 347-
741 358 (2008).
- 742 71 Lazaridis, I. *et al.* Ancient human genomes suggest three ancestral populations for
743 present-day Europeans. *Nature* **513**, 409-413 (2014).
- 744 72 The UK10 Consortium *et al.* The UK10K project identifies rare variants in health and
745 disease. *Nature* **526**, 82-90 (2015).
- 746 73 Browning, B. L. & Browning, S. R. Genotype Imputation with Millions of Reference
747 Samples. *Am J Hum Genet* **98**, 116-126 (2016).
- 748 74 Gamba, C. *et al.* Genome flux and stasis in a five millennium transect of European
749 prehistory. *Nat Commun* **5**, 5257 (2014).
- 750 75 Barrett, J. C., Fry, B., Maller, J. & Daly, M. J. Haploview: analysis and visualization of
751 LD and haplotype maps. *Bioinformatics* **21**, 263-265 (2005).
- 752 76 Paradis, E. pegas: an R package for population genetics with an integrated-modular
753 approach. *Bioinformatics* **26**, 419-420 (2010).
- 754 77 Dannemann, M., Andres, A. M. & Kelso, J. Introgression of Neandertal- and Denisovan-
755 like Haplotypes Contributes to Adaptive Variation in Human Toll-like Receptors. *Am J*
756 *Hum Genet* **98**, 22-33 (2016).
- 757 78 Tajima, F. Statistical method for testing the neutral mutation hypothesis by DNA
758 polymorphism. *Genetics* **123**, 585-595 (1989).

- 759 79 Fay, J. C. & Wu, C. I. Hitchhiking under positive Darwinian selection. *Genetics* **155**,
760 1405-1413 (2000).
- 761 80 Voight, B. F., Kudravalli, S., Wen, X. & Pritchard, J. K. A map of recent positive
762 selection in the human genome. *PLoS Biol* **4**, e72 (2006).
- 763 81 Szpiech, Z. A. & Hernandez, R. D. selscan: an efficient multithreaded program to
764 perform EHH-based scans for positive selection. *Mol Biol Evol* **31**, 2824-2827 (2014).
- 765 82 Gautier, M. & Vitalis, R. rehh: an R package to detect footprints of selection in genome-
766 wide SNP data from haplotype structure. *Bioinformatics* **28**, 1176-1177 (2012).
- 767 83 Price, A. L. *et al.* Principal components analysis corrects for stratification in genome-
768 wide association studies. *Nat Genet* **38**, 904-909 (2006).
- 769 84 Wang, J. *et al.* Factorbook.org: a Wiki-based database for transcription factor-binding
770 data generated by the ENCODE consortium. *Nucleic Acids Res* **41**, D171-176 (2013).

771 **Acknowledgements**

772 We thank Montgomery Slatkin and Joshua Schraiber for their help in running their software,
773 David Reich and Iain Mathieson for making their data publicly available, Leonardo Arbiza,
774 Charles Liang, Daniel (Alex) Marburgh, Edward Li, Kumar Kothapalli, Tom Brenna, and all
775 members of the Keinan lab for helpful discussion and comments on the manuscript. This work
776 was supported by the National Institutes of Health (Grants R01HG006849 and R01GM108805 to
777 AK) and the Edward Mallinckrodt, Jr. Foundation (AK).

778 This study makes use of data generated by the UK10K Consortium, derived from samples from
779 UK10K_COHORT_ALSPAC and UK10K_COHORT_TWINSUK. A full list of the
780 investigators who contributed to the generation of the data is available from www.UK10K.org.
781 Funding for UK10K was provided by the Wellcome Trust under award WT091310.

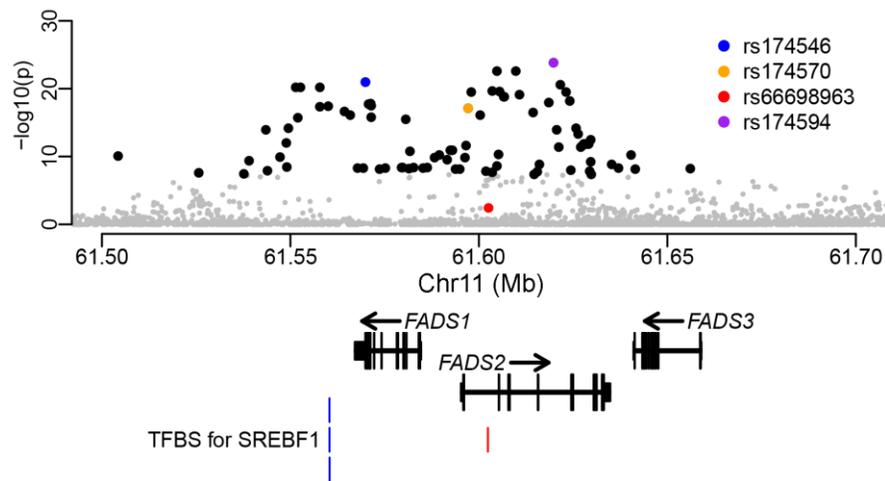
782 The collections and methods for the Population Reference Sample (POPRES) are described by
783 Nelson *et al.* (2008). The datasets used for the analyses described in this manuscript were
784 obtained from dbGaP at [http://www.ncbi.nlm.nih.gov/projects/gap/cgi-](http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000145.v1.p1)
785 [bin/study.cgi?study_id=phs000145.v1.p1](http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000145.v1.p1) through dbGaP accession number phs000145.v1.p1.

786 **Author contributions**

787 A.K. and K.Y. conceived and designed the project. K.Y. performed data collection and analysis,
788 with contributions from D.W. and F.G.. K.Y. and A.K. interpreted the results, with contribution
789 from O.B. on the anthropological perspective. K.Y. and A.K. wrote the manuscript. All authors
790 read, edited and approved the final version of the manuscript.

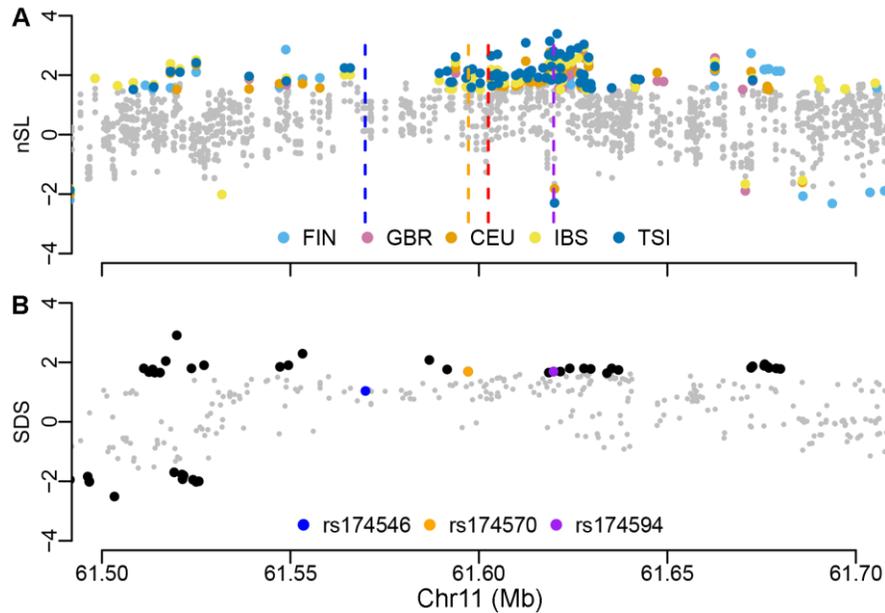
791 **Competing interests**

792 The authors declare no competing interests.



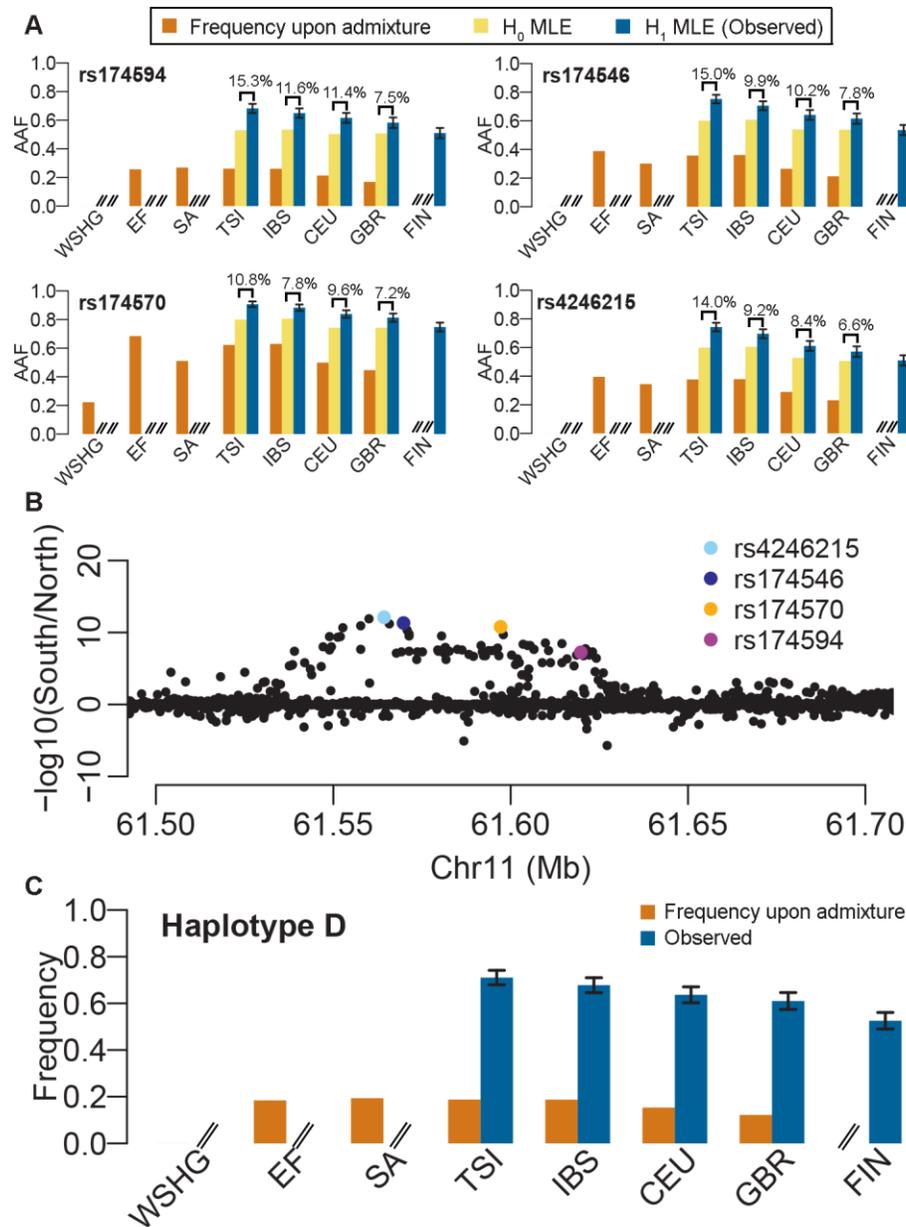
793

794 **Fig. 1. Ancient DNA-based test for recent positive selection.** The y axis indicates genomic
795 control corrected p values at a negative logarithm scale. Variants under genome-wide
796 significance level ($5e-8$) are in gray except for highlighted ones. Four variants are highlighted:
797 the most significant SNP (purple); the top SNP reported by Mathieson *et al.*¹¹ (blue); one of the
798 top adaptive SNPs reported in Greenlandic Inuit⁹ (orange); the indel reported to be targeted by
799 positive selection in populations with historical plant-based diets⁸ (red). The overall pattern is
800 consistent with that previously described¹¹ (Supplementary Fig. S31). At the bottom are the
801 representative transcript models for the three *FADS* genes and the four transcription factor
802 binding sites (TFBS) for SREBF1 from ENCODE⁸⁴ (blue) and another previous study¹⁰ (red).



803

804 **Fig. 2. Tests for recent positive selection based solely on modern DNA.** (A) Haplotype-based
805 selection test (nSL^{21}) in modern Europeans from 1000GP. The test was performed separately for
806 each of the five European groups. Only variants with significant values are shown with
807 population-specific colors as indicated in the legend. The positions for four variants of interest
808 were indicated with vertical dashed lines, colored as in Fig. 1. For presentation purpose, the sign
809 was set so that being positive indicates that the adaptive allele revealed by nSL is consistent with
810 that revealed by the aDNA-based test in Fig. 1. Original statistics for 1000GP and UK10K are
811 shown in Supplementary Figs. S5 and S7. The five 1000GP European populations are: CEU –
812 Utah Residents (CEPH) with Northern and Western Ancestry; FIN – Finnish in Finland; GBR –
813 British in England and Scotland; IBS – Iberian Population in Spain; TSI – Toscani in Italia. (B)
814 Singleton Density Score (SDS^{22}) in modern Europeans from UK10K. Variants under
815 significance level are in gray except for highlighted ones. Three variants of interest were
816 highlighted with colors as indicated in the legend. The indel rs66698963 was not present in the
817 original UK10K data set. The sign of SDS was set as in nSL. Original statistics are shown in
818 Supplementary Fig. S8.

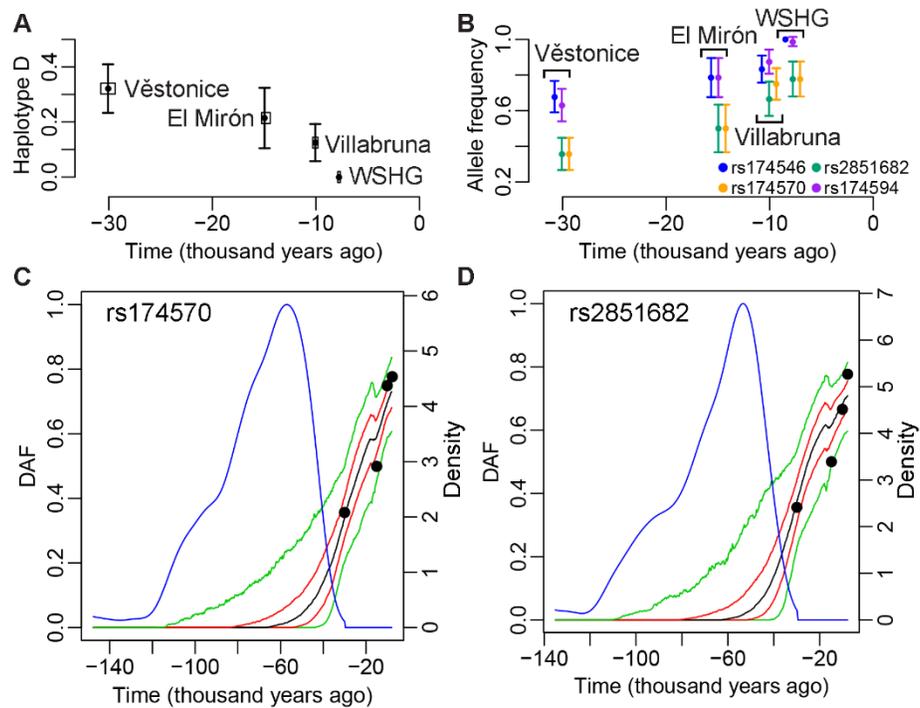


819

820 **Fig. 3. Varying selection and frequency patterns between Southern and Northern Europe.**

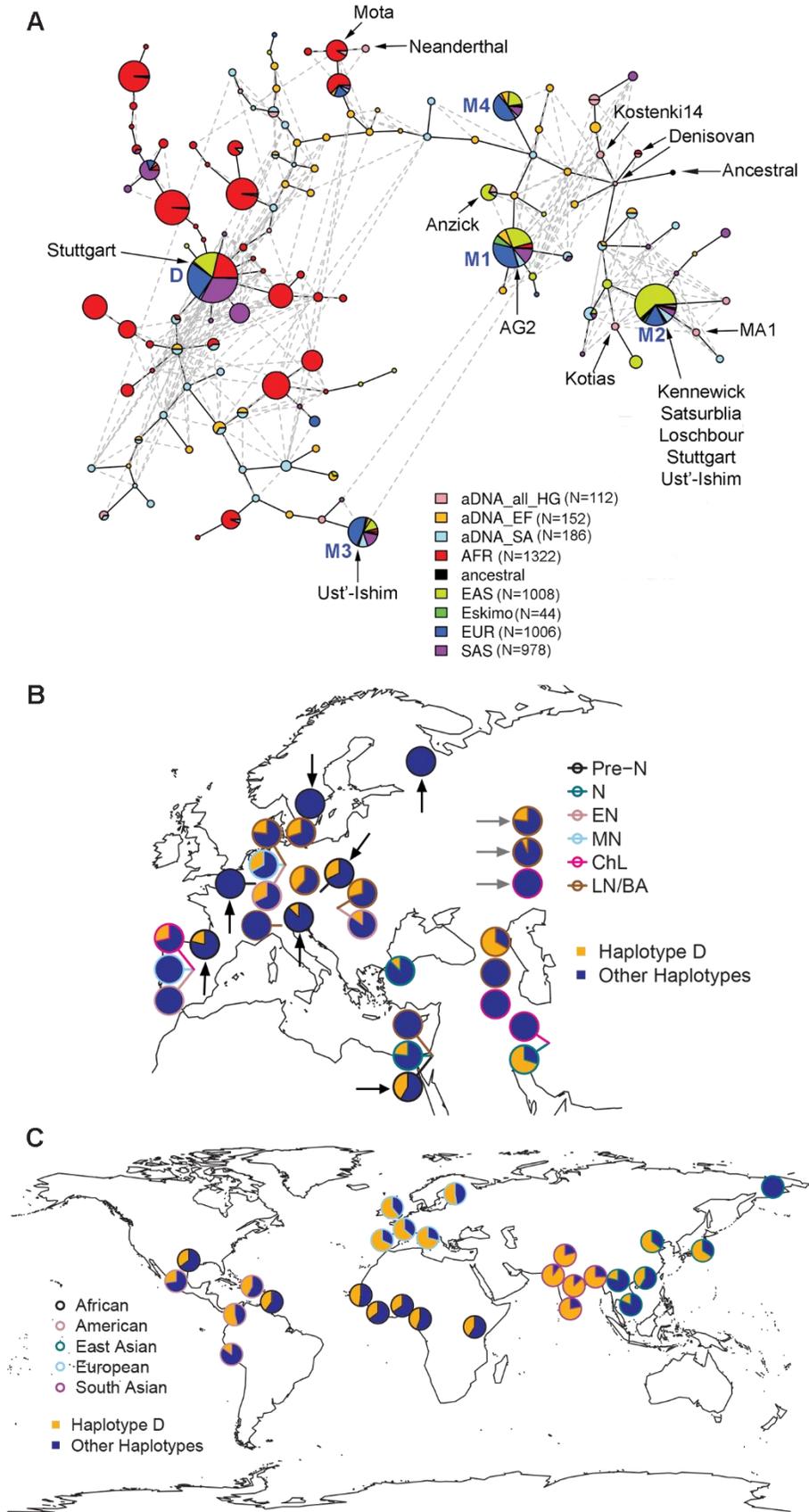
821 (A) South-North frequency gradient for adaptive alleles of four representative SNPs under
 822 different scenarios. AAF refers to adaptive allele frequency. Orange bars represent frequencies
 823 upon admixture, which were directly observed in ancient groups and predicted for extant
 824 populations based on linear mixture of frequencies in ancient groups. Yellow bars represent
 825 frequencies estimated under H_0 . Estimates for ancient groups were not shown because they are
 826 not relevant here. Blue bars represent frequencies estimated under H_1 , whose only constraint is
 827 the observed data and therefore the MLEs are just the observed means. The estimates for ancient
 828 groups are the same as their frequencies upon admixture and are omitted on the plot. The
 829 absolute difference between H_0 and H_1 estimates are indicated above the corresponding bars.
 830 Please note that the frequencies upon admixture in WSHG are 0 for rs174594, rs174546 and

831 rs4246215 and no bars were plotted. **(B)** Comparison of aDNA-based selection signals between
832 Southern and Northern Europe. aDNA-based selection tests were performed separately for
833 Southern (TSI and IBS) and Northern (CEU and GBR) Europeans. For each variant, the p values
834 from these two tests were compared at a $-\log_{10}$ scale (y axis). SNPs of interest were colored as
835 indicated. **(C)** South-North frequency gradient for the adaptive haplotype in extant populations.
836 The two frequency types are just as in (A). The frequency upon admixture for WSHG is 0. In (A)
837 and (C), FIN has only observed values. If values are not shown or not available, signs of “//” are
838 indicated at corresponding positions. Error bars stand for standard errors.

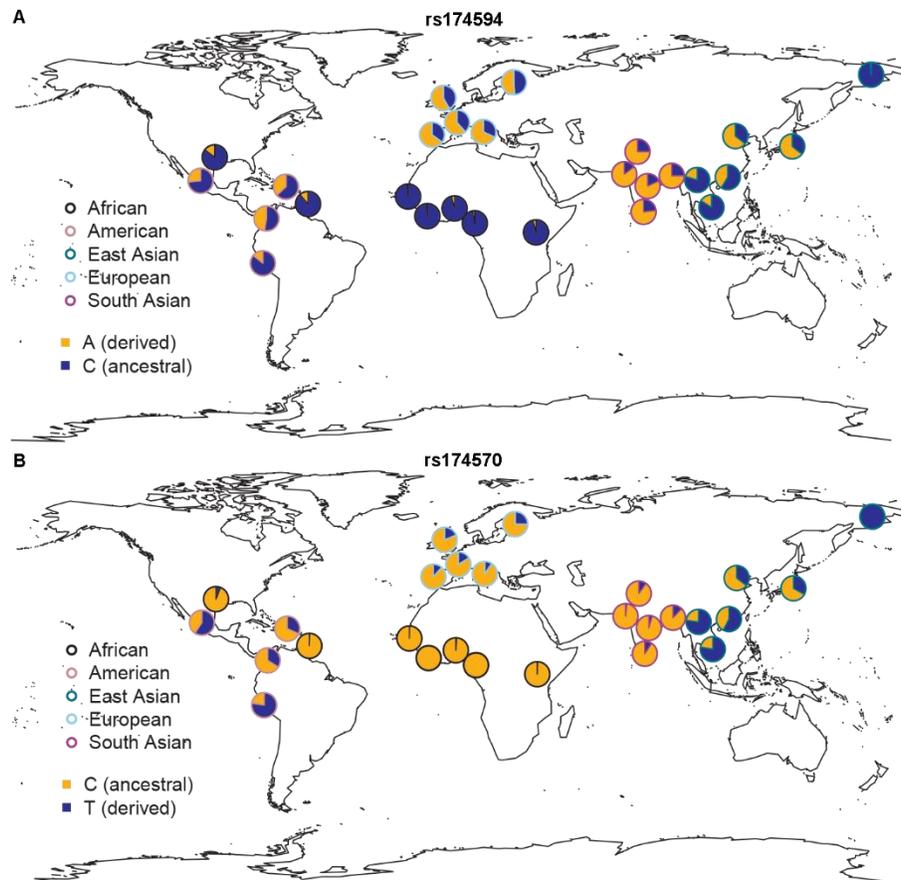


839

840 **Fig. 4. Temporal frequency pattern and selection signals in pre-Neolithic European hunter-**
841 **gatherers. (A)** The frequency of haplotype D over time in four groups of hunter-gatherers.
842 Frequency for each group is plotted as a black point at the median age of samples. The horizontal
843 box surrounding the point represents the medians of lower- and upper-bound estimates of sample
844 ages. Error bars are standard errors. Group names are indicated next to their frequencies. The
845 frequency for WSHG is 0. **(B)** Allele frequencies for four SNPs. It has similar format as in (A)
846 except that small arbitrary values were added on their x coordinates in order to visualize all
847 SNPs, which were colored as indicated in the legend. The alleles chosen are the ones increasing
848 frequency over time. They are derived alleles for rs174570 and rs2851682, and ancestral alleles
849 for rs174546 and rs174594. **(C)** and **(D)** Posterior distribution on the derived allele frequency
850 path for rs174570 and rs2851682, respectively. The sampled frequencies are indicated with black
851 points, which are the same point estimates as in (B). The median, 25% and 75% quantiles, and
852 5% and 95% quantiles of the posterior distribution are indicated respectively with black, red and
853 green lines. The posterior distribution on the age of derived allele is shown with a blue line, with
854 values on the right y axis.



856 **Fig. 5. Haplotype network and geographical frequency distribution.** (A) Haplotype network
857 for 1000G samples (2,157 individuals, excluding admixed American samples), 22 modern
858 Eskimos and 225 aDNAs. Each pie chart represents one haplotype and its size is proportional to
859 $\log_2(\# \text{ of haplotype})$ plus a minimum size to visualize rare haplotypes. Sections in the pie provide
860 the breakdown by groups. Detailed haplotype frequencies are in Supplementary Table S2. The
861 edges connecting haplotypes are of arbitrary length. Haplotypes for some well-known ancient
862 samples are labelled. The top five haplotypes in modern Europeans, referred to as D, M1, M2,
863 M3, and M4 from the most to least frequent common, are indicated with their names in blue. (B)
864 Frequency of haplotype D in Eurasian ancient DNAs. Each pie represents one sampled group
865 and is placed at the sampling location or nearby with a line pointing at the sampling location.
866 The color of the pie chart border indicates the archaeological period. If multiple samples of
867 different periods were collected at the same geographical location, these samples are ordered
868 vertically with the older samples at the bottom. Hunter-gatherer groups are indicated with black
869 arrows and pastoralist groups with gray arrows, while others are farmers. Geographical locations
870 for some hunter-gatherer groups (e.g. the Věstonice, El Mirón and Villabruna clusters) are only
871 from representative samples. Detailed frequencies are in Supplementary Table S3. Pre-N: Pre-
872 Neolithic; N: Neolithic; EN: Early Neolithic; MN: Mid-Neolithic; ChL: Chalcolithic; LN/BA:
873 Late Neolithic/Bronze Age. (C) Frequency of haplotype D in present-day global populations. All
874 26 populations from 1000GP and one Eskimo group are included. The color of the pie chart
875 border represents the genetic ancestry. It is noteworthy that there are two samples in America
876 that are actually of African ancestry. Detailed frequencies are in Supplementary Table S4.



877

878 **Fig. 6. Geographical frequency distribution for SNPs rs174594 and rs174570 in present-day**
879 **global populations.** Adaptive alleles in recent European history are colored in orange. All 26
880 populations from 1000GP and one Eskimo group are included. The color of the pie chart border
881 represents the genetic ancestry. It is noteworthy that there are two samples in America that are
882 actually of African ancestry. Similar global patterns were observed with HGDP samples
883 (Supplementary Figs. S19 and S21).