



Published in final edited form as:

Nature. 2015 December 24; 528(7583): 499–503. doi:10.1038/nature16152.

## Genome-wide patterns of selection in 230 ancient Eurasians

Iain Mathieson<sup>1</sup>, Iosif Lazaridis<sup>1,2</sup>, Nadin Rohland<sup>1,2</sup>, Swapan Mallick<sup>1,2,3</sup>, Nick Patterson<sup>2</sup>, Songül Alpaslan Roodenberg<sup>4</sup>, Eadaoin Harney<sup>1,3</sup>, Kristin Stewardson<sup>1,3</sup>, Daniel Fernandes<sup>5</sup>, Mario Novak<sup>5,6</sup>, Kendra Sirak<sup>5,7</sup>, Cristina Gamba<sup>5,8,9</sup>, Eppie R. Jones<sup>9</sup>, Bastien Llamas<sup>10</sup>, Stanislav Dryomov<sup>11,12</sup>, Joseph Pickrel<sup>11,13</sup>, Juan Luís Arsuaga<sup>14,15</sup>, José María Bermúdez de Castro<sup>16</sup>, Eudald Carbonell<sup>17,18</sup>, Fokke Gerritsen<sup>19</sup>, Aleksandr Khokhlov<sup>20</sup>, Pavel Kuznetsov<sup>20</sup>, Marina Lozano<sup>17,18</sup>, Harald Meller<sup>21</sup>, Oleg Mochalov<sup>20</sup>, Vayacheslav Moiseyev<sup>22</sup>, Manuel A. Rojo Guerra<sup>23</sup>, Jacob Roodenberg<sup>24</sup>, Josep Maria Vergès<sup>17,18</sup>, Johannes Krause<sup>25,26</sup>, Alan Cooper<sup>10</sup>, Kurt W. Alt<sup>21,27,28,29</sup>, Dorcas Brown<sup>30</sup>, David Anthony<sup>30</sup>, Carles Lalueza-Fox<sup>31</sup>, Wolfgang Haak<sup>10,25,+</sup>, Ron Pinhasi<sup>5,+</sup>, and David Reich<sup>1,2,3,+</sup>

<sup>1</sup>Department of Genetics, Harvard Medical School, Boston, Massachusetts 02115, USA <sup>2</sup>Broad Institute of MIT and Harvard, Cambridge Massachusetts 02142, USA <sup>3</sup>Howard Hughes Medical Institute, Harvard Medical School, Boston, Massachusetts 02115, USA <sup>4</sup>Independent researcher, Santpoort-Noord, The Netherlands <sup>5</sup>School of Archaeology and Earth Institute, Belfield, University College Dublin, Dublin 4, Ireland <sup>6</sup>Institute for Anthropological Research, Zagreb 10000, Croatia <sup>7</sup>Department of Anthropology, Emory University, Atlanta, Georgia, USA <sup>8</sup>Smurfit Institute of Genetics, Trinity College Dublin, Dublin 2, Ireland <sup>9</sup>Australian Centre for Ancient DNA, School of Earth and Environmental Sciences & Environment Institute, University of Adelaide, Adelaide, South Australia 5005, Australia <sup>10</sup>Laboratory of Human Molecular Genetics, Institute of Molecular and Cellular Biology, Siberian Branch of the Russian Academy of Sciences, Novosibirsk 630090, Russia <sup>11</sup>Department of Paleolithic Archaeology, Institute of Archaeology and Ethnography, Siberian Branch of the Russian Academy of Sciences, Novosibirsk 630090, Russia <sup>12</sup>Centro Mixto UCM-ISCIll de Evolución y Comportamiento Humanos, Madrid, Spain <sup>13</sup>Departamento de Paleontología, Facultad Ciencias Geológicas, Universidad Complutense de

Reprints and permissions information are available at [www.nature.com/reprints](http://www.nature.com/reprints).

Correspondence and requests for materials should be addressed to I.M. (iain\_mathieson@hms.harvard.edu), W.H. (haak@shh.mpg.de), R.P. (ron.pinhasi@ucd.ie) or D.R. (reich@genetics.med.harvard.edu).

<sup>+</sup>These authors contributed equally to this work

<sup>8</sup>Current address: Centre for GeoGenetics, Natural History Museum of Denmark, University of Copenhagen, Øster Voldgade 5–7, 1350 Copenhagen, Denmark

<sup>13</sup>Current Address: New York Genome Center, New York NY, USA

The aligned sequences are available through the European Nucleotide Archive under accession number PRJEB11450 The Human Origins genotype datasets including ancient individuals can be found at ([http://genetics.med.harvard.edu/reichlab/Reich\\_Lab/Datasets.html](http://genetics.med.harvard.edu/reichlab/Reich_Lab/Datasets.html)).

The authors declare no competing financial interests.

Readers are welcome to comment on the online version of the paper.

### Author Contributions

W.H., R.P. and D.R. supervised the study. S.A.R., J.L.A., J.M.B., E.C., F.G., A.K., P.K., M.L., H.M., O.M., V.M., M.A.R., J.R., J.M.V., J.K., A.C., K.W.A., D.B., D.A., C.L., W.H., R.P. and D.R. assembled archaeological material. I.M., I.L., N.R., S.M., N.P., S.D., J.P., W.H. and D.R. analysed genetic data. N.R., E.H., K.S., D.F., M.N., K.S., C.G., E.R.J., B.L., C.L. and W.H. performed wet laboratory ancient DNA work. I.M., I.L. and D.R. wrote the manuscript with input from all co-authors.

Madrid, Spain <sup>16</sup>Centro Nacional de Investigación sobre Evolución Humana (CENIEH), 09002 Burgos, Spain <sup>17</sup>IPHES. Institut Català de Paleoecologia Humana i Evolució Social, Campus Sescelades-URV, 43007. Tarragona, Spain <sup>18</sup>Area de Prehistoria, Universitat Rovira i Virgili (URV), 43002 Tarragona, Spain <sup>19</sup>Netherlands Institute in Turkey, Istiklal Caddesi, Nur-i Ziya Sokak 5, Beyo lu, Istanbul, Turkey <sup>20</sup>Volga State Academy of Social Sciences and Humanities, Samara 443099, Russia <sup>21</sup>State Office for Heritage Management and Archaeology Saxony-Anhalt and State Museum of Prehistory, D-06114 Halle, Germany <sup>22</sup>Peter the Great Museum of Anthropology and Ethnography (Kunstkamera) RAS, St Petersburg, 199034, Russia <sup>23</sup>Department of Prehistory and Archaeology, University of Valladolid, Spain <sup>24</sup>The Netherlands Institute for the Near East, Leiden, RA-2300, The Netherlands <sup>25</sup>Max Planck Institute for the Science of Human History, D-07745 Jena, Germany <sup>26</sup>Institute for Archaeological Sciences, University of Tübingen, D-72070 Tübingen, Germany <sup>27</sup>Danube Private University, A-3500 Krems, Austria <sup>28</sup>Institute for Prehistory and Archaeological Science, University of Basel, CH-4003 Basel, Switzerland <sup>29</sup>Institute of Anthropology, Johannes Gutenberg University of Mainz, D-55128 Mainz, Germany <sup>30</sup>Anthropology Department, Hartwick College, Oneonta, New York 13820, USA <sup>31</sup>Institute of Evolutionary Biology (CSIC-Universitat Pompeu Fabra), Barcelona, Spain

## Abstract

Ancient DNA makes it possible to directly witness natural selection by analyzing samples from populations before, during and after adaptation events. Here we report the first scan for selection using ancient DNA, capitalizing on the largest genome-wide dataset yet assembled: 230 West Eurasians dating to between 6500 and 1000 BCE, including 163 with newly reported data. The new samples include the first genome-wide data from the Anatolian Neolithic culture whose genetic material we extracted from the DNA-rich petrous bone and who we show were members of the population that was the source of Europe's first farmers. We also report a complete transect of the steppe region in Samara between 5500 and 1200 BCE that allows us to recognize admixture from at least two external sources into steppe populations during this period. We detect selection at loci associated with diet, pigmentation and immunity, and two independent episodes of selection on height.

## Introduction

The arrival of farming in Europe around 8,500 years ago necessitated adaptation to new environments, pathogens, diets, and social organizations. While indirect evidence of this adaptation can be detected in patterns of genetic variation in present-day people<sup>1</sup>, these patterns are only echoes of past events, which are difficult to date and interpret, and are often confounded by neutral processes. Ancient DNA provides a more direct view, and should be a transformative technology for studies of selection just as it has transformed studies of history. Until now, however, the large sample sizes required to detect selection have meant that ancient DNA studies have concentrated on characterizing effects at parts of the genome already believed to have been affected by selection<sup>2–5</sup>.

## Genome-wide ancient DNA from West Eurasia

We assembled genome-wide data from 230 ancient individuals who lived in West Eurasia from 6500 to 1000 BCE (Fig. 1a, Extended Data Table 1, Supplementary Data Table 1, Supplementary Information section 1). To obtain this dataset, we combined published data from 67 samples from relevant periods and cultures<sup>4–6</sup>, with 163 samples for which we report new data, of which 83 have never previously been analyzed (the remaining 80 samples include 67 whose targeted single nucleotide polymorphism (SNP) coverage we triple from 390k to 1240k<sup>7</sup>; and 13 with shotgun data whose data quality we increase using our enrichment strategy<sup>3,8</sup>). The 163 samples for which we report new data are drawn from 270 distinct individuals who we screened for evidence of authentic DNA<sup>7</sup>. We used in-solution hybridization with synthesized oligonucleotide probes to enrich promising libraries for more than 1.2 million SNPs (“1240k capture”, Methods). The targeted sites include nearly all SNPs on the Affymetrix Human Origins and Illumina 610-Quad arrays, 49,711 SNPs on chromosome X and 32,681 on chromosome Y, and 47,384 SNPs with evidence of functional importance. We merged libraries from the same individual and filtered out samples with low coverage or evidence of contamination to obtain the final set of individuals. The advantage of 1240k capture is that it gives access to genome-wide data from ancient samples with small fractions of human DNA and increases efficiency by targeting sites in the human genome that will actually be analyzed. The effectiveness of the approach can be seen by comparing our results to the largest previously published ancient DNA study, which used a shotgun sequencing strategy<sup>5</sup>. Our median coverage on analyzed SNPs is ~4-times higher even while the mean number of reads generated per sample is 36-times lower (Extended Data Fig. 1).

## Insight into population transformations

To learn about the history of archaeological cultures for which genome-wide data is reported for the first time here, we studied either 1,055,209 autosomal SNPs when analyzing 230 ancient individuals alone, or 592,169 SNPs when co-analyzing them with 2,345 present-day individuals genotyped on the Human Origins array<sup>4</sup>. We removed 13 samples either as outliers in ancestry relative to others of the same archaeologically determined culture, or first-degree relatives (Supplementary Data Table 1).

Our sample of 26 Anatolian Neolithic individuals represents the first genome-wide ancient DNA data from the eastern Mediterranean. Our success at analyzing such a large number of samples is likely due to the fact that at the Barcın site—the source of 21 of the working samples—we sampled from the cochlea of the petrous bone<sup>9</sup>, which has been shown to increase the amount of DNA obtained by up to two orders of magnitude relative to teeth (the next-most-promising tissue)<sup>3</sup>. Principal component (PCA) and ADMIXTURE<sup>10</sup> analysis, shows that the Anatolian Neolithic samples do not resemble any present-day Near Eastern populations but are shifted towards Europe, clustering with Neolithic European farmers (EEF) from Germany, Hungary, and Spain<sup>7</sup> (Fig. 1b, Extended Data Fig. 2). Further evidence that the Anatolian Neolithic and EEF were related comes from the high frequency (47%; n=15) of Y-chromosome haplogroup G2a typical of ancient EEF samples<sup>7</sup> (Supplementary Data Table 1), and the low fixation index ( $F_{ST}$ ; 0.005–0.016) between

Neolithic Anatolians and EEF (Supplementary Data Table 2). These results support the hypothesis<sup>7</sup> of a common ancestral population of EEF prior to their dispersal along distinct inland/central European and coastal/Mediterranean routes. The EEF are slightly more shifted to Europe in the PCA than are the Anatolian Neolithic (Fig. 1b) and have significantly more admixture from Western hunter-gatherers (WHG), shown by  $f_4$ -statistics ( $|Z| > 6$  standard errors from 0) and negative  $f_3$ -statistics ( $|Z| > 4$ )<sup>11</sup> (Extended Data Table 2). We estimate that the EEF have 7–11% more WHG admixture than their Anatolian relatives (Extended Data Fig. 2, Supplementary Information section 2).

The Iberian Chalcolithic individuals from El Mirador cave are genetically similar to the Middle Neolithic Iberians who preceded them (Fig. 1b; Extended Data Fig. 2), and have more WHG ancestry than their Early Neolithic predecessors<sup>7</sup> ( $|Z| > 10$ ) (Extended Data Table 2). However, they do not have a significantly different proportion of WHG ancestry (we estimate 23–28%) than the Middle Neolithic Iberians (Extended Data Fig. 2). Chalcolithic Iberians have no evidence of steppe ancestry (Fig. 1b, Extended Data Fig. 2), in contrast to central Europeans of the same period<sup>5,7</sup>. Thus, the “Ancient North Eurasian”-related ancestry that is ubiquitous across present-day Europe<sup>4,7</sup> arrived in Iberia later than in Central Europe (Supplementary Information section 2).

To understand population transformations in the Eurasian steppe, we analyzed a time transect of 37 samples from the Samara region spanning ~5600–300 BCE and including the Eastern Hunter-gatherer (EHG), Eneolithic, Yamnaya, Poltavka, Potapovka and Srubnaya cultures. Admixture between populations of Near Eastern ancestry and the EHG<sup>7</sup> began as early as the Eneolithic (5200–4000 BCE), with some individuals resembling EHG and some resembling Yamnaya (Fig. 1b; Extended Data Fig. 2). The Yamnaya from Samara and Kalmykia, the Afanasievo people from the Altai (3300–3000 BCE), and the Poltavka Middle Bronze Age (2900–2200 BCE) population that followed the Yamnaya in Samara, are all genetically homogeneous, forming a tight “Bronze Age steppe” cluster in PCA (Fig. 1b), sharing predominantly R1b Y-chromosomes<sup>5,7</sup> (Supplementary Data Table 1), and having 48–58% ancestry from an Armenian-like Near Eastern source (Extended Data Table 2) without additional Anatolian Neolithic or Early European Farmer (EEF) ancestry<sup>7</sup> (Extended Data Fig. 2). After the Poltavka period, population change occurred in Samara: the Late Bronze Age Srubnaya have ~17% Anatolian Neolithic or EEF ancestry (Extended Data Fig. 2). Previous work documented that such ancestry appeared east of the Urals beginning at least by the time of the Sintashta culture, and suggested that it reflected an eastward migration from the Corded Ware peoples of central Europe<sup>5</sup>. However, the fact that the Srubnaya also harbored such ancestry indicates that the Anatolian Neolithic or EEF ancestry could have come into the steppe from a more eastern source. Further evidence that migrations originating as far west as central Europe may not have had an important impact on the Late Bronze Age steppe comes from the fact that the Srubnaya possess exclusively ( $n=6$ ) R1a Y-chromosomes (Supplementary Data Table 1), and four of them (and one Poltavka male) belonged to haplogroup R1a-Z93 which is common in central/south Asians<sup>12</sup>, very rare in present-day Europeans, and absent in all ancient central Europeans studied to date.

## Twelve signals of selection

To study selection, we created a dataset of 1,084,781 autosomal SNPs in 617 samples by merging 213 ancient samples with genome-wide sequencing data from four populations of European ancestry from the 1,000 Genomes Project<sup>13</sup>. Most present-day Europeans can be modeled as a mixture of three ancient populations related to Mesolithic hunter-gatherers (WHG), early farmers (EEF) and steppe pastoralists (Yamnaya)<sup>4,7</sup>, and so to scan for selection, we divided our samples into three groups based on which of these populations they clustered with most closely (Fig. 1b, Extended Data Table 1). We estimated mixture proportions for the present-day European ancestry populations and tested every SNP to evaluate whether its present-day frequencies were consistent with this model. We corrected for test statistic inflation by applying a genomic control correction analogous to that used to correct for population structure in genome-wide association studies<sup>14</sup>. Of ~1 million non-monomorphic autosomal SNPs, the ~50,000 in the set of potentially functional SNPs were significantly more inconsistent with the model than neutral SNPs (Fig. 2), suggesting pervasive selection on polymorphisms of functional importance. Using a conservative significance threshold of  $p=5.0\times 10^{-8}$ , and a genomic control correction of 1.38, we identified 12 loci that contained at least three SNPs achieving genome-wide significance within 1 Mb of the most associated SNP (Fig. 2, Extended Data Table 3, Extended Data Fig. 3, Supplementary Data Table 3).

The strongest signal of selection is at the SNP (rs4988235) responsible for lactase persistence in Europe<sup>15,16</sup>. Our data (Fig. 3) strengthens previous reports that an appreciable frequency of lactase persistence in Europe only dates to the last four thousand years<sup>3,5,17</sup>. The allele's earliest appearance in our data is in a central European Bell Beaker sample (individual I0112) that lived between approximately 2300 and 2200 BCE. Two other independent signals related to diet are located on chromosome 11 near *FADS1* and *DHCR7*. *FADS1* and *FADS2* are involved in fatty acid metabolism, and variation at this locus is associated with plasma lipid and fatty acid concentration<sup>18</sup>. The selected allele of the most significant SNP (rs174546) is associated with decreased triglyceride levels<sup>18</sup>. This locus has experienced independent selection in non-European populations<sup>13,19,20</sup> and is likely to be a critical component of adaptation to different diets. Variants at *DHCR7* and *NADSYN1* are associated with circulating vitamin D levels<sup>21</sup> and the most associated SNP in our analysis, rs7940244, is highly differentiated across closely related Northern European populations<sup>22,23</sup>, suggesting selection related to variation in dietary or environmental sources of vitamin D.

Two signals have a potential link to celiac disease. One occurs at the ergothioneine transporter *SLC22A4* that is hypothesized to have experienced a selective sweep to protect against ergothioneine deficiency in agricultural diets<sup>24</sup>. Common variants at this locus are associated with increased risk for ulcerative colitis, celiac disease, and irritable bowel disease and may have hitchhiked to high frequency as a result of this sweep<sup>24–26</sup>. However the specific variant (rs1050152, L503F) that was thought to be the target did not reach high frequency until relatively recently (Extended Data Fig. 4). The signal at *ATXN2/SH2B3*—also associated with celiac disease<sup>25</sup>—shows a similar pattern (Extended Data Fig. 4).

The second strongest signal in our analysis is at the derived allele of rs16891982 in *SLC45A2*, which contributes to light skin pigmentation and is almost fixed in present-day Europeans but occurred at much lower frequency in ancient populations. In contrast, the derived allele of *SLC24A5* that is the other major determinant of light skin pigmentation in modern Europe appears fixed in the Anatolian Neolithic, suggesting that its rapid increase in frequency to around 0.9 in Early Neolithic Europe was mostly due to migration (Extended Data Fig. 4). Another pigmentation signal is at *GRM5*, where SNPs are associated with pigmentation possibly through a regulatory effect on nearby *TYR*<sup>27</sup>. We also find evidence of selection for the derived allele of rs12913832 at *HERC2/OCA2*, which appears to be fixed in Mesolithic hunter-gatherers, and is the primary determinant of blue eye color in present-day Europeans<sup>28,29</sup>. In contrast to the other loci, the range of frequencies in modern populations is within that of ancient populations (Fig. 3). The frequency increases with higher latitude, suggesting a complex pattern of environmental selection.

The *TLR1-TLR6-TLR10* gene cluster is a known target of selection in Europe, possibly related to resistance to leprosy, tuberculosis or other mycobacteria<sup>30–32</sup>. There is also a strong signal of selection at the major histocompatibility complex (MHC) on chromosome 6. The strongest signal is at rs2269424 near the genes *PPT2* and *EGFL8* but there are at least six other apparently independent signals in the MHC (Extended Data Fig. 3); and the entire region is significantly more associated than the genome-wide average (residual inflation of 2.07 in the region on chromosome 6 between 29–34 Mb after genome-wide genomic control correction). This could be the result of multiple sweeps, balancing selection, or increased drift due to background selection reducing effective population size in this gene-rich region.

We find a surprise in six Scandinavian hunter-gatherers (SHG) from the Motala site in southern Sweden. In three out of six samples, we observe the haplotype carrying the derived allele of rs3827760 in the *EDAR* gene (Extended Data Fig. 5), which affects tooth morphology and hair thickness<sup>33,34</sup>, has been the subject of a selective sweep in East Asia<sup>35</sup>, and today is at high frequency in East Asians and Native Americans. The *EDAR* derived allele is largely absent in present-day Europe except in Scandinavia, plausibly due to Siberian movements into the region millennia after the date of the Motala samples. The SHG have no evidence of East Asian ancestry<sup>4,7</sup>, suggesting that the *EDAR* derived allele may not have originated not in East Asians as previously suggested<sup>35</sup>. A second surprise is that, unlike closely related western hunter-gatherers, the Motala samples have predominantly derived pigmentation alleles at *SLC45A2* and *SLC24A5*.

## Evidence of selection on height

We also tested for selection on complex traits. The best-documented example of this process in humans is height, for which the differences between Northern and Southern Europe have driven by selection<sup>36</sup>. To test for this signal in our data, we used a statistic that tests whether trait-affecting alleles are both highly correlated and more differentiated, compared to randomly sampled alleles<sup>37</sup>. We predicted genetic heights for each population and applied the test to all populations together, as well as to pairs of populations (Fig. 4). Using 180 height-associated SNPs<sup>38</sup> (restricted to 169 where we successfully targeted at least two chromosomes in each population), we detect a significant signal of directional selection on



height ( $p=0.002$ ). Applying this to pairs of populations allows us to detect two independent signals. First, the Iberian Neolithic and Chalcolithic samples show selection for reduced height relative to both the Anatolian Neolithic ( $p=0.042$ ) and the Central European Early and Middle Neolithic ( $p=0.003$ ). Second, we detect a signal for increased height in the steppe populations ( $p=0.030$  relative to the Central European Early and Middle Neolithic). These results suggest that the modern South-North gradient in height across Europe is due to both increased steppe ancestry in northern populations, and selection for decreased height in Early Neolithic migrants to southern Europe. We did not observe any other significant signals of polygenetic selection in five other complex traits we tested: body mass index<sup>39</sup> ( $p=0.20$ ), waist-to-hip ratio<sup>40</sup> ( $p=0.51$ ), type 2 diabetes<sup>41</sup> ( $p=0.37$ ), inflammatory bowel disease<sup>26</sup> ( $p=0.17$ ) and lipid levels<sup>18</sup> ( $p=0.50$ ).

## Future studies of selection with ancient DNA

Our results, which take advantage of the massive increase in sample size enabled by optimized techniques for sampling from the petrous bone as well as in-solution enrichment methods for targeted sets of SNPs, show how ancient DNA can be used to perform a genome-wide scan for selection, and demonstrate selection on loci related to pigmentation, diet and immunity, painting a picture of Neolithic populations adapting to settled agricultural life at high latitudes. For most of the signals we detect, allele frequencies of modern Europeans are outside the range of any ancient populations, indicating that phenotypically, Europeans of four thousand years ago were different in important respects from Europeans today despite having overall similar ancestry. An important direction for future research is to increase the sample size for European selection scans (Extended Data Fig. 6), and to apply this approach to regions beyond Europe and to nonhuman species.

## Methods

### Ancient DNA analysis

We screened 433 next generation sequencing libraries from 270 distinct samples for authentic ancient DNA using previously reported protocols<sup>7</sup>. All libraries that we included in nuclear genome analysis were treated with uracil-DNA-glycosylase (UDG) to reduce characteristic errors of ancient DNA<sup>42</sup>.

We performed in-solution enrichment for a targeted set of 1,237,207 SNPs using previously reported protocols<sup>4,7,43</sup>. The targeted SNP set merges 394,577 SNPs first reported in Ref. 7 (390k capture), and 842,630 SNPs first reported in ref.<sup>44</sup> (840k capture). For 67 samples for which we newly report data in this study, there was pre-existing 390k capture data<sup>7</sup>. For these samples, we only performed 840k capture and merged the resulting sequences with previously generated 390k data. For the remaining samples, we pooled the 390k and 840k reagents together to produce a single enrichment reagent. We attempted to sequence each enriched library up to the point where we estimated that it was economically inefficient to sequence further. Specifically, we iteratively sequenced more and more from each sample and only stopped when we estimated that the expected increase in the number of targeted SNPs hit at least once would be less than one for every 100 new read pairs generated. After sequencing, we filtered out samples with <30,000 targeted SNPs covered at least once, with

evidence of contamination based on mitochondrial DNA polymorphism<sup>43</sup>, an appreciable rate of heterozygosity on chromosome X despite being male<sup>45</sup>, or an atypical ratio of X to Y sequences.

Of the targeted SNPs, 47,384 are “potentially functional” sites chosen as follows (with some overlap): 1,290 SNPs identified as targets of selection in Europeans by the Composite of Multiple Signals (CMS) test<sup>1</sup>; 21,723 SNPs identified as significant hits by genome-wide association studies, or with known phenotypic effect (GWAS); 1,289 SNPs with extremely differentiated frequencies between HapMap populations<sup>46</sup> (HiDiff); 9,116 immunochip SNPs chosen for study of immune phenotypes (Immune); 347 SNPs phenotypically relevant to South America (mostly altitude adaptation SNPs in *EGLN1* and *EPAS1*), 5,387 SNPs which tag HLA haplotypes and 13,672 expression quantitative trait loci<sup>47</sup> (eQTL).

### Population history analysis

We used two datasets for population history analysis. “*HO*” consists of 592,169 SNPs, taking the intersection of the SNP targets and the Human Origins SNP array<sup>4</sup>; we used this dataset for co-analysis of present-day and ancient samples. “*HOIII*” consists of 1,055,209 SNPs that additionally includes sites from the Illumina genotype array<sup>48</sup>; we used this dataset for analyses only involving the ancient samples.

On the *HO* dataset, we carried out principal components analysis in *smartpca*<sup>49</sup> using a set of 777 West Eurasian individuals<sup>4</sup>, and projected the ancient individuals with the option “*lsqproject: YES*”. We carried out ADMIXTURE analysis on a set of 2,345 present-day individuals and the ancient samples after pruning for LD in PLINK 1.9 (<https://www.cog-genomics.org/plink2>)<sup>50</sup> with parameters “*-indep-pairwise 200 25 0.4*”. We varied the number of ancestral populations between  $K=2$  and  $K=20$ , and used cross-validation (*-cv*) to identify the value of  $K=17$  to plot in Extended Data Fig. 2f.

We used ADMIXTOOLS<sup>11</sup> to compute  $f$ -statistics, determining standard errors with a Block Jackknife and default parameters. We used the option “*inbreed: YES*” when computing  $f_3$ -statistics of the form  $f_3(\textit{Ancient}, \textit{Ref}_1, \textit{Ref}_2)$  as the *Ancient* samples are represented by randomly sampled alleles rather than by diploid genotypes. For the same reason, we estimated  $F_{ST}$  genetic distances between populations on the *HO* dataset with at least two individuals in *smartpca* also using the “*inbreed: YES*” option.

We estimated ancestral proportions as in Supplementary Information section 9 of Ref. 7, using a method that fits mixture proportions on a *Test* population as a mixture of  $N$  *Reference* populations by using  $f_4$ -statistics of the form  $f_4(\textit{Test} \textit{or} \textit{Ref}, O_1; O_2, O_3)$  that exploit allele frequency correlations of the *Test* or *Reference* populations with triples of *Outgroup* populations. We used a set of 15 world outgroup populations<sup>4,7</sup>. In Extended Data Fig. 2, we added WHG and EHG as outgroups for those analyses in which they are not used as reference populations.

We determined sex by examining the ratio of aligned reads to the sex chromosomes<sup>51</sup>. We assigned Y-chromosome haplogroups to males using version 9.1.129 of the nomenclature of the International Society of Genetic Genealogy ([www.isogg.org](http://www.isogg.org)), restricting analysis using



*samtools*<sup>52</sup> to sites with map quality and base quality of at least 30, and excluding 2 bases at the ends of each sequenced fragment.

### Genome-wide scan for selection

For most ancient samples, we did not have sufficient coverage to make reliable diploid calls. We therefore used the counts of sequences covering each SNP to compute the likelihood of the allele frequency in each population. Suppose that at a particular site, for each population we have  $M$  samples with sequence level data, and  $N$  samples for which we had hard genotype calls (Loschbour, Stuttgart and the 1,000 Genomes samples). For samples  $i=1..N$ , with genotype data, we observe  $X$  copies of the reference allele out of  $2N$  total chromosomes. For each of samples  $i=(N+1)..(N+M)$ , with sequence level data, we observe  $R_i$  sequences with the reference allele out of  $T_i$  total sequences. Then, dropping the subscript  $i$  for brevity, the likelihood of the population reference allele frequency,  $p$  given data

$D = \{X, N, \bar{R}, \bar{T}\}$  is given by

$$L(p; D) = B(X, 2N, p) \prod_{i=N+1}^{N+M} \left\{ p^2 B(R_i, T_i, 1 - \varepsilon) + 2p(1 - p) B(R_i, T_i, 0.5) + (1 - p)^2 B(R_i, T_i, \varepsilon) \right\},$$

where  $B(k, n, p) = \binom{n}{k} p^k (1 - p)^{n-k}$  is the binomial probability distribution and  $\varepsilon$  is a small probability of error, which we set to 0.001. We write  $\ell(p; D)$  for the log-likelihood. To estimate allele frequencies, for example in Fig. 3 or for the polygenic selection test, we maximized this likelihood numerically for each population.

To scan for selection across the genome, we used the following test. Consider a single SNP. Assume that we can model the allele frequencies  $\mathbf{p}_{mod}$  in  $A$  modern populations as a linear combination of allele frequencies in  $B$  ancient populations  $\mathbf{p}_{anc}$ . That is,  $\mathbf{p}_{mod} = \mathbf{C} \mathbf{p}_{anc}$  where  $\mathbf{C}$  is an  $A$  by  $B$  matrix with rows summing to 1. We have data  $D_j$  from population  $j$  which is some combination of sequence counts and genotypes as described above. Then, writing  $\bar{\mathbf{p}} = [\mathbf{p}_{anc}, \mathbf{p}_{mod}] = [\mathbf{p}_1 \dots \mathbf{p}_{A+B}]$  the log-likelihood of the allele frequencies equals the sum of the log-likelihoods for each population.

$$\ell(\bar{\mathbf{p}}, \bar{D}) = \sum_{j=1}^{A+B} \ell(\mathbf{p}_j; D_j)$$

To detect deviations in allele frequency from expectation, we test the null hypothesis  $H_0$ :  $\mathbf{p}_{mod} = \mathbf{C} \mathbf{p}_{anc}$  against the alternative  $H_1$ :  $\mathbf{p}_{mod}$  unconstrained. We numerically maximize this likelihood in both the constrained and unconstrained model and use the fact that twice the difference in log-likelihood is approximately  $\chi_A^2$  distributed to compute a test statistic and P-value.

We defined the ancient source populations by the ‘‘Selection group 1’’ label in Extended Data Table 1 and Supplementary Table 1 and used the 1000 Genomes CEU, GBR, IBS and

TSI as the present-day populations. We removed SNPs that were monomorphic in all four of these modern populations as well as in 1000 Genomes Yoruba (YRI). We do not use FIN as one of the modern populations, because they do not fit this three-population model well. We estimate the proportions of (HG, EF, SA) to be  $CEU=(0.196, 0.257, 0.547)$ ,  $GBR=(0.362, 0.229, 0.409)$ ,  $IBS=(0, 0.686, 0.314)$  and  $TSI=(0, 0.645, 0.355)$ . In practice we found that there was substantial inflation in the test statistic, most likely due to unmodeled ancestry or additional drift. To address this, we applied a genomic control correction<sup>14</sup>, dividing all the test statistics by a constant,  $\lambda$ , chosen so that the median p-value matched the median of the null  $\chi_4^2$  distribution. Excluding sites in the potentially functional set, we estimated  $\lambda=1.38$  and used this value as a correction throughout. One limitation of this test is that, although it identifies likely signals of selection, it cannot provide much information about the strength or date of selection. If the ancestral populations in the model are, in fact, close to the real ancestral populations, then any selection must have occurred after the first admixture event (in this case, after 6500 BCE), but if the ancestral populations are misspecified, even this might not be true.

To estimate power, we randomly sampled allele counts from the full dataset, restricting to polymorphic sites with a mean frequency across all populations of  $<0.1$ . We then simulated what would happen if the allele had been under selection in all of the modern populations by simulating a Wright-Fisher trajectory with selection for 50, 100 or 200 generations, starting at the observed frequency. We took the final frequency from this simulation, sampled observations to replace the actual observations in that population, and counted the proportion of simulations that gave a genome-wide significant result after GC correction (Extended Data Fig. 6a). We resampled sequence counts for the observed distribution for each population to simulate the effect of increasing sample size, assuming that the coverage and distribution of the sequences remained the same (Extended Data Fig. 6b).

We investigated how the genomic control correction responded when we simulated small amounts of admixture from a highly diverged population (Yoruba; 1000 Genomes YRI) into a randomly chosen modern population. The genomic inflation factor increases from around 1.38 to around 1.51 with 10% admixture, but there is little reduction in power (Extended Fig. 6c). Finally, we investigated how robust the test was to misspecification of the mixture matrix  $C$ . We reran the power simulations using a matrix  $C'=pC+(1-p)R$  for  $p \in [0, 1]$  where  $R$  was a random matrix chosen so that for each modern population, the mixture proportions of the three ancient populations were jointly normally distributed on  $[0, 1]$ . Increasing  $p$  increases the genomic inflation factor and reduces power, demonstrating the advantage of explicitly modeling the ancestries of the modern populations (Extended Fig. 6d).

### Test for polygenic selection

We implemented the test for polygenic selection described by Ref. 37. This evaluates whether trait-associated alleles, weighted by their effect size, are over-dispersed compared to randomly sampled alleles, in the directions associated with the effects measured by genome-wide association studies (GWAS). For each trait, we obtained a list of significant SNP associations and effect estimates from GWAS data, and then applied the test both to all

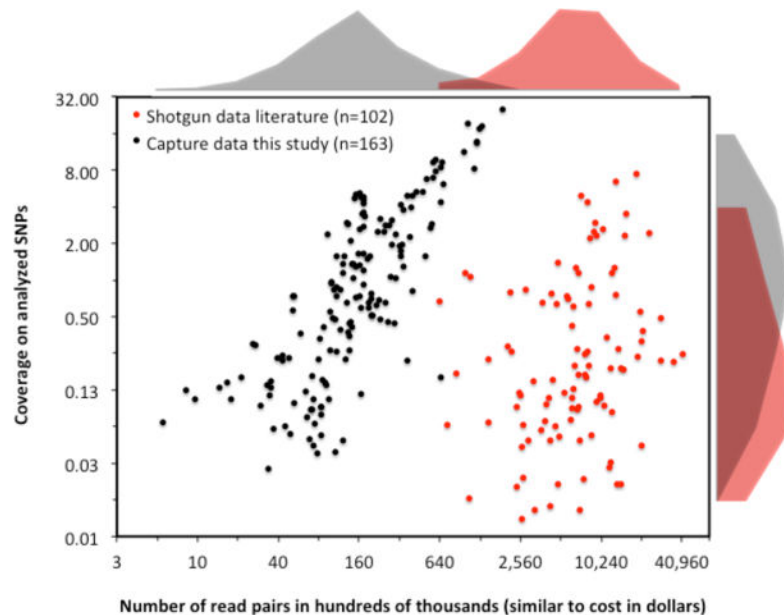
populations combined and to selected pairs of populations. We restricted the list of GWAS associations to 169 SNPs where we observed at least two chromosomes in all tested populations (selection population 2). We estimated frequencies in each population by computing the MLE, using the likelihood described above. For each test, we sampled SNPs frequency matched in 20 bins, computed the test statistic  $Q_X$  and for ease of comparison, converted these to Z scores, signed according to the direction of the genetic effects. Theoretically  $Q_X$  has a  $\chi^2$  distribution but in practice, it is over-dispersed. Therefore, we report bootstrap p-values computed by sampling 10,000 sets of frequency matched SNPs.

To estimate population-level genetic height in Fig. 4A, we assumed a uniform prior on [0,1] for the distribution of all height-associated alleles, and then sampled from the posterior joint frequency distribution of the alleles, assuming they were independent, using a Metropolis-Hastings sampler with a  $N(0,0.001)$  proposal density. We then multiplied the sampled allele frequencies by the effect sizes to get a distribution of genetic height.

### Code availability

Code implementing the selection analysis is available at [https://github.com/mathii/europe\\_selection](https://github.com/mathii/europe_selection).

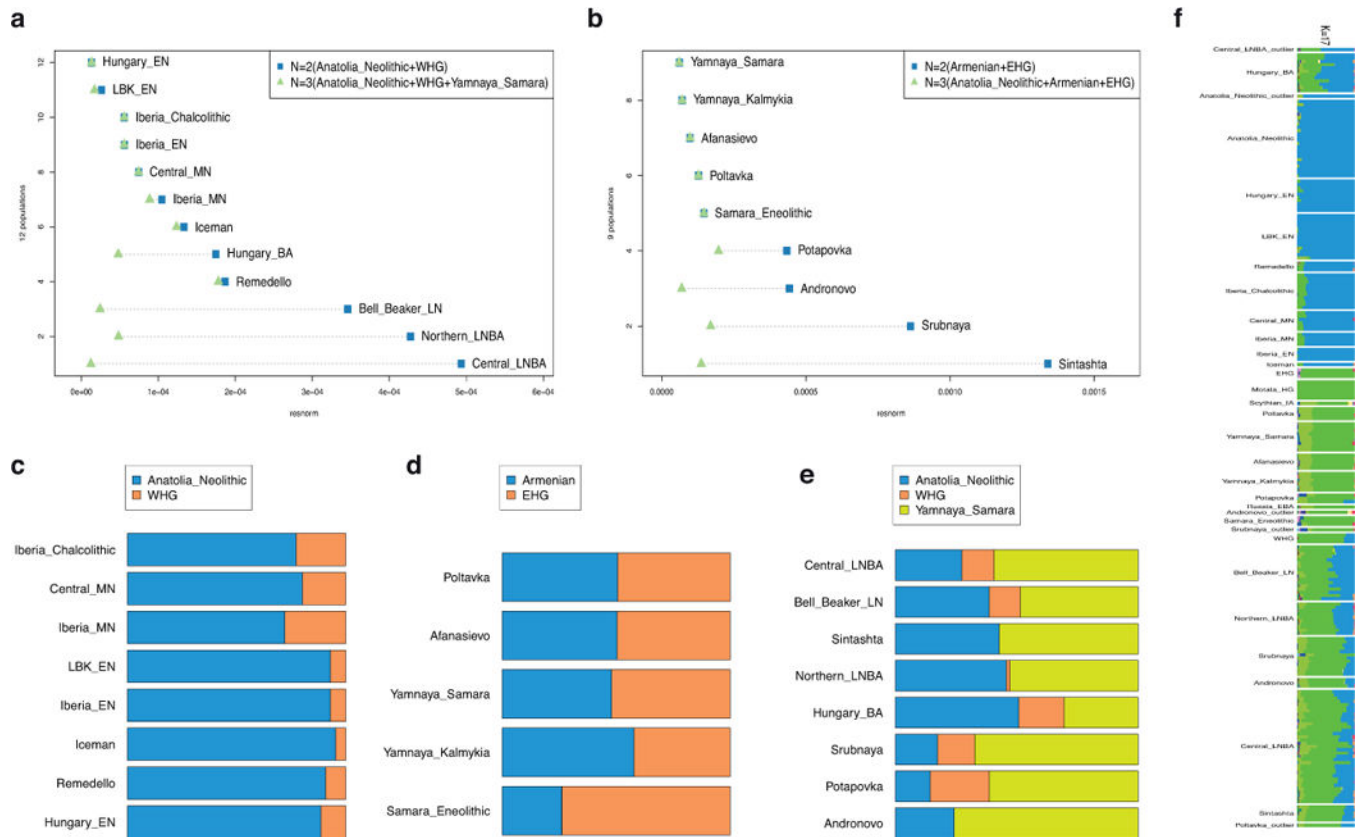
### Extended Data



	Shotgun data (Literature)	Capture data (Newly reported)
Average number of raw read pairs (hundreds of thousands, similar to cost in dollars)	8730	241
Median coverage on analyzed SNPs	0.19	0.75

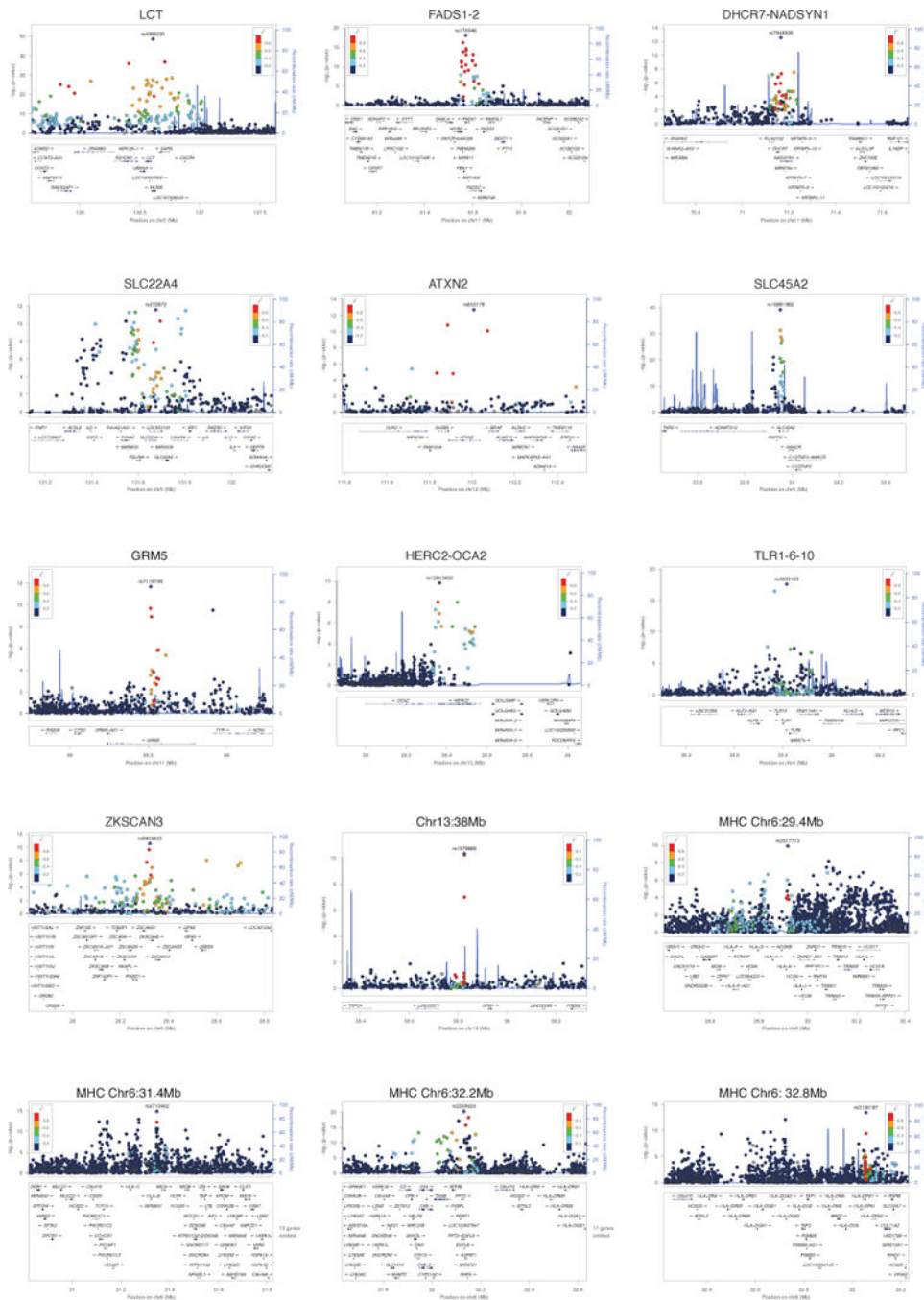
**Extended Data Figure 1. Efficiency and cost-effectiveness of 1240k capture**

We plot the number of raw sequences against the mean coverage of analyzed SNPs after removal of duplicates, comparing the 163 samples for which capture data are reported in this study, against the 102 samples analyzed by shotgun sequencing in ref.<sup>5</sup> We caution that the true cost is more than that of sequencing alone.



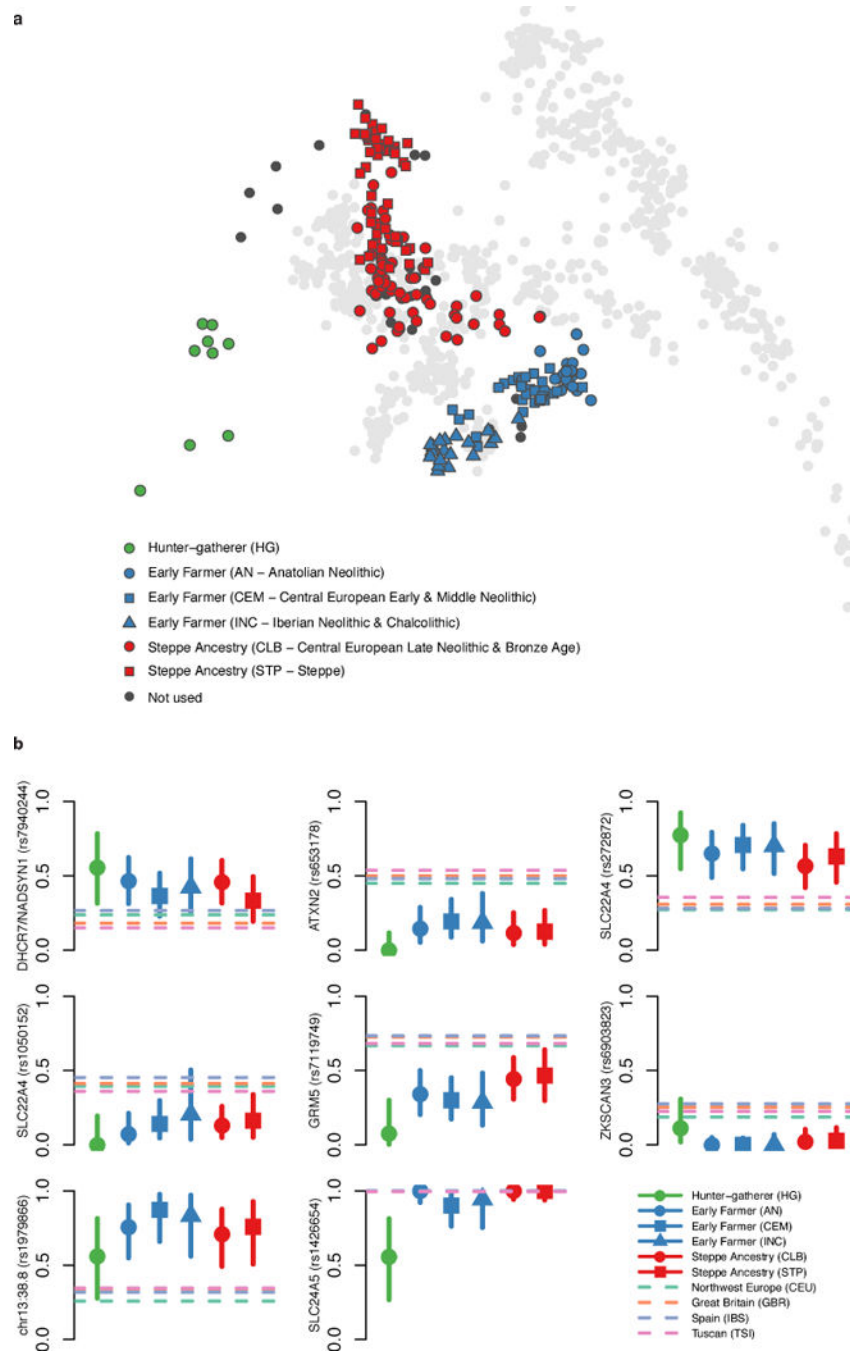
### Extended Data Figure 2. Early isolation and later admixture between farmers and steppe populations

**A:** Mainland European populations later than 3000 BCE are better modeled with steppe ancestry as a 3<sup>rd</sup> ancestral population. **B:** Later (post-Poltavka) steppe populations are better modeled with Anatolian Neolithic as a 3<sup>rd</sup> ancestral population. **C:** Estimated mixture proportions of mainland European populations without steppe ancestry. **D:** Estimated mixture proportions of Eurasian steppe populations without Anatolian Neolithic ancestry. **E:** Estimated mixture proportions of later populations with both steppe and Anatolian Neolithic ancestry. **F:** ADMIXTURE plot at  $k=17$  showing population differences over time and space.



**Extended Data Figure 3. Regional association plots**

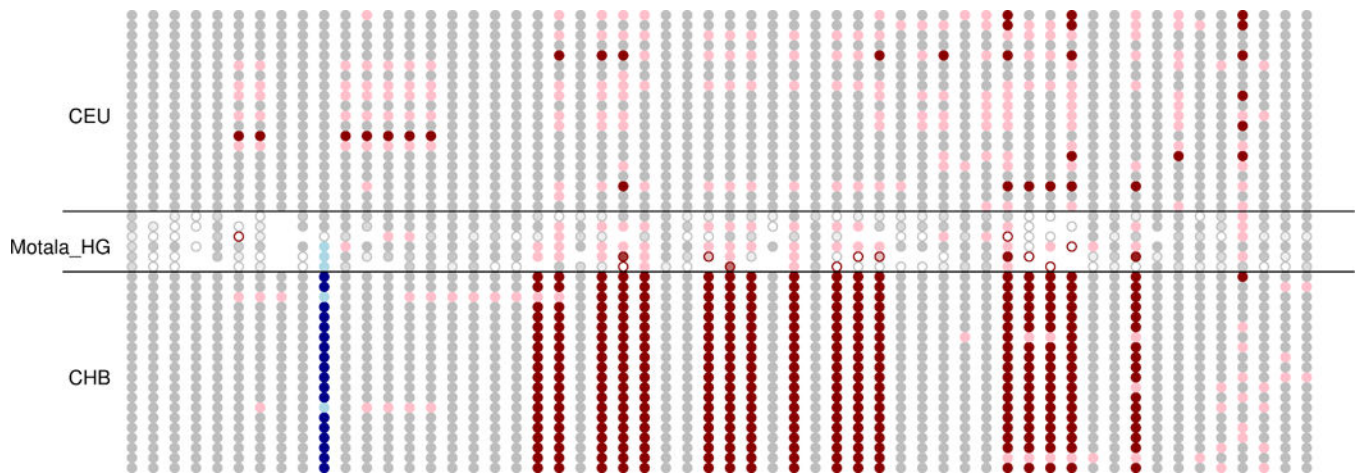
LocusZoom<sup>60</sup> plots for genome-wide significant signals. Points show the  $-\log_{10}$  P-value for each SNP, colored according to their LD with the most associated SNP. The blue line shows the recombination rate, with scale on right hand axis. Genes are shown in the lower panel of each subplot.



**Extended Data Figure 4. PCA of selection populations and derived allele frequencies for genome-wide significant signals**

**A:** Ancient samples projected onto principal components of modern samples, as in Fig. 1, but labeled according to selection populations defined in Extended Data Table 1. **B:** Allele frequency plots as in Fig. 3 – for *SLC22A4* we show both rs272872, which is our strongest signal, and rs1050152, which was previously hypothesized to be under selection – and we also show *SLC24A5*, which is not genome-wide significant but is discussed in the main text.





**Extended Data Figure 5. Motala haplotypes carrying the derived, selected EDAR allele**

This figure compares the genotypes at all sites within 150kb of rs3827760 (in blue) for the 6 Motala samples and 20 randomly chosen CHB (Chinese from Beijing) and CEU (Utah residents with northern and western European ancestry) samples. Each row is a sample and each column is a SNP. Grey means homozygous for the major (in CEU) allele. Pink denotes heterozygous and red homozygous for the other allele. For the Motala samples, an open circle means that there is only a single sequence otherwise the circle is colored according to the number of sequences observed. Three of the Motala samples are heterozygous for rs3827760 and the derived allele lies on the same haplotype background as in present-day East Asians. The only other ancient samples with evidence of the derived EDAR allele in this dataset are two Afanasievo samples dating to 3300-3000 BCE, and one Scythian dating to 400-200 BCE (not shown).

Lead SNP	Chromosome	Position (hg19)	P Value	Range (Mb)	Genes	Potential function
rs4988235	2	136,608,646	3.19E-49	135.3-137.3	MCM6,LCT	Lactase persistence <sup>*15</sup>
rs16891982	5	33,951,693	7.05E-40	33.8-34.0	SLC45A2	Skin pigmentation <sup>*53</sup>
rs2269424	6	32,132,233	5.41E-21	29.9-33.1	MHC region	Immunity*
rs174546	11	61,569,830	8.18E-19	61.5-61.6	FADS1,FADS2	Fatty acid metabolism <sup>*18,54,55</sup>
rs4833103	4	38,815,502	2.58E-18	38.7-38.8	TLR1,TLR6,TLR10	Immunity <sup>*31,32</sup>
rs653178	12	112,007,756	1.96E-13	111.9-112.6	ATXN2,SH2B3	Unknown
rs7944926	11	71,165,625	2.86E-13	71.2-71.2	DHCR7,NADSYN1	Vitamin D metabolism <sup>*21,56</sup>
rs7119749	11	88,515,022	2.03E-12	88.5-88.9	GRM5	Skin pigmentation <sup>*27</sup>
rs272872	5	131,675,864	2.56E-12	131.4-131.8	SLC22A4	Ergothioneine transport <sup>*57</sup>
rs6903823	6	28,322,296	2.96E-11	28.3-28.7	ZKSCAN3,ZSCAN31	Autophagy <sup>*58</sup> , Lung function <sup>*59</sup>
rs1979866	13	38,825,900	4.60E-11	38.1-38.8	-	Unknown
rs12913832	11	28,365,618	1.5E-10	29.9-30.1	HERC2,OCA2	Eye color <sup>*28,29</sup>

**Extended Data Figure 6. Estimated power of the selection scan**

**A:** Estimated power for different selection coefficients for a SNP that is selected in all populations for either 50, 100 or 200 generations. **B:** Effect of increasing sample size, showing estimated power for a SNP selected for 100 generations, with different amounts of

data, relative to the main text. **C**: Effect of admixture from Yoruba (YRI) into one of the modern populations, showing the effect on the genomic inflation factor (blue, left axis) and the power to detect selection on a SNP selected for 100 generations with a selection coefficient of 0.02. **D**: Effect of mis-specification of the mixture proportions. Here 0 on the x-axis corresponds to the proportions we used, and 1 corresponds to a random mixture matrix.

**Extended Data Table 2**

Key  $f_4$ -statistics used to support claims about population history.

A	B	C	D	$f_4(A, B, C, D)$	Z	Number of SNPs	Interpretation
Anatolia_Neolithic	LBK_EN	WHG	Chimp	-0.00114	-6.8	1003751w1	Early European Farmers had more WHG ancestry than Anatolian Neolithic
Anatolia_Neolithic	Hungary_EN	WHG	Chimp	-0.00212	-11.9	929553	
Anatolia_Neolithic	Iberia_EN	WHG	Chimp	-0.00244	-9.6	904437	
Iberia_EN	Iberia_Chalcolithic	WHG	Chimp	-0.00311	-10.5	802471	Iberian Chalcolithic had more WHG ancestry than Iberian Early Neolithic
Iberia_MN	Iberia_Chalcolithic	WHG	Chimp	0.00010	0.3	779905	Iberian Chalcolithic did not have more WHG ancestry than Iberian Middle Neolithic
EHG	Samara_Eneolithic	MA1	Chimp	0.00140	2.3	463388	First dilution of Ancient North Eurasian ancestry (prior to the Bronze Age Yamnaya culture)
EHG	Yamnaya_Samara	MA1	Chimp	0.00513	10.6	645211	
Samara_Eneolithic	Yamnaya_Samara	MA1	Chimp	0.00366	7.6	482492	
EHG	Yamnaya_Samara	Armenian	Chimp	-0.00191	-6.1	547370	Contribution of Near Eastern ancestry to the Bronze Age Yamnaya culture
EHG	Yamnaya_Kalmykia	Armenian	Chimp	-0.00180	-5.4	536989	
Samara_Eneolithic	Yamnaya_Samara	Armenian	Chimp	-0.00100	-3.3	405599	
EHG	Poltavka	Armenian	Chimp	-0.00175	-4.9	541983	
Yamnaya_Samara	Yamnaya_Kalmykia	MA1	Chimp	-0.00010	-0.3	675630	Stability of Ancient North Eurasian ancestry between Early Bronze Age Yamnaya from Kalmykia and Samara, and the Middle Bronze Age Poltavka
Yamnaya_Samara	Poltavka	MA1	Chimp	-0.00014	-0.4	673726	
Yamnaya_Kalmykia	Poltavka	MA1	Chimp	0.00012	0.3	659346	
Yamnaya_Samara	Srubnaya	MA1	Chimp	0.00151	5.1	691149	Second dilution of Ancient North Eurasian ancestry (prior to the Late Bronze Age Srubnaya culture)
Yamnaya_Kalmykia	Srubnaya	MA1	Chimp	0.00161	4.8	676735	
Poltavka	Srubnaya	MA1	Chimp	0.00164	4.5	674756	

A		B	C	D	$f_4(A, B, C, D)$	Z	Number of SNPs	Interpretation
Yamnaya_Samara		Srubnaya	LBK_EN	Chimp	-0.00225	-11.4	974659	Arrival of Early European Farmer-related ancestry prior to the Late Bronze Age Srubnaya culture. Statistics with Anatolia_Neolithic instead of LBK_EN are similar ( $Z < -8$ , not shown).
Yamnaya_Kalmykia		Srubnaya	LBK_EN	Chimp	-0.00264	-11.4	951827	
Poltavka		Srubnaya	LBK_EN	Chimp	-0.00210	-9.0	948968	
EHG	Yamnaya_Samara	Armenian	LBK_EN		-0.00080	-5.0	559478	Different source of dilution of Ancient North Eurasian ancestry prior to the Yamnaya (Near Eastern) vs. prior to the Srubnaya (Early European Farmer-related)
EHG	Yamnaya_Kalmykia	Armenian	LBK_EN		-0.00086	-5.2	548882	
EHG	Poltavka	Armenian	LBK_EN		-0.00069	-4.1	553996	
Yamnaya_Samara		Srubnaya	Armenian	LBK_EN	0.00138	13.1	585240	
Yamnaya_Kalmykia		Srubnaya	Armenian	LBK_EN	0.00142	11.3	574333	
Poltavka		Srubnaya	Armenian	LBK_EN	0.00134	10.7	577082	
Ref <sub>1</sub>	Ref <sub>2</sub>	Test		$f_3(\text{Test}; \text{Ref 1}, \text{Ref 2})$	Z	Number of SNPs	Interpretation	
WHG	Anatolia_Neolithic	Hungary_EN		-0.00412	-6.7	548445	Early European farmers were formed by admixture between Anatolia Neolithic and WHG (the non-significant signal in the Iberian_EN may be due to genetic drift specific to this population)	
WHG	Anatolia_Neolithic	LBK_EN		-0.00257	-4.6	654357		
WHG	Anatolia_Neolithic	Iberia_EN		0.00179	1.4	389101		
EHG	Armenian	Poltavka		-0.00539	-3.9	213055	Early and Middle Bronze Age steppe pastoralists were formed by admixture between EHG and a population of Near Eastern ancestry	
EHG	Armenian	Yamnaya_Kalmykia		-0.00537	-4.2	213996		
EHG	Armenian	Yamnaya_Samara		-0.00586	-6.2	276568		
LBK_EN	Yamnaya_Samara	Srubnaya		-0.00630	-11.2	584111	Srubnaya was formed by admixture between populations related to Yamnaya and Early European Farmers	

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgments

We thank Paul de Bakker, Joachim Burger, Christos Economou, Elin Fornander, Qiaomei Fu, Fredrik Hallgren, Karola Kirsanow, Alissa Mittnik, Iñigo Olalde, Adam Powell, Pontus Skoglund, Shervin Tabrizi, and Arti Tandon for discussions, suggestions about SNPs to include, or contribution to sample preparation or data curation. We thank Svante Pääbo, Matthias Meyer, Qiaomei Fu, and Birgit Nickel for collaboration in developing the 1240k capture reagent. We thank Julio Manuel Vidal Encinas and María Encina Prada for allowing us to resample La Braña 1 I.M. was supported by the Human Frontier Science Program LT001095/2014-L. C.G. was supported by the Irish Research Council for Humanities and Social Sciences (IRCHSS). A.K., P.K. and O.M. were supported by RFBR №15-06-01916 and RFH№15-11-63008 and O.M. by a state grant of the Ministry of education and science of Russia Federation #33.1195.2014/k. J.K. was supported by ERC starting grant APGREID and DFG grant KR 4015/1-1. K.W.A. was supported by DFG grant AL 287/14-1. C.L.-F. was supported by a BFU2015-64699-P. grant from the Spanish government. W.H. and B.L. were supported by Australian Research Council DP130102158. R.P. was supported by ERC starting grant ADNABIOARC (263441), and an Irish Research Council ERC support grant. D.R. was supported by U.S. National Science Foundation HOMINID grant BCS-1032255, U.S. National Institutes of Health grant GM100233, and the Howard Hughes Medical Institute.

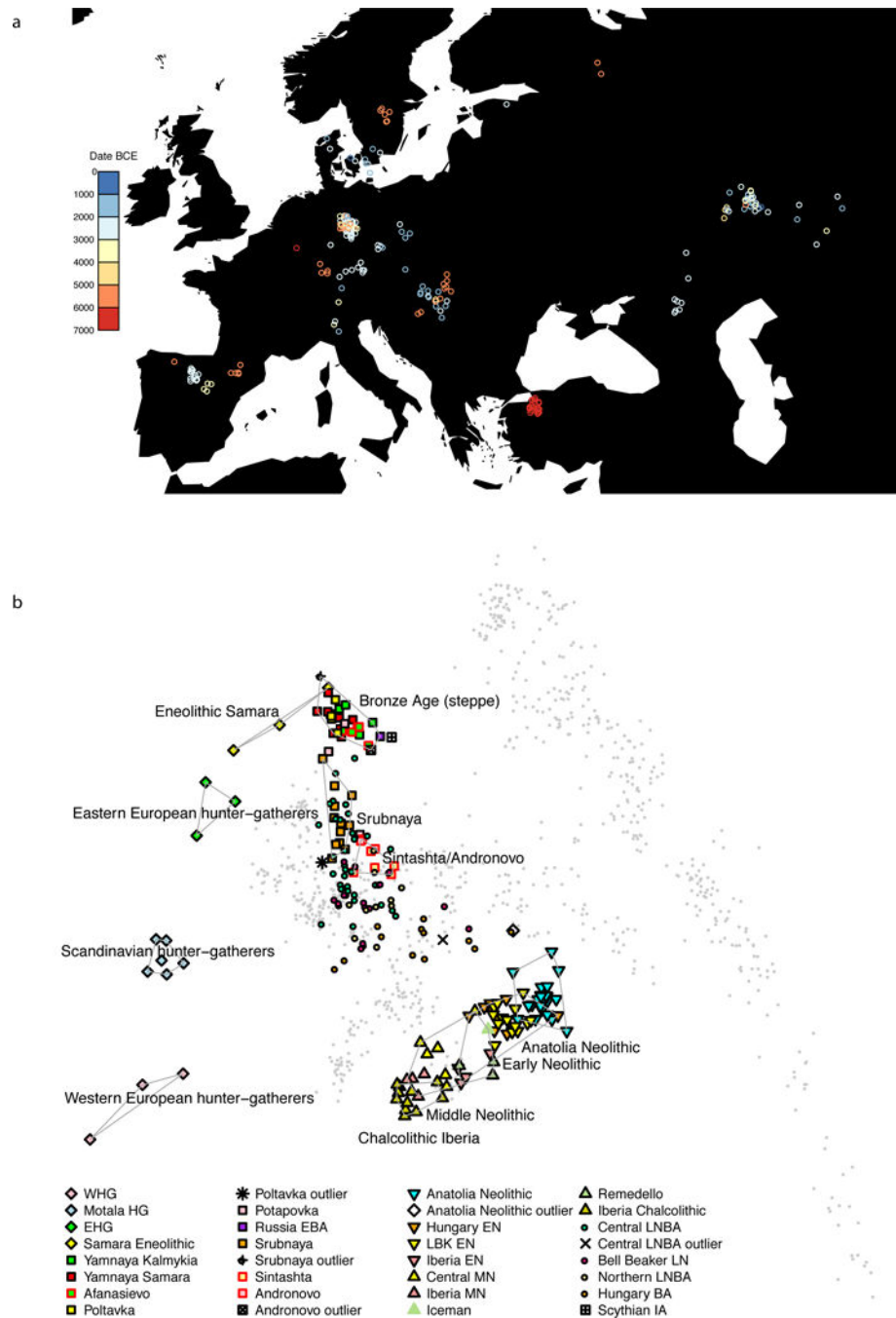
## References

- Grossman SR, et al. Identifying recent adaptations in large-scale genomic data. *Cell*. 2013; 152:703–713. [PubMed: 23415221]
- Wilde S, et al. Direct evidence for positive selection of skin, hair, and eye pigmentation in Europeans during the last 5,000 y. *Proc Natl Acad Sci U S A*. 2014; 111:4832–4837. [PubMed: 24616518]
- Gamba C, et al. Genome flux and stasis in a five millennium transect of European prehistory. *Nat Commun*. 2014; 5:5257. [PubMed: 25334030]
- Lazaridis I, et al. Ancient human genomes suggest three ancestral populations for present-day Europeans. *Nature*. 2014; 513:409–413. [PubMed: 25230663]
- Allentoft ME, et al. Population genomics of Bronze Age Eurasia. *Nature*. 2015; 522:167–172. [PubMed: 26062507]
- Keller A, et al. New insights into the Tyrolean Iceman's origin and phenotype as inferred by whole-genome sequencing. *Nat Commun*. 2012; 3:698. [PubMed: 22426219]
- Haak W, et al. Massive migration from the steppe was a source for Indo-European languages in Europe. *Nature*. 2015; 522:207–211. [PubMed: 25731166]
- Olalde I, et al. Derived immune and ancestral pigmentation alleles in a 7,000-year-old Mesolithic European. *Nature*. 2014; 507:225–228. [PubMed: 24463515]
- Pinhasi R, et al. Optimal Ancient DNA Yields from the Inner Ear Part of the Human Petrous Bone. *PLoS One*. 2015; 10:e0129102. [PubMed: 26086078]
- Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in unrelated individuals. *Genome Res*. 2009; 19:1655–1664. [PubMed: 19648217]
- Patterson N, et al. Ancient admixture in human history. *Genetics*. 2012; 192:1065–1093. [PubMed: 22960212]
- Underhill PA, et al. The phylogenetic and geographic structure of Y-chromosome haplogroup R1a. *Eur J Hum Genet*. 2015; 23:124–131. [PubMed: 24667786]
- The 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature*. 2015; 526:68–74. [PubMed: 26432245]
- Devlin B, Roeder K. Genomic control for association studies. *Biometrics*. 1999; 55:997–1004. [PubMed: 11315092]
- Enattah NS, et al. Identification of a variant associated with adult-type hypolactasia. *Nat Genet*. 2002; 30:233–237. [PubMed: 11788828]
- Bersaglieri T, et al. Genetic signatures of strong recent positive selection at the lactase gene. *Am J Hum Genet*. 2004; 74:1111–1120. [PubMed: 15114531]
- Burger J, Kirchner M, Bramanti B, Haak W, Thomas MG. Absence of the lactase-persistence-associated allele in early Neolithic Europeans. *Proc Natl Acad Sci U S A*. 2007; 104:3736–3741. [PubMed: 17360422]

18. Teslovich TM, et al. Biological, clinical and population relevance of 95 loci for blood lipids. *Nature*. 2010; 466:707–713. [PubMed: 20686565]
19. Fumagalli M, et al. Greenlandic Inuit show genetic signatures of diet and climate adaptation. *Science*. 2015; 349:1343–1347. [PubMed: 26383953]
20. Mathias RA, et al. Adaptive evolution of the FADS gene cluster within Africa. *PLoS One*. 2012; 7:e44926. [PubMed: 23028684]
21. Wang TJ, et al. Common genetic determinants of vitamin D insufficiency: a genome-wide association study. *Lancet*. 2010; 376:180–188. [PubMed: 20541252]
22. Price AL. The impact of divergence time on the nature of population structure: an example from Iceland. *PLoS Genet*. 2009; 5:e1000505. [PubMed: 19503599]
23. Wellcome Trust Case Control, C. Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature*. 2007; 447:661–678. [PubMed: 17554300]
24. Huff CD, et al. Crohn's disease and genetic hitchhiking at IBD5. *Mol Biol Evol*. 2012; 29:101–111. [PubMed: 21816865]
25. Hunt KA, et al. Newly identified genetic risk variants for celiac disease related to the immune response. *Nat Genet*. 2008; 40:395–402. [PubMed: 18311140]
26. Jostins L, et al. Host-microbe interactions have shaped the genetic architecture of inflammatory bowel disease. *Nature*. 2012; 491:119–124. [PubMed: 23128233]
27. Beleza S, et al. Genetic architecture of skin and eye color in an African-European admixed population. *PLoS Genet*. 2013; 9:e1003372. [PubMed: 23555287]
28. Sturm RA, et al. A single SNP in an evolutionary conserved region within intron 86 of the HERC2 gene determines human blue-brown eye color. *Am J Hum Genet*. 2008; 82:424–431. [PubMed: 18252222]
29. Eiberg H, et al. Blue eye color in humans may be caused by a perfectly associated founder mutation in a regulatory element located within the HERC2 gene inhibiting OCA2 expression. *Hum Genet*. 2008; 123:177–187. [PubMed: 18172690]
30. Barreiro LB, et al. Evolutionary dynamics of human Toll-like receptors and their different contributions to host defense. *PLoS Genet*. 2009; 5:e1000562. [PubMed: 19609346]
31. Uciechowski P, et al. Susceptibility to tuberculosis is associated with TLR1 polymorphisms resulting in a lack of TLR1 cell surface expression. *J Leukoc Biol*. 2011; 90:377–388. [PubMed: 21642391]
32. Wong SH, et al. Leprosy and the adaptation of human toll-like receptor 1. *PLoS Pathog*. 2010; 6:e1000979. [PubMed: 20617178]
33. Fujimoto A, et al. A scan for genetic determinants of human hair morphology: EDAR is associated with Asian hair thickness. *Hum Mol Genet*. 2008; 17:835–843. [PubMed: 18065779]
34. Kimura R, et al. A common variation in EDAR is a genetic determinant of shovel-shaped incisors. *Am J Hum Genet*. 2009; 85:528–535. [PubMed: 19804850]
35. Kamberov YG, et al. Modeling recent human evolution in mice by expression of a selected EDAR variant. *Cell*. 2013; 152:691–702. [PubMed: 23415220]
36. Turchin MC, et al. Evidence of widespread selection on standing variation in Europe at height-associated SNPs. *Nat Genet*. 2012; 44:1015–1019. [PubMed: 22902787]
37. Berg JJ, Coop G. A population genetic signal of polygenic adaptation. *PLoS Genet*. 2014; 10:e1004412. [PubMed: 25102153]
38. Lango Allen H, et al. Hundreds of variants clustered in genomic loci and biological pathways affect human height. *Nature*. 2010; 467:832–838. [PubMed: 20881960]
39. Speliotes EK, et al. Association analyses of 249,796 individuals reveal 18 new loci associated with body mass index. *Nat Genet*. 2010; 42:937–948. [PubMed: 20935630]
40. Heid IM, et al. Meta-analysis identifies 13 new loci associated with waist-hip ratio and reveals sexual dimorphism in the genetic basis of fat distribution. *Nat Genet*. 2010; 42:949–960. [PubMed: 20935629]
41. Morris AP, et al. Large-scale association analysis provides insights into the genetic architecture and pathophysiology of type 2 diabetes. *Nat Genet*. 2012; 44:981–990. [PubMed: 22885922]

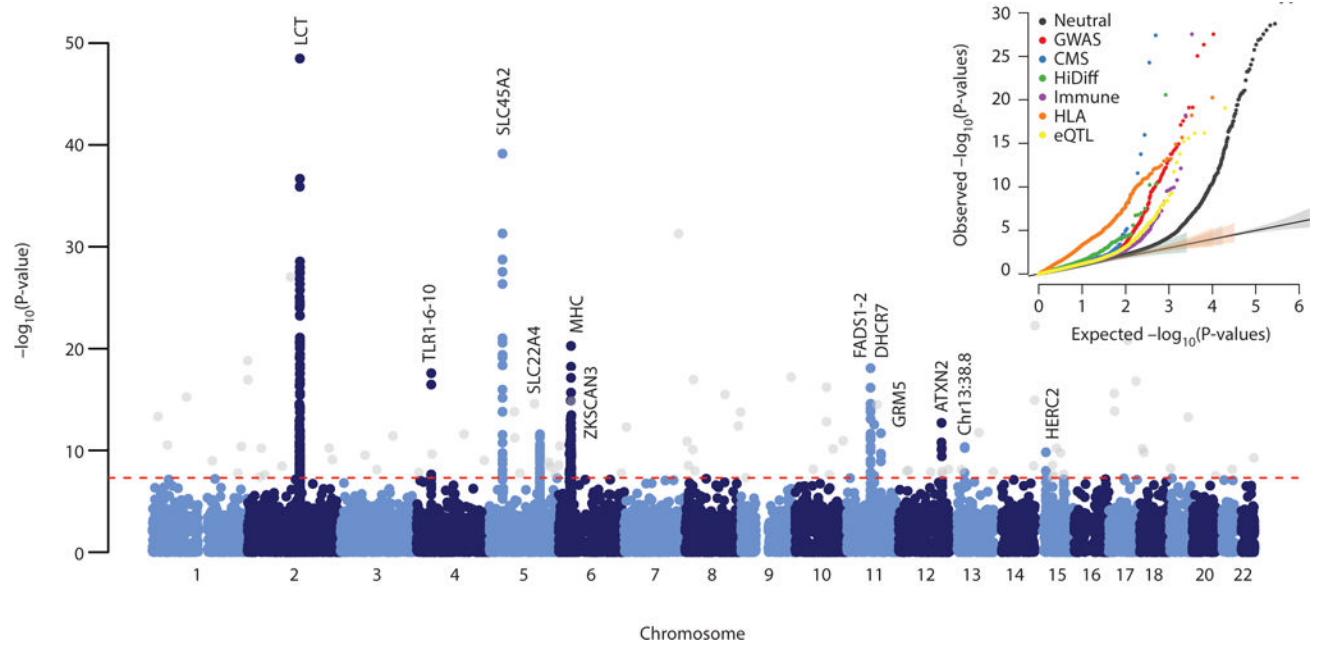
42. Briggs AW. Removal of deaminated cytosines and detection of in vivo methylation in ancient DNA. *Nucleic Acids Res.* 2010; 38:e87. [PubMed: 20028723]
43. Fu Q, et al. DNA analysis of an early modern human from Tianyuan Cave, China. *Proc Natl Acad Sci U S A.* 2013; 110:2223–2227. [PubMed: 23341637]
44. Fu Q, et al. An early modern human from Romania with a recent Neanderthal ancestor. *Nature.* 2015
45. Korneliussen TS, Albrechtsen A, Nielsen R. ANGSD: Analysis of Next Generation Sequencing Data. *BMC Bioinformatics.* 2014; 15:356. [PubMed: 25420514]
46. International HapMap Consortium. A second generation human haplotype map of over 3.1 million SNPs. *Nature.* 2007; 449:851–861. [PubMed: 17943122]
47. Lappalainen T, et al. Transcriptome and genome sequencing uncovers functional variation in humans. *Nature.* 2013; 501:506–511. [PubMed: 24037378]
48. Li JZ, et al. Worldwide human relationships inferred from genome-wide patterns of variation. *Science.* 2008; 319:1100–1104. [PubMed: 18292342]
49. Loh PR, et al. Inferring admixture histories of human populations using linkage disequilibrium. *Genetics.* 2013; 193:1233–1254. [PubMed: 23410830]
50. Chang CC, et al. Second-generation PLINK: rising to the challenge of larger and richer datasets. *GigaScience.* 2015; 4
51. Skoglund P, Storå J, Götherström A, Jakobsson M. Accurate sex identification of ancient human remains using DNA shotgun sequencing. *JAS.* 2013; 40:4477–4482.
52. Li H, et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics.* 2009; 25:2078–2079. [PubMed: 19505943]
53. Norton HL. Genetic evidence for the convergent evolution of light skin in Europeans and East Asians. *Mol Biol Evol.* 2007; 24:710–722. [PubMed: 17182896]
54. Bokor S, et al. Single nucleotide polymorphisms in the FADS gene cluster are associated with delta-5 and delta-6 desaturase activities estimated by serum fatty acid ratios. *J Lipid Res.* 2010; 51:2325–2333. [PubMed: 20427696]
55. Tanaka T, et al. Genome-wide association study of plasma polyunsaturated fatty acids in the InCHIANTI Study. *PLoS Genet.* 2009; 5:e1000338. [PubMed: 19148276]
56. Ahn J, et al. Genome-wide association study of circulating vitamin D levels. *Hum Mol Genet.* 2010; 19:2739–2745. [PubMed: 20418485]
57. Grundemann D, et al. Discovery of the ergothioneine transporter. *Proc Natl Acad Sci U S A.* 2005; 102:5256–5261. [PubMed: 15795384]
58. Chauhan S, et al. ZKSCAN3 is a master transcriptional repressor of autophagy. *Mol Cell.* 2013; 50:16–28. [PubMed: 23434374]
59. Soler Artigas M, et al. Genome-wide association and large-scale follow up identifies 16 new loci influencing lung function. *Nat Genet.* 2011; 43:1082–1090. [PubMed: 21946350]
60. Pruim RJ, et al. LocusZoom: regional visualization of genome-wide association scan results. *Bioinformatics.* 2010; 26:2336–2337. [PubMed: 20634204]





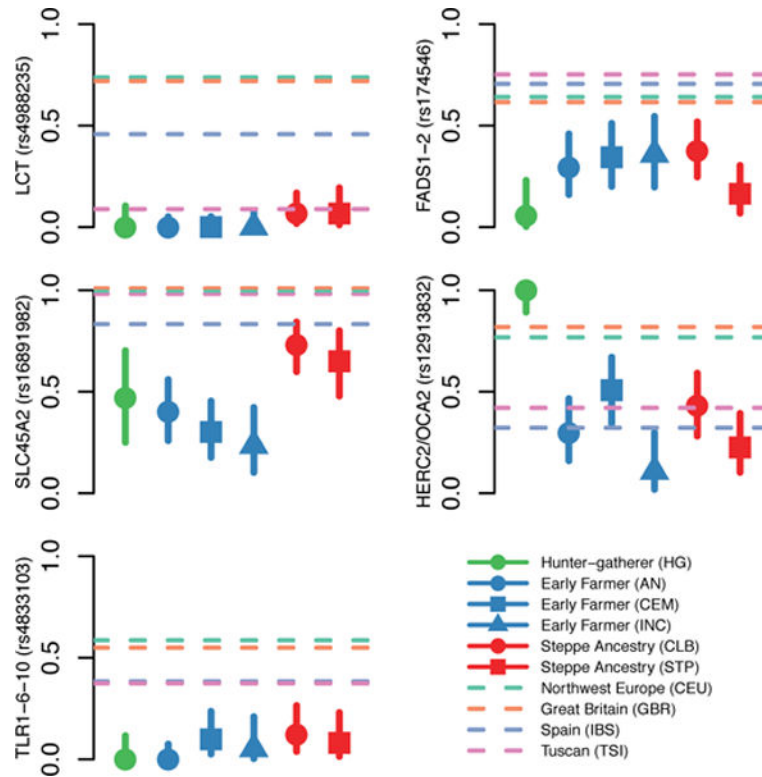
**Figure 1. Population relationships of samples**

**A.** Locations color-coded by date, with a random jitter added for visibility (8 Afanasievo and Andronovo samples lie further east and are not shown). **B:** Principal component analysis of 777 modern West Eurasian samples (grey), with 221 ancient samples projected onto the first two principal component axes and labeled by culture. **Abbreviations:** [E/M/L]N Early/Middle/Late Neolithic, LBK *Linearbandkeramik*, [E/W]HG Eastern/Western hunter-gatherer, [E]BA [Early] Bronze Age, IA Iron Age.



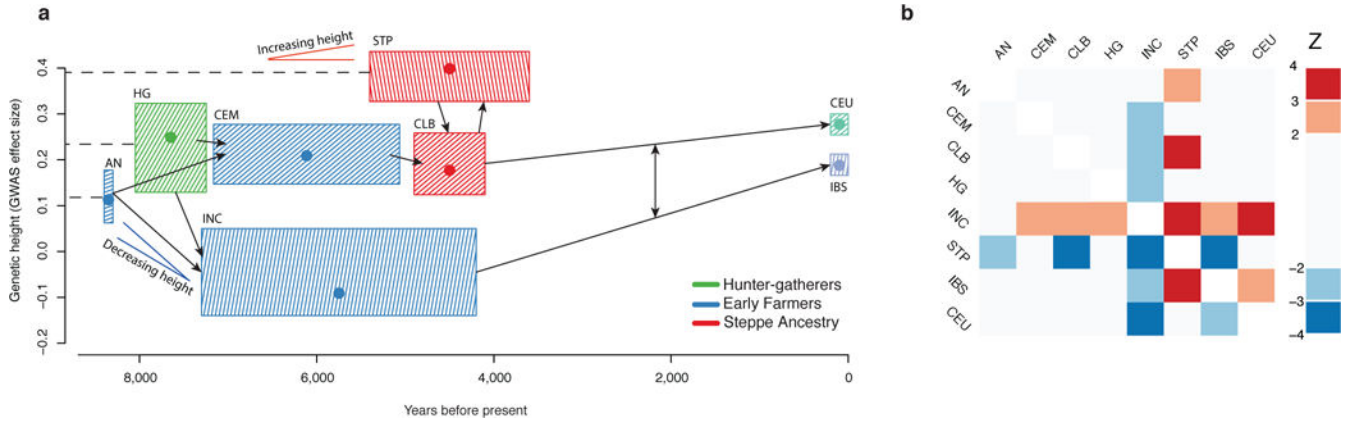
**Figure 2. Genome-wide scan for selection**

GC-corrected  $-\log_{10}$  p-value for each marker (Methods). The red dashed line represents a genome-wide significance level of  $0.5 \times 10^{-8}$ . Genome-wide significant points filtered because there were fewer than two other genome-wide significant points within 1Mb are shown in grey. Inset: QQ plots for corrected  $-\log_{10}$  P-values for different categories of potentially functional SNPs (Methods). Truncated at  $-\log_{10}(\text{P-value})=30$ . All curves are significantly different from neutral expectation.



**Figure 3. Allele frequencies for five genome-wide significant signals of selection**

Dots and solid lines show maximum likelihood frequency estimates and a 1.9-log-likelihood support interval for the derived allele frequency in each ancient population. Horizontal dashed lines show allele frequencies in the four modern 1000 Genomes populations. Abbreviations for ancient populations; AEN: Anatolian Neolithic; HG: hunter-gatherer; CEM: Central European Early and Middle Neolithic; INC: Iberian Neolithic and Chalcolithic; CLB: Central European Late Neolithic and Bronze Age; STP: Steppe. The Hunter-Gatherer, Early Farmer and Steppe Ancestry classifications correspond approximately to the three populations used in the genome-wide scan with some differences (See Extended Data Table 1 for details).



**Figure 4. Polygenic selection on height**

**A.** Estimated genetic heights. Boxes show 0.05–0.95 posterior densities for population mean genetic height (Methods). Dots show the maximum likelihood point estimate. Arrows show major population relationships, dashed lines represent ancestral populations. V’s show potential independent selection events. **B:** Z scores for the pairwise polygenic selection test. Positive if the column population is taller than the row population. Abbreviations; AN: Anatolian Neolithic; HG: hunter-gatherer; CEM: Central European Early and Middle Neolithic; INC: Iberian Neolithic and Chalcolithic; CLB: Central European Late Neolithic and Bronze Age; STP: Steppe; CEU: Utah residents with northern and western European ancestry; IBS: Iberian population in Spain.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

## Extended Data Table 1

## 230 ancient individuals analyzed in this study

**Population:** samples grouped by a combination of date, source, archaeology and genetics. **Date range:** approximate date range of samples in this group. **N:** Number of individuals sampled. **Out:** Number of PCA outliers (marked with an asterisk if used in selection analysis). **Rel:** Number of related individuals removed. **Eff N Chr:** Average over sites of the effective number of chromosomes when we use genotype likelihoods. Computed as 2 per called site for samples with genotype calls, or  $2 - 0.5^{(c-1)}$  for samples with read depth  $c$ . **Selection population 1:** Coarse population labels (marked with a caret if not used in genome-wide scan) **Selection population 2:** Fine population labels. **Abbreviations:** [E/M/L]N Early/Middle/Late Neolithic, LBK *Linearbandkeramik*, [E/S/W]HG Eastern/Scandinavian/Western hunter-gatherer, [E]BA [Early] Bronze Age, IA Iron Age.

By population	Population	Date range	N	Out	Rel	Eff N Chr	Selection population 1	Selection population 2
	WHG	8.2–8.0 kya	3	0	0	4.66	HG	HG
	Motala_HG	7.9–7.5 kya	6	0	0	5.19	HG	HG
	Anatolia_Neolithic	8.4–8.3 kya	24	1	1	22.49	EF	AN
	Hungary_EN	7.7–7.7 kya	10	0	0	8.81	EF	CEM
	LBK_EN	7.5–7.1 kya	15	0	0	11.15	EF	CEM
	Central_MN	5.9–5.8 kya	6	0	0	3.66	EF	CEM
	Iberia_EN	7.3–7.2 kya	4	0	1	3.54	EF	INC
	Iberia_MN	5.9–5.6 kya	4	0	0	3.47	EF	INC
	Iberia_Chalcolithic	4.8–4.2 kya	12	0	2	5.93	EF	INC
	Remedello	5.5–5.1 kya	3	0	0	0.93	EF	–
	Iceman	5.4–5.1 kya	1	0	0	1.90	EF	–
	Central_LNBA	4.9–4.6 kya	35	1	2	17.55	SA	CLB
	Yamnaya_Samara	5.4–4.9 kya	9	0	0	6.55	SA	STP
	Yamnaya_Kalmykia	5.3–4.7 kya	6	0	0	3.50	SA	STP
	Afanasievo	5.3–5.0 kya	5	0	0	3.01	SA	STP
	Poltavka	4.9–4.7 kya	4	1*	0	4.28	SA	STP
	Sintashta	4.3–4.1 kya	5	0	0	2.35	SA	STP
	Potapovka	4.2–4.1 kya	3	0	0	0.66	SA	STP
	Srubnaya	3.9–3.6 kya	12	1*	1	7.68	SA	STP
	Andronovo	3.8–3.6 kya	3	1*	0	3.87	SA	STP
	Russia_EBA	4.9–4.5 kya	1	0	0	0.21	SA	–
	Northern_LNBA	4.9–4.5 kya	10	0	0	3.81	SA	–

By population	Population	Date range	N	Out	Rel	Eff N	Chr	Selection population 1	Selection population 2
	Bell_Becker_LN	4.5–4.5 kya	17	0	1	6.64	SA <sup>^</sup>	SA <sup>^</sup>	CLB
	Hungary_BA	4.2–4.1 kya	12	0	0	4.18	SA <sup>^</sup>	SA <sup>^</sup>	CLB
	EHG	7.7–7.6 kya	3	0	0	2.15	–	–	–
	Samara_Eneolithic	7.2–6.0 kya	3	0	0	1.07	–	–	–
	Scythian_IA	2.4–2.2 kya	1	0	0	1.26	–	–	–
By selection population	Selection population 1	Date range	N	Out	Rel	Eff N	Chr	Description	
	EF	8.4–4.2 kya	79	0	0	61.88		Early Farmer	
	HG	8.2–7.5 kya	9	0	0	9.85		Hunter-gatherer	
	SA	5.4–3.6 kya	93	3	0	52.14		Steppe Ancestry	
By selection population 2	Selection population 2	Date range	N	Out	Rel	Eff N	Chr	Description	
	AN	8.4–8.3 kya	24	0	0	22.49		Anatolian Neolithic	
	CEM	7.7–5.8 kya	31	0	0	23.62		Central European Early and Middle Neolithic	
	INC	7.3–4.2 kya	20	0	0	12.95		Iberian Neolithic and Chalcolithic	
	HG	8.2–7.5 kya	9	0	0	9.85		Hunter-gatherer	
	CLB	4.9–4.1 kya	64	0	0	28.38		Central European Late Neolithic and Bronze Age	
	STP	5.4–3.6 kya	47	3	0	30.58		Steppe	



Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

**Twelve genome-wide significant signals of selection**

**Extended Data Table 3**

**Chromosome/Position/Range:** Co-ordinates (hg19) of the SNP with the most significant signal, and the approximate range in which genome-wide significant SNPs are found. **Genes:** Genes in which the top SNP is located, and selected nearby genes. **Potential function:** Function of the gene, or specific trait under selection. Marked with an asterisk if the signal was still genome-wide significant in an analysis that

used only the populations that correspond best to the three ancestral populations (WHG, Anatolian Neolithic and Bronze Age steppe), resulting in a less powerful test with the effective number of chromosomes analyzed at the average SNP reduced from 1.25 to 50, a genomic control correction of 1.32, and five genome-wide significant loci that are a subset of the original twelve.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

