

# Massive migration from the steppe was a source for Indo-European languages in Europe

Wolfgang Haak<sup>1\*</sup>, Iosif Lazaridis<sup>2,3\*</sup>, Nick Patterson<sup>3</sup>, Nadin Rohland<sup>2,3</sup>, Swapan Mallick<sup>2,3,4</sup>, Bastien Llamas<sup>1</sup>, Guido Brandt<sup>5</sup>, Susanne Nordenfelt<sup>2,3</sup>, Eadaoin Harney<sup>2,3,4</sup>, Kristin Stewardson<sup>2,3,4</sup>, Qiaomei Fu<sup>2,3,6,7</sup>, Alissa Mittnik<sup>8</sup>, Eszter Bánffy<sup>9,10</sup>, Christos Economou<sup>11</sup>, Michael Francken<sup>12</sup>, Susanne Friederich<sup>13</sup>, Rafael Garrido Pena<sup>14</sup>, Fredrik Hallgren<sup>15</sup>, Valery Khartanovich<sup>16</sup>, Aleksandr Khokhlov<sup>17</sup>, Michael Kunst<sup>18</sup>, Pavel Kuznetsov<sup>17</sup>, Harald Meller<sup>13</sup>, Oleg Mochalov<sup>17</sup>, Vayacheslav Moiseyev<sup>16</sup>, Nicole Nicklisch<sup>5,13,19</sup>, Sandra L. Pichler<sup>20</sup>, Roberto Risch<sup>21</sup>, Manuel A. Rojo Guerra<sup>22</sup>, Christina Roth<sup>5</sup>, Anna Szécsényi-Nagy<sup>5,9</sup>, Joachim Wahl<sup>23</sup>, Matthias Meyer<sup>6</sup>, Johannes Krause<sup>8,12,24</sup>, Dorcas Brown<sup>25</sup>, David Anthony<sup>25</sup>, Alan Cooper<sup>1</sup>, Kurt Werner Alt<sup>5,13,19,20</sup> & David Reich<sup>2,3,4</sup>

**We generated genome-wide data from 69 Europeans who lived between 8,000–3,000 years ago by enriching ancient DNA libraries for a target set of almost 400,000 polymorphisms. Enrichment of these positions decreases the sequencing required for genome-wide ancient DNA analysis by a median of around 250-fold, allowing us to study an order of magnitude more individuals than previous studies<sup>1–8</sup> and to obtain new insights about the past. We show that the populations of Western and Far Eastern Europe followed opposite trajectories between 8,000–5,000 years ago. At the beginning of the Neolithic period in Europe, ~8,000–7,000 years ago, closely related groups of early farmers appeared in Germany, Hungary and Spain, different from indigenous hunter-gatherers, whereas Russia was inhabited by a distinctive population of hunter-gatherers with high affinity to a ~24,000-year-old Siberian<sup>6</sup>. By ~6,000–5,000 years ago, farmers throughout much of Europe had more hunter-gatherer ancestry than their predecessors, but in Russia, the Yamnaya steppe herders of this time were descended not only from the preceding eastern European hunter-gatherers, but also from a population of Near Eastern ancestry. Western and Eastern Europe came into contact ~4,500 years ago, as the Late Neolithic Corded Ware people from Germany traced ~75% of their ancestry to the Yamnaya, documenting a massive migration into the heartland of Europe from its eastern periphery. This steppe ancestry persisted in all sampled central Europeans until at least ~3,000 years ago, and is ubiquitous in present-day Europeans. These results provide support for a steppe origin<sup>9</sup> of at least some of the Indo-European languages of Europe.**

Genome-wide analysis of ancient DNA has emerged as a transformative technology for studying prehistory, providing information that is comparable in power to archaeology and linguistics. Realizing its promise, however, requires collecting genome-wide data from an adequate number of individuals to characterize population changes over time, which means not only sampling a succession of archaeological cultures<sup>2</sup>, but also multiple individuals per culture. To make analysis of large numbers of ancient DNA samples practical, we used in-solution hybridization capture<sup>10,11</sup> to enrich next generation sequencing libraries for a

target set of 394,577 single nucleotide polymorphisms (SNPs) ('390k capture'), 354,212 of which are autosomal SNPs that have also been genotyped using the Affymetrix Human Origins array in 2,345 humans from 203 populations<sup>4,12</sup>. This reduces the amount of sequencing required to obtain genome-wide data by a minimum of 45-fold and a median of 262-fold (Supplementary Data 1). This strategy allows us to report genomic scale data on more than twice the number of ancient Eurasians as has been presented in the entire preceding literature<sup>1–8</sup> (Extended Data Table 1).

We used this technology to study population transformations in Europe. We began by preparing 212 DNA libraries from 119 ancient samples in dedicated clean rooms, and testing these by light shotgun sequencing and mitochondrial genome capture (Supplementary Information section 1, Supplementary Data 1). We restricted the analysis to libraries with molecular signatures of authentic ancient DNA (elevated damage in the terminal nucleotide), negligible evidence of contamination based on mismatches to the mitochondrial consensus<sup>13</sup> and, where available, a mitochondrial DNA haplogroup that matched previous results using PCR<sup>4,14,15</sup> (Supplementary Information section 2). For 123 libraries prepared in the presence of uracil-DNA-glycosylase<sup>16</sup> to reduce errors due to ancient DNA damage<sup>17</sup>, we performed 390k capture, carried out paired-end sequencing and mapped the data to the human genome. We restricted analysis to 94 libraries from 69 samples that had at least 0.06-fold average target coverage (average of 3.8-fold) and used majority rule to call an allele at each SNP covered at least once (Supplementary Data 1). After combining our data (Supplementary Information section 3) with 25 ancient samples from the literature — three Upper Paleolithic samples from Russia<sup>1,6,7</sup>, seven people of European hunter-gatherer ancestry<sup>2,4,5,8</sup>, and fifteen European farmers<sup>2,3,4,8</sup> — we had data from 94 ancient Europeans. Geographically, these came from Germany ( $n = 41$ ), Spain ( $n = 10$ ), Russia ( $n = 14$ ), Sweden ( $n = 12$ ), Hungary ( $n = 15$ ), Italy ( $n = 1$ ) and Luxembourg ( $n = 1$ ) (Extended Data Table 2). Following the central European chronology, these included 19 hunter-gatherers (~43,000–2,600 BC), 28 Early Neolithic farmers (~6,000–4,000 BC), 11 Middle Neolithic farmers (~4,000–3,000 BC) including

<sup>1</sup>Australian Centre for Ancient DNA, School of Earth and Environmental Sciences & Environment Institute, University of Adelaide, Adelaide, South Australia 5005, Australia. <sup>2</sup>Department of Genetics, Harvard Medical School, Boston, Massachusetts 02115, USA. <sup>3</sup>Broad Institute of Harvard and MIT, Cambridge, Massachusetts 02142, USA. <sup>4</sup>Howard Hughes Medical Institute, Harvard Medical School, Boston, Massachusetts 02115, USA. <sup>5</sup>Institute of Anthropology, Johannes Gutenberg University of Mainz, D-55128 Mainz, Germany. <sup>6</sup>Max Planck Institute for Evolutionary Anthropology, D-04103 Leipzig, Germany. <sup>7</sup>Key Laboratory of Vertebrate Evolution and Human Origins of Chinese Academy of Sciences, IVPP, CAS, Beijing 100049, China. <sup>8</sup>Institute for Archaeological Sciences, University of Tübingen, D-72070 Tübingen, Germany. <sup>9</sup>Institute of Archaeology, Research Centre for the Humanities, Hungarian Academy of Science, H-1014 Budapest, Hungary. <sup>10</sup>Römisch Germanische Kommission (RGK) Frankfurt, D-60325 Frankfurt, Germany. <sup>11</sup>Archaeological Research Laboratory, Stockholm University, 114 18 Stockholm, Sweden. <sup>12</sup>Departments of Paleoanthropology and Archaeogenetics, Senckenberg Center for Human Evolution and Paleoenvironment, University of Tübingen, D-72070 Tübingen, Germany. <sup>13</sup>State Office for Heritage Management and Archaeology Saxony-Anhalt and State Museum of Prehistory, D-06114 Halle, Germany. <sup>14</sup>Departamento de Prehistoria y Arqueología, Facultad de Filosofía y Letras, Universidad Autónoma de Madrid, E-28049 Madrid, Spain. <sup>15</sup>The Cultural Heritage Foundation, Västerås 722 12, Sweden. <sup>16</sup>Peter the Great Museum of Anthropology and Ethnography (Kunstkamera) RAS, St Petersburg 199034, Russia. <sup>17</sup>Volga State Academy of Social Sciences and Humanities, Samara 443099, Russia. <sup>18</sup>Deutsches Archäologisches Institut, Abteilung Madrid, E-28002 Madrid, Spain. <sup>19</sup>Danube Private University, A-3500 Krems, Austria. <sup>20</sup>Institute for Prehistory and Archaeological Science, University of Basel, CH-4003 Basel, Switzerland. <sup>21</sup>Departamento de Prehistoria, Universitat Autònoma de Barcelona, E-08193 Barcelona, Spain. <sup>22</sup>Departamento de Prehistoria y Arqueología, Universidad de Valladolid, E-47002 Valladolid, Spain. <sup>23</sup>State Office for Cultural Heritage Management Baden-Württemberg, Osteology, D-78467 Konstanz, Germany. <sup>24</sup>Max Planck Institute for the Science of Human History, D-07745 Jena, Germany. <sup>25</sup>Anthropology Department, Hartwick College, Oneonta, New York 13820, USA.

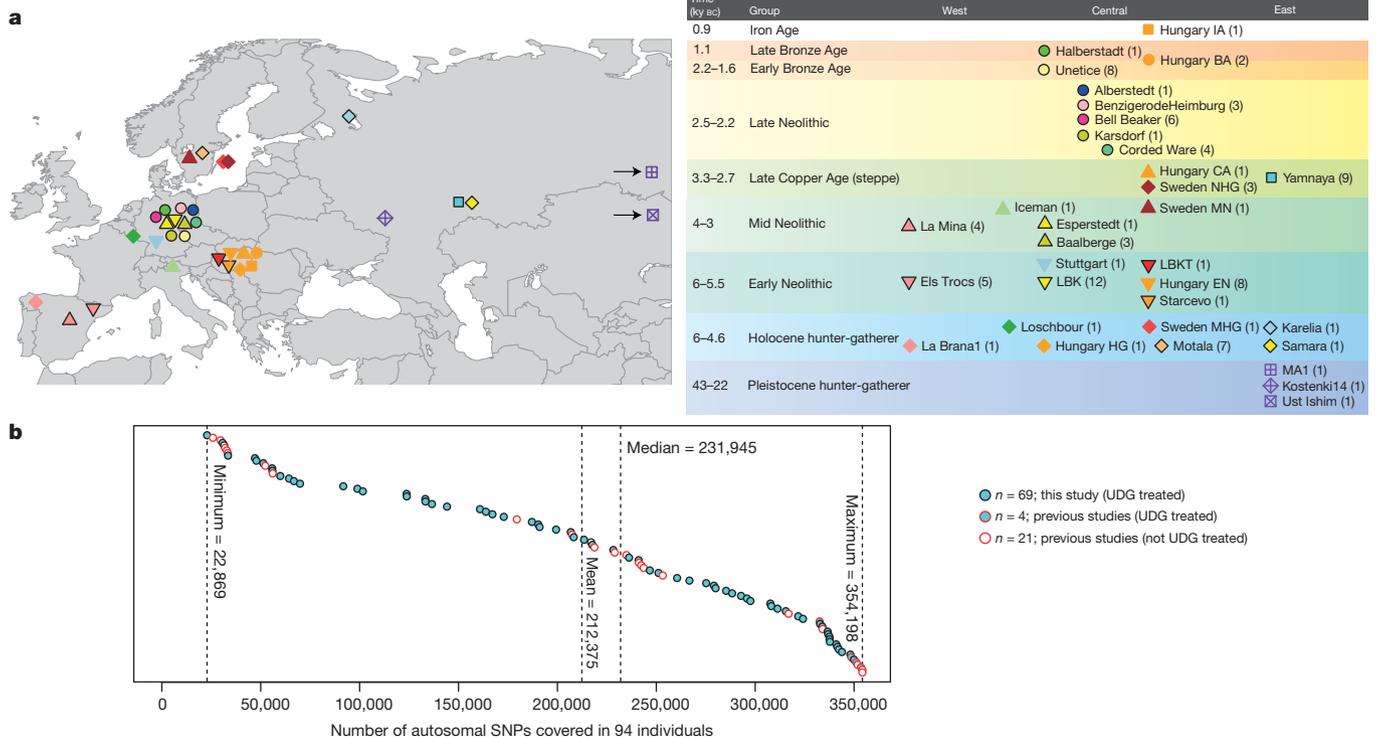
\*These authors contributed equally to this work.

the Tyrolean Iceman<sup>3</sup>, 9 Late Copper/Early Bronze Age individuals (Yamnaya: ~3,300–2,700 BC), 15 Late Neolithic individuals (~2,500–2,200 BC), 9 Early Bronze Age individuals (~2,200–1,500 BC), two Late Bronze Age individuals (~1,200–1,100 BC) and one Iron Age individual (~900 BC). Two individuals were excluded from analyses as they were related to others from the same population. The average number of SNPs covered at least once was 212,375 and the minimum was 22,869 (Fig. 1).

We determined that 34 of the 69 newly analysed individuals were male and used 2,258 Y chromosome SNP targets included in the capture to obtain high resolution Y chromosome haplogroup calls (Supplementary Information section 4). **Outside Russia, and before the Late Neolithic period, only a single R1b individual was found** (early Neolithic Spain) in the combined literature ( $n = 70$ ). **By contrast, haplogroups R1a and R1b were found in 60% of Late Neolithic/Bronze Age Europeans outside Russia ( $n = 10$ ), and in 100% of the samples from European Russia from all periods (7,500–2,700 BC;  $n = 9$ ). R1a and R1b are the most common haplogroups in many European populations today<sup>18,19</sup>, and our results suggest that they spread into Europe from the East after 3,000 BC.** Two hunter-gatherers from Russia included in our study belonged to R1a (Karelia) and R1b (Samara), the earliest documented ancient samples of either haplogroup discovered to date. These two hunter-gatherers did not belong to the derived lineages M417 within R1a and M269 within R1b that are predominant in Europeans today<sup>18,19</sup>, but all 7 Yamnaya males did belong to the M269 subclade<sup>18</sup> of haplogroup R1b.

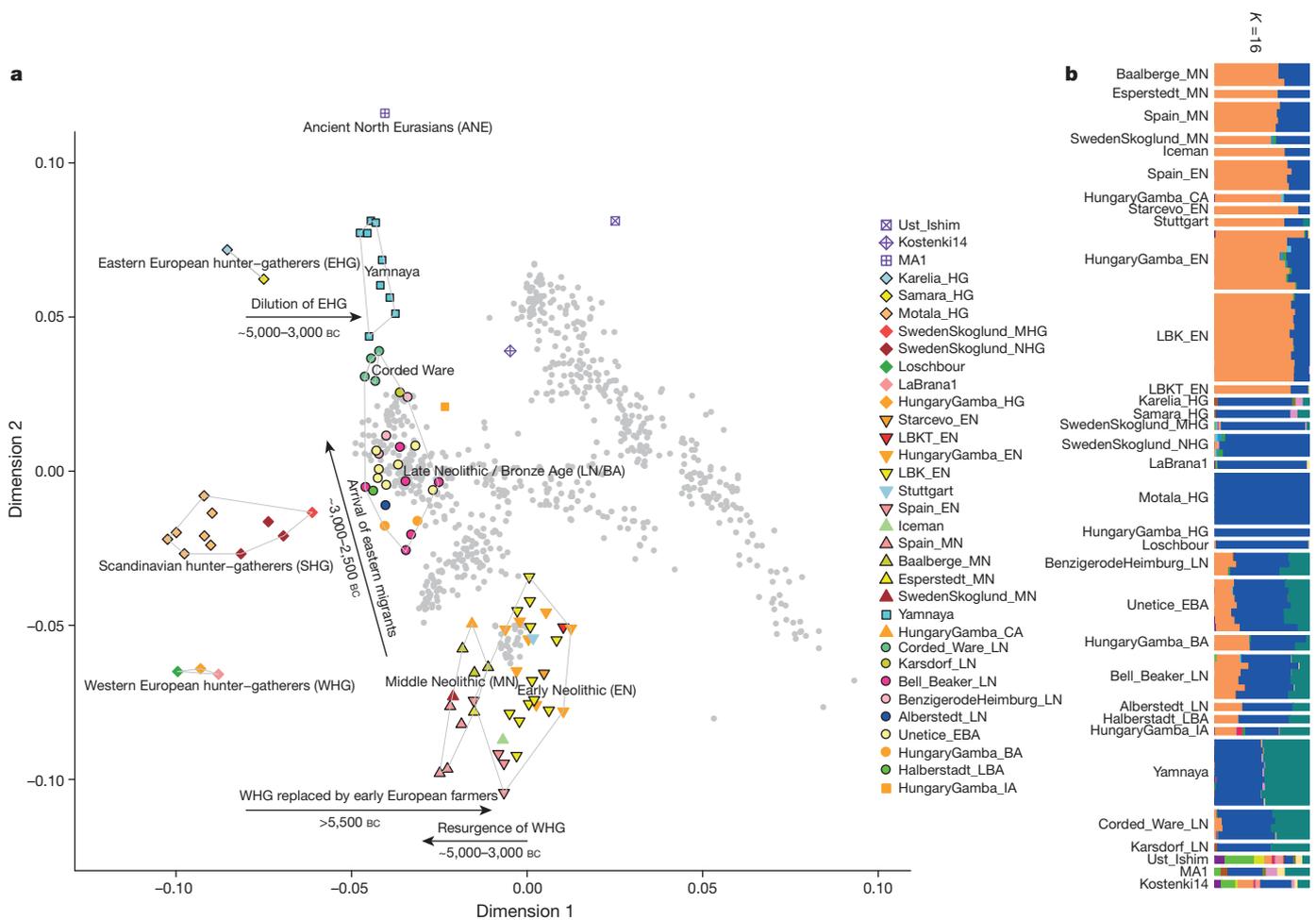
Principal components analysis (PCA) of all ancient individuals along with 777 present-day West Eurasians<sup>4</sup> (Fig. 2a, Supplementary Information section 5) replicates the **positioning of present-day Europeans between the Near East and European hunter-gatherers<sup>4,20</sup>**, and the clustering of early farmers from across Europe with present day Sardinians<sup>3,4</sup>, suggesting that **farming expansions across the Mediterranean to Spain and via the Danubian route to Hungary and Germany descended from a common stock.** By adding samples from later periods and additional locations, we also observe several new patterns. All samples from Russia have affinity to the ~24,000-year-old MA1 (ref. 6), the type specimen for

the Ancient North Eurasians (ANE) who contributed to both Europeans<sup>4</sup> and Native Americans<sup>4,6,8</sup>. The two hunter-gatherers from Russia (Karelia in the northwest of the country and Samara on the steppe near the Urals) form an ‘eastern European hunter-gatherer’ (EHG) cluster at one end of a hunter-gatherer cline across Europe; people of hunter-gatherer ancestry from Luxembourg, Spain, and Hungary sit at the opposite ‘western European hunter-gatherer’<sup>24</sup> (WHG) end, while the hunter-gatherers from Sweden<sup>4,8</sup> (SHG) are intermediate. **Against this background of differentiated European hunter-gatherers and homogeneous early farmers, multiple population turnovers transpired in all parts of Europe included in our study. Middle Neolithic Europeans from Germany, Spain, Hungary, and Sweden from the period ~4,000–3,000 BC are intermediate between the earlier farmers and the WHG, suggesting an increase of WHG ancestry throughout much of Europe. By contrast, in Russia, the later Yamnaya steppe herders of ~3,000 BC plot between the EHG and the present-day Near East/Caucasus, suggesting a decrease of EHG ancestry during the same time period.** The Late Neolithic and Bronze Age samples from Germany and Hungary<sup>7</sup> are distinct from the preceding Middle Neolithic and plot between them and the Yamnaya. This pattern is also seen in ADMIXTURE analysis (Fig. 2b, Supplementary Information section 6), which implies that the Yamnaya have ancestry from populations related to the Caucasus and South Asia that is largely absent in 38 Early or Middle Neolithic farmers but present in all 25 Late Neolithic or Bronze Age individuals. This ancestry appears in Central Europe for the first time in our series with the Corded Ware around 2,500 BC (Supplementary Information section 6, Fig. 2b). **The Corded Ware shared elements of material culture with steppe groups such as the Yamnaya although whether this reflects movements of people has been contentious<sup>21</sup>. Our genetic data provide direct evidence of migration and suggest that it was relatively sudden.** The Corded Ware are genetically closest to the Yamnaya ~2,600 km away, as inferred both from PCA and ADMIXTURE (Fig. 2) and  $F_{ST}$  (0.011 ± 0.002) (Extended Data Table 3). If continuous gene flow from the east, rather than migration, had occurred, we would expect successive cultures in Europe to become increasingly differentiated from the Middle Neolithic, but



**Figure 1 | Location and SNP coverage of samples included in this study.** **a**, Geographic location and time-scale (central European chronology) of the 69 newly analysed ancient individuals from this study (black outline) and 25 from

the literature for which shotgun sequencing data was available (no outline). **b**, Number of SNPs covered at least once in the analysis data set of 94 individuals.



**Figure 2 | Population transformations in Europe.** **a**, PCA analysis. **b**, ADMIXTURE analysis. The full ADMIXTURE analysis including present-day humans is shown in Supplementary Information section 6.

instead, the Corded Ware are both the earliest and most strongly differentiated from the Middle Neolithic population.

'Outgroup'  $f_3$  statistics<sup>6</sup> (Supplementary Information section 7), which measure shared genetic drift between a pair of populations (Extended Data Fig. 1), support the clustering of hunter-gatherers, Early/Middle Neolithic, and Late Neolithic/Bronze Age populations into different groups as in the PCA (Fig. 2a). We also analysed  $f_4$  statistics, which allow us to test whether pairs of populations are consistent with descent from common ancestral populations, and to assess significance using a normally distributed  $Z$  score. Early European farmers from the Early and Middle Neolithic were closely related but not identical. This is reflected in the fact that Loschbour, a WHG individual from Luxembourg<sup>4</sup>, shared more alleles with post-4,000 BC European farmers from Germany, Spain, Hungary, Sweden and Italy than with early farmers of Germany, Spain, and Hungary, documenting an increase of hunter-gatherer ancestry in multiple regions of Europe during the course of the Neolithic. The two EHG form a clade with respect to all other present-day and ancient populations ( $|Z| < 1.9$ ), and MA1 shares more alleles with them ( $|Z| > 4.7$ ) than with other ancient or modern populations, suggesting that they may be a source for the ANE ancestry in present Europeans<sup>4,12,22</sup> as they are geographically and temporally more proximate than Upper Paleolithic Siberians. The Yamnaya differ from the EHG by sharing fewer alleles with MA1 ( $|Z| = 6.7$ ) suggesting a dilution of ANE ancestry between 5,000–3,000 BC on the European steppe. This was likely due to admixture of EHG with a population related to present-day Near Easterners, as the most negative  $f_3$  statistic in the Yamnaya (giving unambiguous evidence of admixture) is observed when we model them as a mixture of EHG and present-day Near Eastern populations like Armenians ( $Z = -6.3$ ;

Supplementary Information section 7). The Late Neolithic/Bronze Age groups of central Europe share more alleles with Yamnaya than the Middle Neolithic populations do ( $|Z| = 12.4$ ) and more alleles with the Middle Neolithic than the Yamnaya do ( $|Z| = 12.5$ ), and have a negative  $f_3$  statistic with the Middle Neolithic and Yamnaya as references ( $Z = -20.7$ ), indicating that they were descended from a mixture of the local European populations and new migrants from the east. Moreover, the Yamnaya share more alleles with the Corded Ware ( $|Z| \geq 3.6$ ) than with any other Late Neolithic/Early Bronze Age group with at least two individuals (Supplementary Information section 7), indicating that they had more eastern ancestry, consistent with the PCA and ADMIXTURE patterns (Fig. 2).

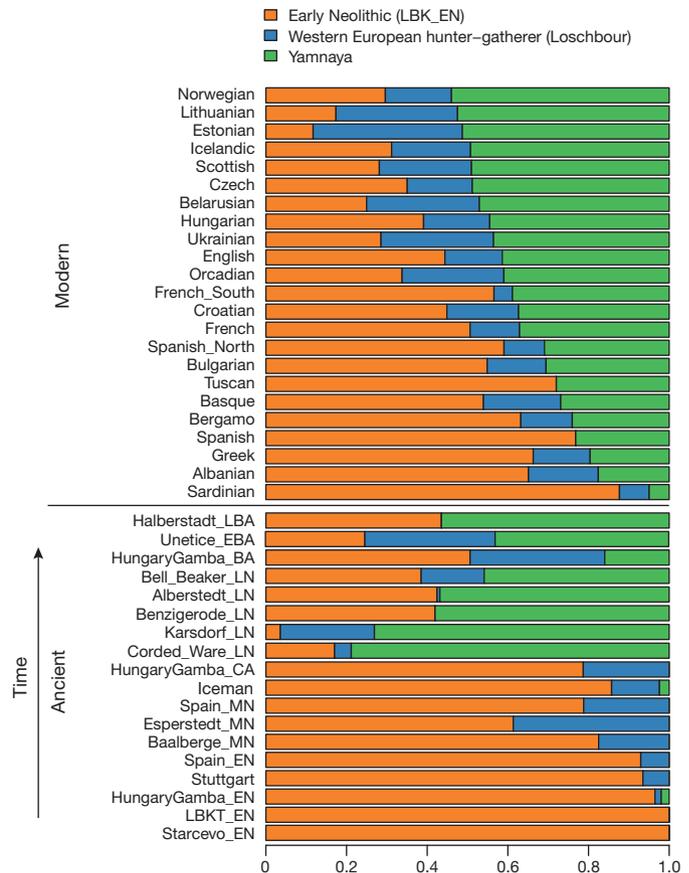
Modelling of the ancient samples shows that while Karelia is genetically intermediate between Loschbour and MA1, the topology that considers Karelia as a mixture of these two elements is not the only one that can fit the data (Supplementary Information section 8). To avoid biasing our inferences by fitting an incorrect model, we developed new statistical methods that are substantial extensions of a previously reported approach<sup>4</sup>, which allow us to obtain precise estimates of the proportion of mixture in later Europeans without requiring a formal model for the relationship among the ancestral populations. The method (Supplementary Information section 9) is based on the idea that if a Test population has ancestry related to reference populations  $Ref_1, Ref_2, \dots, Ref_N$  in proportions  $\alpha_1, \alpha_2, \dots, \alpha_N$ , and the references are themselves differentially related to a triple of outgroup populations  $A, B, C$ , then:

$$f_4(\text{Test}, A; B, C) = \sum_{i=1}^N \alpha_i f_4(\text{Ref}_i, A; B, C)$$

By using a large number of outgroup populations we can fit the admixture coefficients  $\alpha_i$  and estimate mixture proportions (Supplementary Information section 9, Extended Data Fig. 2). Using 15 outgroups from Africa, Asia, Oceania and the Americas, we obtain good fits as assessed by a formal test (Supplementary Information section 10), and estimate that the Middle Neolithic populations of Germany and Spain have ~18–34% more WHG-related ancestry than Early Neolithic populations and that the Late Neolithic and Early Bronze Age populations of Germany have ~22–39% more EHG-related ancestry than the Middle Neolithic ones (Supplementary Information section 9). If we model them as mixtures of Yamnaya-related and Middle Neolithic populations, the inferred degree of population turnover is doubled to 48–80% (Supplementary Information sections 9 and 10).

To distinguish whether a Yamnaya or an EHG source fits the data better, we added ancient samples as outgroups (Supplementary Information section 9). Adding any Early or Middle Neolithic farmer results in EHG-related genetic input into Late Neolithic populations being a poor fit to the data (Supplementary Information section 9); thus, Late Neolithic populations have ancestry that cannot be explained by a mixture of EHG and Middle Neolithic. When using Yamnaya instead of EHG, however, we obtain a good fit (Supplementary Information sections 9 and 10). These results can be explained if the new genetic material that arrived in Germany was a composite of two elements: EHG and a type of Near Eastern ancestry different from that which was introduced by early farmers (also suggested by PCA and ADMIXTURE; Fig. 2, Supplementary Information sections 5 and 6). We estimate that these two elements each contributed about half the ancestry of the Yamnaya (Supplementary Information sections 6 and 9), explaining why the population turnover inferred using Yamnaya as a source is about twice as high compared to the undiluted EHG. The estimate of Yamnaya-related ancestry in the Corded Ware is consistent when using either present populations or ancient Europeans as outgroups (Supplementary Information sections 9 and 10), and is  $73.1 \pm 2.2\%$  when both sets are combined (Supplementary Information section 10). The best proxies for ANE ancestry in Europe<sup>4</sup> were initially Native Americans<sup>12,22</sup>, and then the Siberian MA1 (ref. 6), but both are geographically and temporally too remote for what appears to be a recent migration into Europe<sup>4</sup>. We can now add three new pieces to the puzzle of how ANE ancestry was transmitted to Europe: first by the EHG, then the Yamnaya formed by mixture between EHG and a Near Eastern related population, and then the Corded Ware who were formed by a mixture of the Yamnaya with Middle Neolithic Europeans. We caution that the sampled Yamnaya individuals from Samara might not be directly ancestral to Corded Ware individuals from Germany. It is possible that a more western Yamnaya population, or an earlier (pre-Yamnaya) steppe population may have migrated into central Europe, and future work may uncover more missing links in the chain of transmission of steppe ancestry.

By extending our model to a three-way mixture of WHG, Early Neolithic and Yamnaya, we estimate that the ancestry of the Corded Ware was 79% Yamnaya-like, 4% WHG, and 17% Early Neolithic (Fig. 3). A small contribution of the first farmers is also consistent with uniparentally inherited DNA: for example, mitochondrial DNA haplogroup N1a and Y chromosome haplogroup G2a, common in early central European farmers<sup>14,23</sup>, almost disappear during the Late Neolithic and Bronze Age, when they are largely replaced by Y haplogroups R1a and R1b (Supplementary Information section 4) and mtDNA haplogroups I, T1, U2, U4, U5a, W, and subtypes of H<sup>14,23,24</sup> (Supplementary Information section 2). The uniparental data not only confirm a link to the steppe populations but also suggest that both sexes participated in the migrations (Supplementary Information sections 2 and 4 and Extended Data Table 2). The magnitude of the population turnover that occurred becomes even more evident if one considers the fact that the steppe migrants may well have mixed with eastern European agriculturalists on their way to central Europe. Thus, we cannot exclude a scenario in which the Corded Ware arriving in today's Germany had no ancestry at all from local populations.



**Figure 3 | Admixture proportions.** We estimate mixture proportions using a method that gives unbiased estimates even without an accurate model for the relationships between the test populations and the outgroup populations (Supplementary Information section 9). Population samples are grouped according to chronology (ancient) and Yamnaya ancestry (present-day humans).

Our results support a view of European pre-history punctuated by two major migrations: first, the arrival of the first farmers during the Early Neolithic from the Near East, and second, the arrival of Yamnaya pastoralists during the Late Neolithic from the steppe. Our data further show that both migrations were followed by resurgences of the previous inhabitants: first, during the Middle Neolithic, when hunter-gatherer ancestry rose again after its Early Neolithic decline, and then between the Late Neolithic and the present, when farmer and hunter-gatherer ancestry rose after its Late Neolithic decline. This second resurgence must have started during the Late Neolithic/Bronze Age period itself, as the Bell Beaker and Unetice groups had reduced Yamnaya ancestry compared to the earlier Corded Ware, and comparable levels to that in some present-day Europeans (Fig. 3). Today, Yamnaya related ancestry is lower in southern Europe and higher in northern Europe, and all European populations can be modelled as a three-way mixture of WHG, Early Neolithic, and Yamnaya, whereas some outlier populations show evidence for additional admixture with populations from Siberia and the Near East (Extended Data Fig. 3, Supplementary Information section 9). Further data are needed to determine whether the steppe ancestry arrived in southern Europe at the time of the Late Neolithic/Bronze Age, or is due to migrations in later times from northern Europe<sup>25,26</sup>.

Our results provide new data relevant to debates on the origin and expansion of Indo-European languages in Europe (Supplementary Information section 11). Although the findings from ancient DNA are silent on the question of the languages spoken by preliterate populations, they do carry evidence about processes of migration which are invoked by theories on Indo-European language dispersals. Such theories make predictions about movements of people to account for the spread of

languages and material culture (Extended Data Fig. 4). The technology of ancient DNA makes it possible to reject or confirm the proposed migratory movements, as well as to identify new movements that were not previously known. The best argument for the ‘Anatolian hypothesis’<sup>27</sup> that Indo-European languages arrived in Europe from Anatolia ~8,500 years ago is that major language replacements are thought to require major migrations, and that after the Early Neolithic when farmers established themselves in Europe, the population base was likely to have been so large that later migrations would not have made much of an impact<sup>27,28</sup>. However, our study shows that a later major turnover did occur, and that steppe migrants replaced ~75% of the ancestry of central Europeans. An alternative theory is the ‘steppe hypothesis’, which proposes that early Indo-European speakers were pastoralists of the grasslands north of the Black and Caspian Seas, and that their languages spread into Europe after the invention of wheeled vehicles<sup>9</sup>. Our results make a compelling case for the steppe as a source of at least some of the Indo-European languages in Europe by documenting a massive migration ~4,500 years ago associated with the Yamnaya and Corded Ware cultures, which are identified by proponents of the steppe hypothesis as vectors for the spread of Indo-European languages into Europe. These results challenge the Anatolian hypothesis by showing that not all Indo-European languages in Europe can plausibly derive from the first farmer migrations thousands of years earlier (Supplementary Information section 11). We caution that the location of the proto-Indo-European<sup>9,27,29,30</sup> homeland that also gave rise to the Indo-European languages of Asia, as well as the Indo-European languages of southeastern Europe, cannot be determined from the data reported here (Supplementary Information section 11). Studying the mixture in the Yamnaya themselves, and understanding the genetic relationships among a broader set of ancient and present-day Indo-European speakers, may lead to new insight about the shared homeland.

**Online Content** Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 29 December 2014; accepted 12 February 2015.

Published online 2 March 2015.

- Fu, Q. *et al.* Genome sequence of a 45,000-year-old modern human from western Siberia. *Nature* **514**, 445–449 (2014).
- Gamba, C. *et al.* Genome flux and stasis in a five millennium transect of European prehistory. *Nature Commun.* **5**, 5257 (2014).
- Keller, A. *et al.* New insights into the Tyrolean Iceman’s origin and phenotype as inferred by whole-genome sequencing. *Nature Commun.* **3**, 698 (2012).
- Lazaridis, I. *et al.* Ancient human genomes suggest three ancestral populations for present-day Europeans. *Nature* **513**, 409–413 (2014).
- Olalde, I. *et al.* Derived immune and ancestral pigmentation alleles in a 7,000-year-old Mesolithic European. *Nature* **507**, 225–228 (2014).
- Raghavan, M. *et al.* Upper Palaeolithic Siberian genome reveals dual ancestry of Native Americans. *Nature* **505**, 87–91 (2014).
- Seguin-Orlando, A. *et al.* Genomic structure in Europeans dating back to at least 36,200 years. *Science* **346**, 1113–1118 (2014).
- Skoglund, P. *et al.* Genomic diversity and admixture differs for Stone-Age Scandinavian foragers and farmers. *Science* **344**, 747–750 (2014).
- Anthony, D. W. *The Horse, the Wheel, and Language: How Bronze-Age Riders from the Eurasian Steppes Shaped the Modern World* (Princeton Univ. Press, 2007).
- Fu, Q. *et al.* DNA analysis of an early modern human from Tianyuan Cave, China. *Proc. Natl Acad. Sci. USA* **110**, 2223–2227 (2013).
- Rohland, N., Harney, E., Mallick, S., Nordenfelt, S. & Reich, D. Partial uracil–DNA-glycosylase treatment for screening of ancient DNA. *Phil. Trans. R. Soc. Lond. B* **370**, 20130624 (2015).
- Patterson, N. *et al.* Ancient admixture in human history. *Genetics* **192**, 1065–1093 (2012).
- Fu, Q. *et al.* A revised timescale for human evolution based on ancient mitochondrial genomes. *Curr. Biol.* **23**, 553–559 (2013).
- Brandt, G. *et al.* Ancient DNA reveals key stages in the formation of central European mitochondrial genetic diversity. *Science* **342**, 257–261 (2013).
- Der Sarkissian, C. *et al.* Ancient DNA reveals prehistoric gene-flow from Siberia in the complex human population history of North East Europe. *PLoS Genet.* **9**, e1003296 (2013).
- Briggs, A. W. *et al.* Removal of deaminated cytosines and detection of *in vivo* methylation in ancient DNA. *Nucleic Acids Res.* **38**, e87 (2010).
- Briggs, A. W. *et al.* Patterns of damage in genomic DNA sequences from a Neandertal. *Proc. Natl Acad. Sci. USA* **104**, 14616–14621 (2007).
- Myres, N. M. *et al.* A major Y-chromosome haplogroup R1b Holocene era founder effect in Central and Western Europe. *Eur. J. Hum. Genet.* **19**, 95–101 (2011).
- Underhill, P. A. *et al.* The phylogenetic and geographic structure of Y-chromosome haplogroup R1a. *Eur. J. Hum. Genet.* **23**, 124–131 (2015).
- Skoglund, P. *et al.* Origins and genetic legacy of Neolithic farmers and hunter-gatherers in Europe. *Science* **336**, 466–469 (2012).
- Czebreszuk, J. in *Ancient Europe, 8000 B.C. to A.D. 1000: Encyclopedia of the Barbarian World* (eds Bogucki, P. I. & Crabtree, P. J.) 467–475 (Charles Scribners & Sons, 2003).
- Lipson, M. *et al.* Efficient moment-based inference of admixture parameters and sources of gene flow. *Mol. Biol. Evol.* **30**, 1788–1802 (2013).
- Szécényi-Nagy, A. *et al.* Tracing the genetic origin of Europe’s first farmers reveals insights into their social organization. Preprint at *bioRxiv* <http://dx.doi.org/10.1101/008664> (2014).
- Haak, W. *et al.* Ancient DNA from European early Neolithic farmers reveals their Near Eastern affinities. *PLoS Biol.* **8**, e1000536 (2010).
- Hellenthal, G. *et al.* A genetic atlas of human admixture history. *Science* **343**, 747–751 (2014).
- Ralph, P. & Coop, G. The geography of recent genetic ancestry across Europe. *PLoS Biol.* **11**, e1001555 (2013).
- Renfrew, C. *Archaeology and Language: The Puzzle of Indo-European Origins* (Pimlico, 1987).
- Bellwood, P. *First Farmers: The Origins of Agricultural Societies* (Wiley-Blackwell, 2004).
- Gamkrelidze, T. V. & Ivanov, V. V. The early history of Indo-European languages. *Sci. Am.* **262**, 110–116 (1990).
- Mallory, J. P. In *Search of the Indo-Europeans: Language, Archaeology and Myth* (Thames and Hudson, 1991).

**Supplementary Information** is available in the online version of the paper.

**Acknowledgements** We thank P. Bellwood, J. Burger, P. Heggarty, M. Lipson, C. Renfrew, J. Diamond, S. Pääbo, R. Pinhasi and P. Skoglund for critical comments, and the Initiative for the Science of the Human Past at Harvard for organizing a workshop around the issues touched on by this paper. We thank S. Pääbo for support for establishing the ancient DNA facilities in Boston, and P. Skoglund for detecting the presence of two related individuals in our data set. We thank L. Orlando, T. S. Korneliusen, and C. Gamba for help in obtaining data. We thank Agilent Technologies and G. Frommer for help in developing the capture reagents. We thank C. Der Sarkissian, G. Valverde, L. Papac and B. Nickel for wet laboratory support. We thank archaeologists V. Dresely, R. Ganslmeier, O. Balanvosky, J. Ignacio Royo Guillén, A. Oszás, V. Majerik, T. Paluch, K. Somogyi and V. Voicsek for sharing samples and discussion about archaeological context. This research was supported by an Australian Research Council grant to W.H. and B.L. (DP130102158), and German Research Foundation grants to K.W.A. (AI 287/7-1 and 7-3, AI 287/10-1 and AI 287/14-1) and to H.M. (Me 3245/1-1 and 1-3). D.R. was supported by US National Science Foundation HOMINID grant BCS-1032255, US National Institutes of Health grant GM100233, and the Howard Hughes Medical Institute.

**Author Contributions** W.H., N.P., N.R., J.K., K.W.A. and D.R. supervised the study. W.H., E.B., C.E., M.F., S.F., R.G.P., F.H., V.K., A.K., M.K., P.K., H.M., O.M., V.M., N.N., S.L.P., R.R., M.A.R.G., C.R., A.S.-N., J.W., J.K., D.B., D.A., A.C., K.W.A. and D.R. assembled archaeological material, W.H., I.L., N.P., N.R., S.M., A.M. and D.R. analysed genetic data. I.L., N.P. and D.R. developed methods using *f* statistics for inferring admixture proportions. W.H., N.R., B.L., G.B., S.N., E.H., K.S. and A.M. performed wet laboratory ancient DNA work. I.L., N.R., S.M., B.L., Q.F., M.M. and D.R. developed the 390k capture reagent. W.H., I.L. and D.R. wrote the manuscript with help from all co-authors.

**Author Information** The aligned sequences are available through the European Nucleotide Archive under accession number PRJEB8448. The Human Origins genotype dataset 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)). Reprints and permissions information is available at [www.nature.com/reprints](http://www.nature.com/reprints). The authors declare no competing financial interests. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to D.R. ([reich@genetics.med.harvard.edu](mailto:reich@genetics.med.harvard.edu)).

## METHODS

**Screening of libraries (shotgun sequencing and mitochondrial capture).** The 212 libraries screened in this study (Supplementary Information section 1) from a total of 119 samples (Supplementary Information section 3) were produced at Adelaide ( $n = 151$ ), Tübingen ( $n = 16$ ), and Boston ( $n = 45$ ) (Supplementary Data 1).

The libraries from Adelaide and Boston had internal barcodes directly attached to both sides of the molecules from the DNA extract so that each sequence begins with the barcode<sup>10</sup>. The Adelaide libraries had 5 base pair (bp) barcodes on both sides, while the Boston libraries had 7 bp barcodes. Libraries from Tübingen had no internal barcodes, but were differentiated by the sequence of the indexing primer<sup>31</sup>.

We adapted a reported protocol for enriching for mitochondrial DNA<sup>10</sup>, with the difference that we adjusted the blocking oligonucleotides and PCR primers to fit our libraries with shorter adapters. Over the course of this project, we also lowered the hybridization temperature from 65 °C to 60 °C and performed stringent washes at 55 °C instead of 60 °C<sup>32</sup>.

We used an aliquot of approximately 500 ng of each library for target enrichment of the complete mitochondrial genome in two consecutive rounds<sup>32</sup>, using a bait set for human mtDNA<sup>32</sup>. We performed enrichment in 96-well plates with one library per well, and used a liquid handler (Evolution P3, Perkin Elmer) for the capture and washing steps<sup>33</sup>. We used blocking oligonucleotides in hybridization appropriate to the adapters of the truncated libraries. After either of the two enrichment rounds, we amplified the enriched library molecules with the primer pair that keeps the adapters short (PreHyb) using Herculase Fusion II PCR Polymerase. We performed an indexing PCR of the final capture product using one or two indexing primers<sup>31</sup>. We cleaned up all PCR reactions using SPRI technology<sup>34</sup> and the liquid handler. Libraries from Tübingen were amplified with the primer pair IS5/IS6<sup>31</sup>.

For libraries from Boston and Adelaide, we used a second aliquot of each library for shotgun sequencing after performing an indexing PCR<sup>31</sup>. We used unique index combinations for each library and experiment, allowing us to distinguish shotgun sequencing and mitochondrial DNA capture data, even if both experiments were in the same sequencing run. We sequenced shotgun libraries and mtDNA captured libraries from Tübingen in independent sequencing runs since the index was already attached at the library preparation stage.

We quantified the sequencing pool with the BioAnalyzer (Agilent) and/or the KAPA Library Quantification kit (KAPA Biosystems) and sequenced on Illumina MiSeq, HiSeq2500 or NextSeq500 instruments for 2 × 75, 2 × 100 or 2 × 150 cycles along with the indexing read(s).

**Enrichment for 394,577 SNP targets ('390k capture').** The protocol for enrichment for SNP targets was similar to the mitochondrial DNA capture, with the exception that we used another bait set (390k) and about twice as much library (up to 1,000 ng) compared to the mtDNA capture.

The specific capture reagent used in this study is described for the first time here. To target each SNP, we used a different oligonucleotide probe design compared to ref. 10. We used four 52 base pair probes for each SNP target. One probe ends just before the SNP, and one starts just after. Two probes are centred on the SNP, and are identical except for having the alternate alleles. This probe design avoids systematic bias towards one SNP allele or another. For the template sequence for designing the San and Yoruba panels baits, we used the sequence that was submitted for these same SNPs during the design of the Affymetrix Human Origins SNP array<sup>12</sup>. For SNPs that were both in the San and Yoruba panels, we used the Yoruba template sequence in preference. For all other SNPs, we used the human genome reference sequence as a template. Supplementary Data 2a–d gives the list of SNPs that we targeted, along with details of the probes used. The breakdown of SNPs into different classes is as follows.

124,106 'Yoruba SNPs': all SNPs in 'panel 5' of the Affymetrix Human Origins array (discovered as heterozygous in a Yoruba male: HGDP00927)<sup>12</sup> that passed the probe design criteria specified in ref. 11.

146,135 'San SNPs': all SNPs in 'panel 4' of the Affymetrix Human Origins array (discovered as heterozygous in a San male: HGDP01029)<sup>12</sup> that passed probe design criteria<sup>11</sup>. The full Affymetrix Human Origins array panel 4 contains several tens of thousands of additional SNPs overlapping those from panel 5, but we did not wish to redundantly capture panel 5 SNPs.

98,166 'compatibility SNPs': SNPs that overlap between the Affymetrix Human Origins, the Affymetrix 6.0, and the Illumina 610 Quad arrays, which are not already included in the 'Yoruba SNPs' or 'San SNPs' lists<sup>12</sup> and that also passed the probe design criteria<sup>11</sup>.

26,170 'miscellaneous SNPs': SNPs that did not overlap the Human Origins array. The subset analysed in this study were 2,258 Y chromosome SNPs ([http://isogg.org/tree/ISOGG\\_YDNA\\_SNP\\_Index.html](http://isogg.org/tree/ISOGG_YDNA_SNP_Index.html)) that we used for Y haplogroup determination.

**Processing of sequencing data.** We restricted analysis to read pairs that passed quality control according to the Illumina software ('PF reads').

We assigned read pairs to libraries by searching for matches to the expected index and barcode sequences (if present, as for the Adelaide and Boston libraries). We allowed no more than 1 mismatch per index or barcode, and zero mismatches if there was ambiguity in sequence assignment or if barcodes of 5 bp length were used (Adelaide libraries).

We used Seqprep (<https://github.com/jstjohn/SeqPrep>) to search for overlapping sequence between the forward and reverse read, and restricted to molecules where we could identify a minimum of 15 bp of overlap. We collapsed the two reads into a single sequence, using the consensus nucleotide if both reads agreed, and the read with higher base quality in the case of disagreement. For each merged nucleotide, we assigned the base quality to be the higher of the two reads. We further used Seqprep to search for the expected adaptor sequences at either ends of the merged sequence, and to produce a trimmed sequence for alignment.

We mapped all sequences using BWA-0.6.1 (ref. 35). For mitochondrial analysis we mapped to the mitochondrial genome RSRS<sup>36</sup>. For whole-genome analysis we mapped to the human reference genome hg19. We restricted all analyses to sequences that had a mapping quality of MAPQ ≥ 37.

We sorted all mapped sequences by position, and used a custom script to search for mapped sequences that had the same orientation and start and stop positions. We stripped all but one of these sequences (keeping the best quality one) as duplicates.

**Mitochondrial sequence analysis and assessment of ancient DNA authenticity.** For each library for which we had average coverage of the mitochondrial genome of at least tenfold after removal of duplicated molecules, we built a mitochondrial consensus sequence, assigning haplogroups for each library as described in Supplementary Information section 2.

We used contamMix-1.0.9 to search for evidence of contamination in the mitochondrial DNA<sup>13</sup>. This software estimates the fraction of mitochondrial DNA sequences that match the consensus more closely than a comparison set of 311 worldwide mitochondrial genomes. This is done by taking the consensus sequence of reads aligning to the RSRS mitochondrial genome, and requiring a minimum coverage of 5 after filtering bases where the quality was <30. Raw reads are then realigned to this consensus. In addition, the consensus is multiply aligned with the other 311 mitochondrial genomes using kalign (2.0.4)<sup>37</sup> to build the necessary inputs for contamMix, trimming the first and last 5 bases of every read to mitigate against the confounding factor of ancient damage. This software had difficulty running on data sets with higher coverage, and for these data sets, we down-sampled to 50,000 reads.

For all sequences mapping to the mitochondrial DNA for which the consensus mitochondrial DNA sequence had a cytosine at the terminal nucleotide, we measured the proportion of sequences with a thymine at that position. For population genetic analysis, we only used partially UDG-treated libraries with a minimum of 3% C→T substitutions as recommended by ref. 33. In cases where we used a fully UDG-treated library for 390k analysis, we examined mitochondrial capture data from a non-UDG-treated library made from the same extract, and verified that the non-UDG library had a minimum of 10% C→T at the first nucleotide as recommended by ref. 38. Metrics for the mitochondrial DNA analysis on each library are given in Supplementary Data 1.

**390k capture, sequence analysis and quality control.** For 390k analysis, we restricted to reads that not only mapped to the human reference genome hg19 but that also overlapped the 354,212 autosomal SNPs genotyped on the Human Origins array<sup>4</sup>. We trimmed the last two nucleotides from each sequence because we found that these are highly enriched in ancient DNA damage even for UDG-treated libraries. We further restricted analyses to sites with base quality ≥ 30.

We made no attempt to determine a diploid genotype at each SNP in each sample. Instead, we used a single allele—randomly drawn from the two alleles in the individual—to represent the individual at that site<sup>20,39</sup>. Specifically, we made an allele call at each target SNP using majority rule over all sequences overlapping the SNP. When each of the possible alleles was supported by an equal number of sequences, we picked an allele at random. We set the allele to 'no call' for SNPs at which there was no read coverage.

We restricted population genetic analysis to libraries with a minimum of 0.06-fold average coverage on the 390k SNP targets, and for which there was an unambiguous sex determination based on the ratio of X to Y chromosome reads (Supplementary Information section 4 and Supplementary Data 1). For individuals for whom there were multiple libraries per sample, we performed a series of quality control analysis. First, we used the ADMIXTURE software<sup>40,41</sup> in supervised mode, using Kharia, Onge, Karitiana, Han, French, Mbuti, Ulchi and Eskimo as reference populations. We visually inspected the inferred ancestry components in each individual, and removed individuals with evidence of heterogeneity in inferred ancestry components across libraries. For all possible pairs of libraries for each sample, we also computed statistics of the form  $D(Library_1, Library_2; Probe, Mbuti)$ , where *Probe* is any of a panel of the same set of eight reference

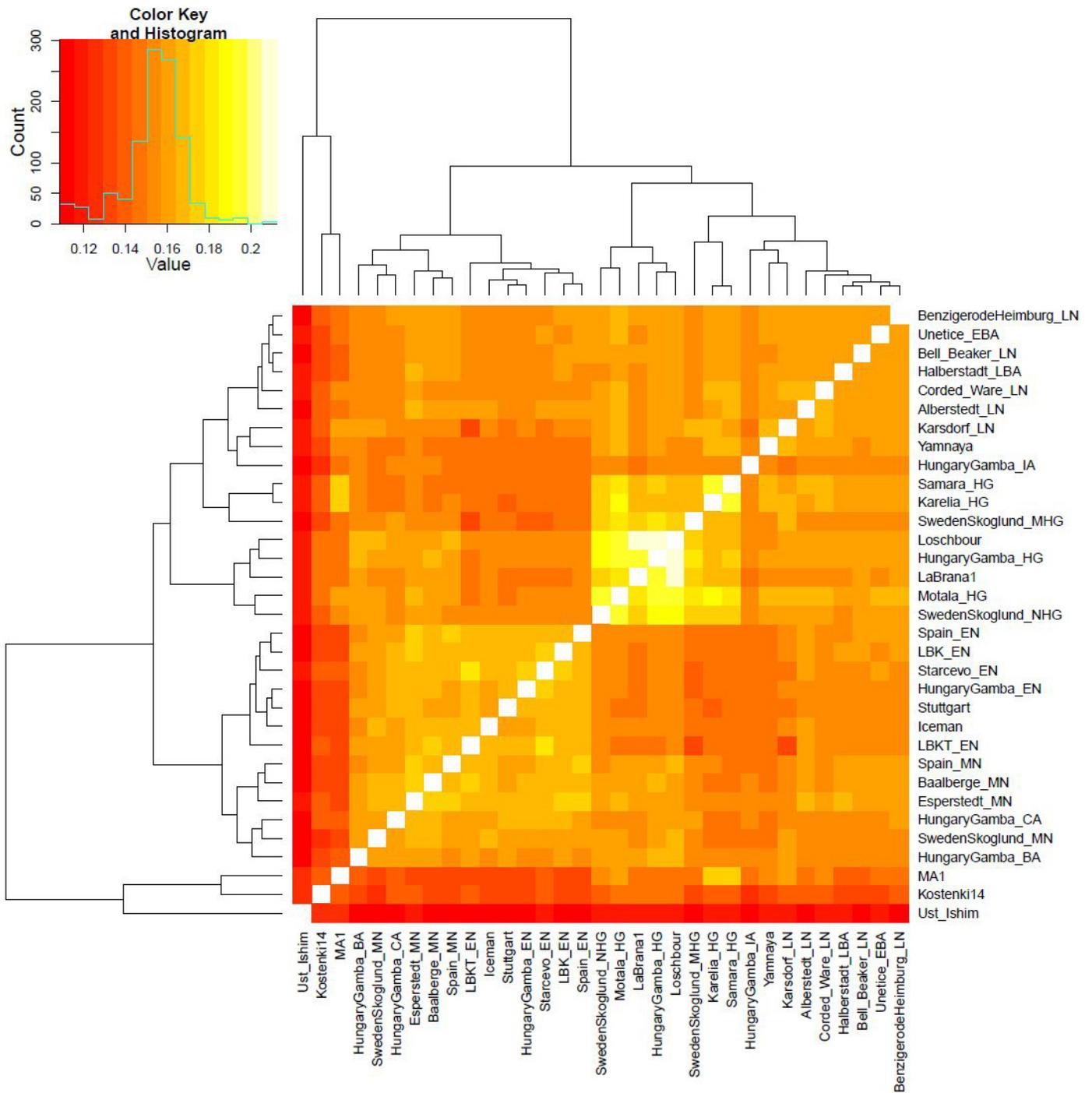
populations), to determine whether there was significant evidence of the *Probe* population being more closely related to one library from an ancient individual than another library from that same individual. None of the individuals that we used had strong evidence of ancestry heterogeneity across libraries. For samples passing quality control for which there were multiple libraries per sample, we merged the sequences into a single BAM.

We called alleles on each merged BAM using the same procedure as for the individual libraries. We used ADMIXTURE<sup>41</sup> as well as PCA as implemented in EIGENSOFT<sup>42</sup> (using the *lsqproject*: YES option to project the ancient samples) to visualize the genetic relationships of each set of samples with the same culture label with respect to 777 diverse present-day West Eurasians<sup>4</sup>. We visually identified outlier individuals, and renamed them for analysis either as outliers or by the name of the site at which they were sampled (Extended Data Table 1). We also identified two pairs of related individuals based on the proportion of sites covered in pairs of ancient samples from the same population that had identical allele calls using PLINK<sup>43</sup>. From each pair of related individuals, we kept the one with the most SNPs.

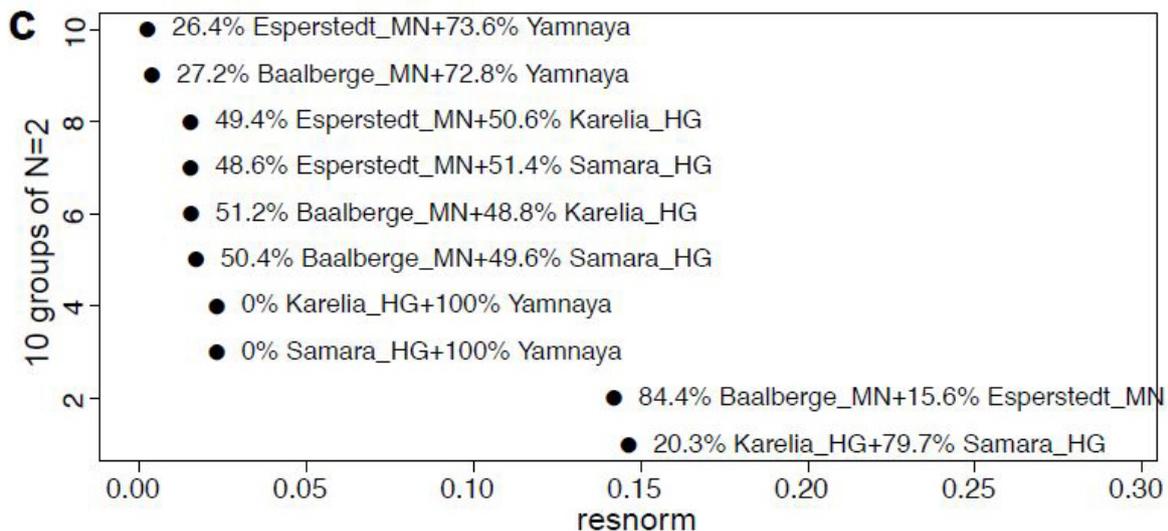
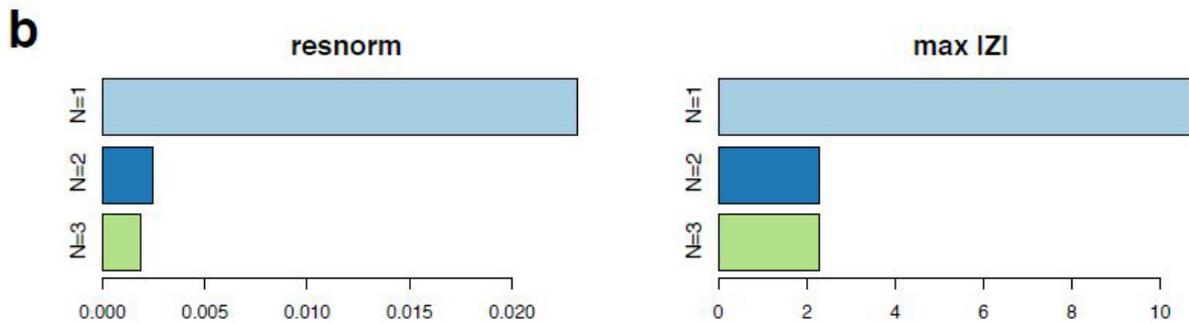
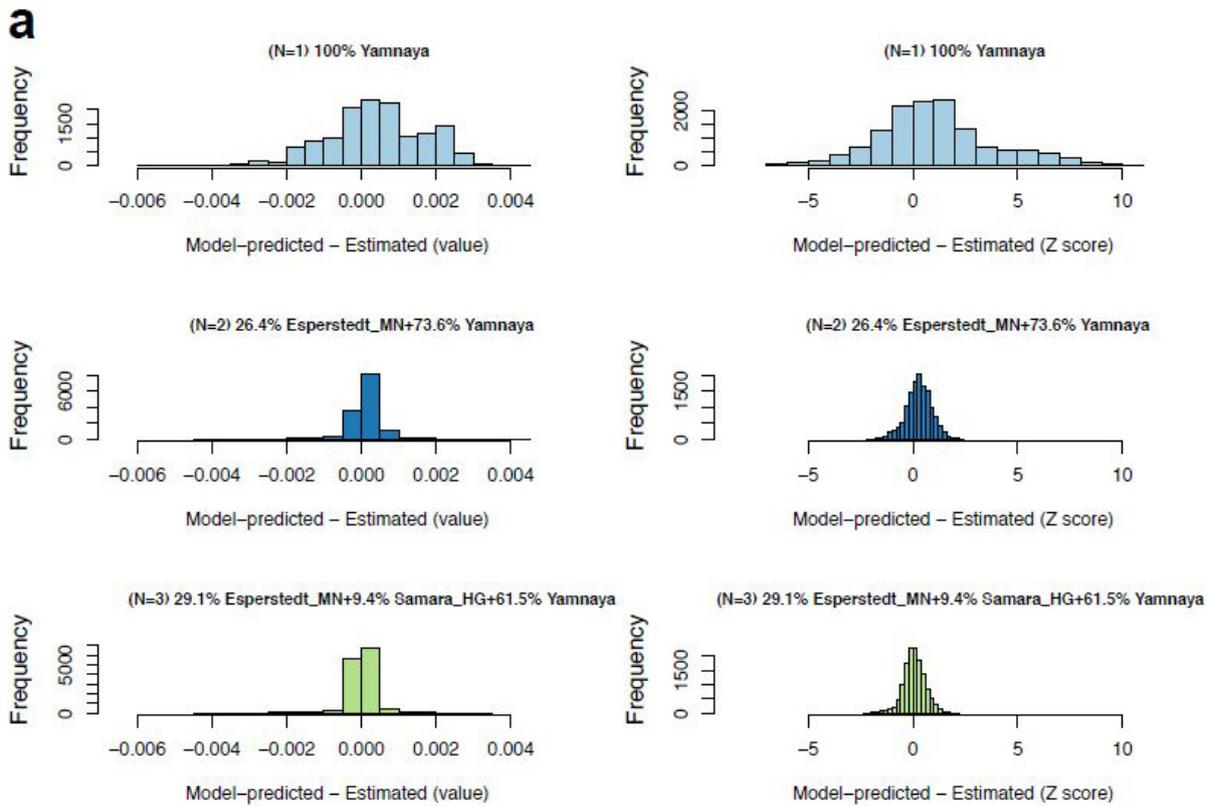
**Population genetic analyses.** We determined genetic sex using the ratio of X and Y chromosome alignments<sup>44</sup> (Supplementary Information section 4), and Y chromosome haplogroup for the male samples (Supplementary Information section 4). We studied population structure (Supplementary Information sections 5 and 6). We used *f* statistics to carry out formal tests of population relationships (Supplementary Information section 6) and built explicit models of population history consistent with the data (Supplementary Information section 7). We estimated mixture proportions in a way that was robust to uncertainty about the exact population history that applied (Supplementary Information section 8). We estimated the minimum number of streams of migration into Europe needed to explain the data (Supplementary Information sections 9 and 10). The estimated mixture proportions shown in Fig. 3 were obtained using the *lsqin* function of Matlab and the optimization method described in Supplementary Information section 9 with 15 world outgroups.

**Sample size.** No statistical methods were used to predetermine sample size.

31. Kircher, M., Sawyer, S. & Meyer, M. Double indexing overcomes inaccuracies in multiplex sequencing on the Illumina platform. *Nucleic Acids Res.* **40**, e3 (2012).
32. Meyer, M. *et al.* A mitochondrial genome sequence of a hominin from Sima de los Huesos. *Nature* **505**, 403–406 (2014).
33. Rohland, N., Harney, E., Mallick, S., Nordenfelt, S. & Reich, D. Partial uracil–DNA–glycosylase treatment for screening of ancient DNA. *Phil. Trans. R. Soc. Lond. B* **370**, 20130624 (2015).
34. Rohland, N. & Reich, D. Cost-effective, high-throughput DNA sequencing libraries for multiplexed target capture. *Genome Res.* **22**, 939–946 (2012).
35. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
36. Behar, D. M. *et al.* A “Copernican” reassessment of the human mitochondrial DNA tree from its root. *Am. J. Hum. Genet.* **90**, 675–684 (2012).
37. Lassmann, T. & Sonnhammer, E. L. L. Kalign—an accurate and fast multiple sequence alignment algorithm. *BMC Bioinformatics* **6**, 298 (2005).
38. Sawyer, S., Krause, J., Guschanski, K., Savolainen, V. & Pääbo, S. Temporal patterns of nucleotide misincorporations and DNA fragmentation in ancient DNA. *PLoS ONE* **7**, e34131 (2012).
39. Green, R. E. *et al.* A draft sequence of the Neandertal genome. *Science* **328**, 710–722 (2010).
40. Alexander, D. H. & Lange, K. Enhancements to the ADMIXTURE algorithm for individual ancestry estimation. *BMC Bioinformatics* **12**, 246 (2011).
41. Alexander, D. H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in unrelated individuals. *Genome Res.* **19**, 1655–1664 (2009).
42. Reich, D., Price, A. L. & Patterson, N. Principal component analysis of genetic data. *Nature Genet.* **40**, 491–492 (2008).
43. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
44. Skoglund, P., Storå, J., Götherström, A. & Jakobsson, M. Accurate sex identification of ancient human remains using DNA shotgun sequencing. *J. Archaeol. Sci.* **40**, 4477–4482 (2013).

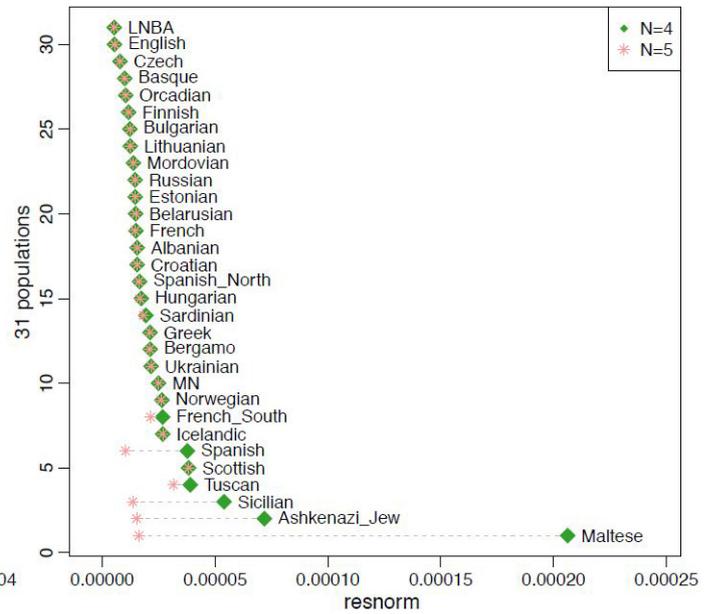
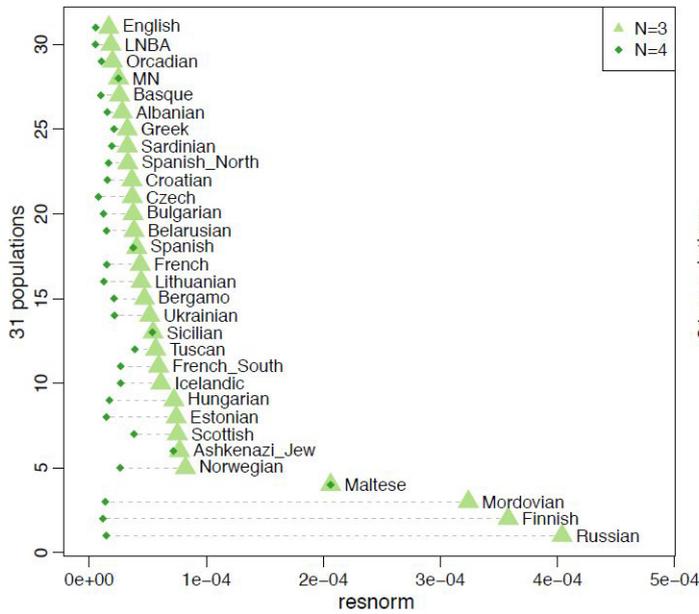
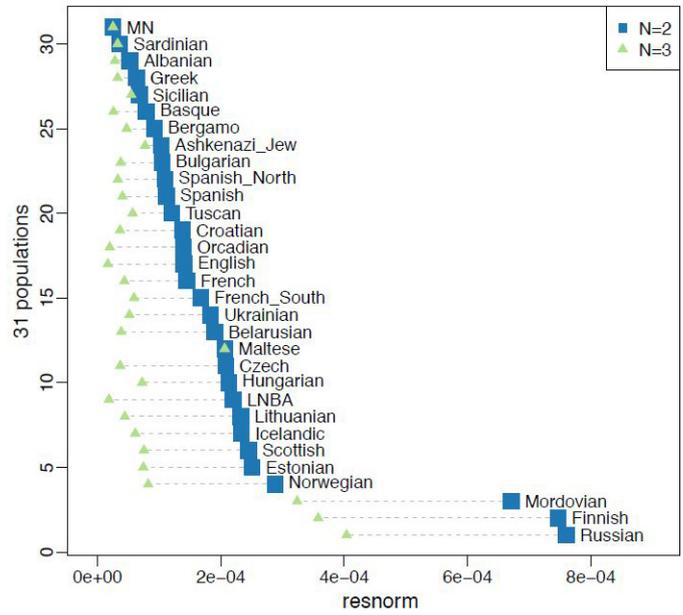
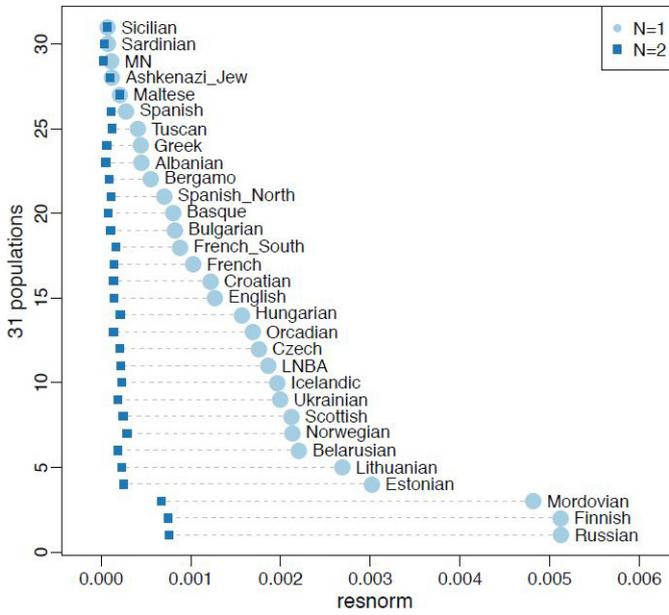


Extended Data Figure 1 | Outgroup  $f_3$  statistic  $f_3(\text{Dinka}; X, Y)$ , measuring the degree of shared drift among pairs of ancient individuals.



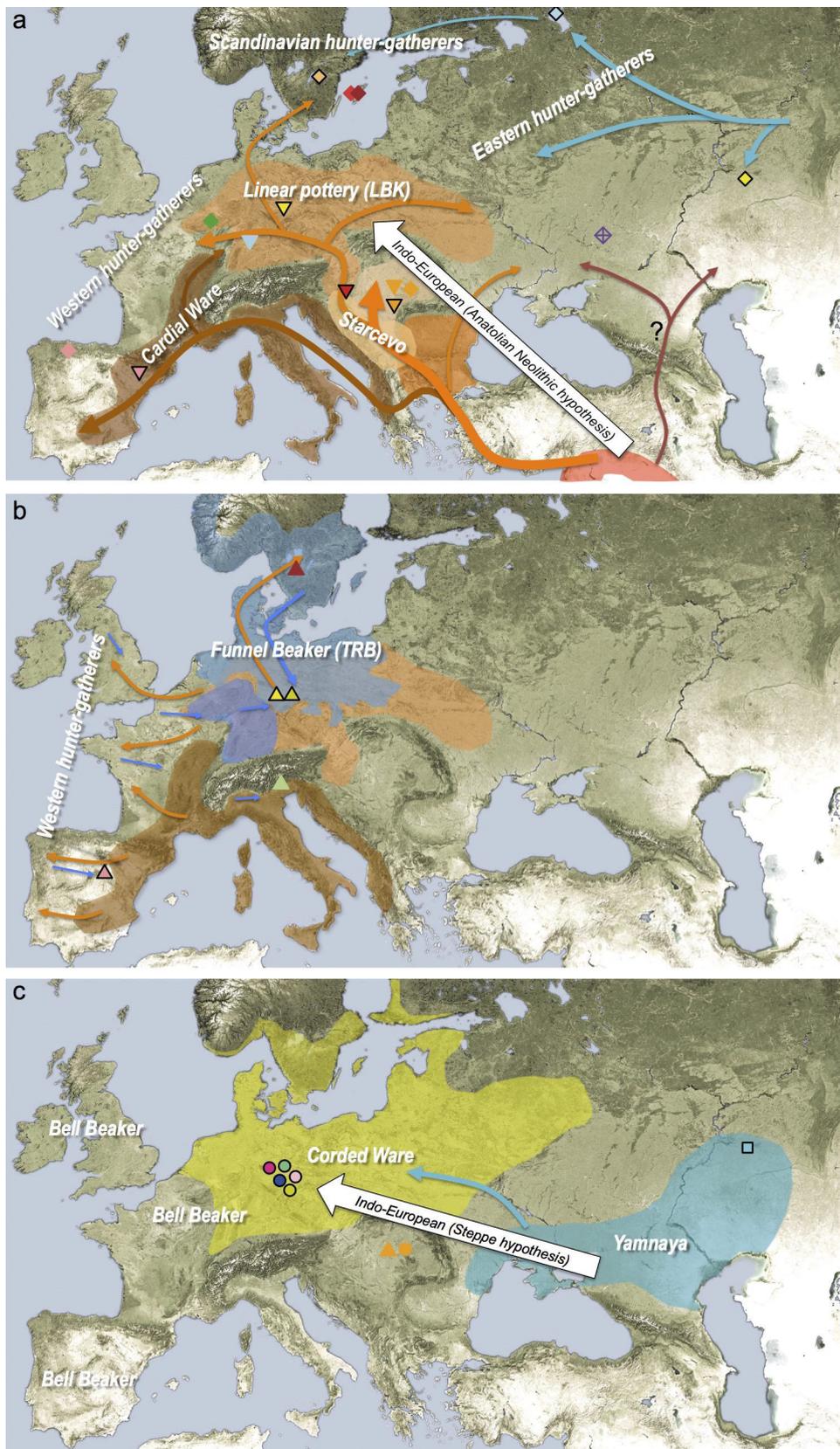
**Extended Data Figure 2 | Modelling Corded Ware as a mixture of  $N = 1, 2,$  or 3 ancestral populations.** **a**, The left column shows a histogram of raw  $f_4$  statistic residuals and on the right Z-scores for the best-fitting (lowest squared 2-norm of the residuals, or *resnorm*) model at each  $N$ . **b**, The data on

the left show *resnorm* and on the right show the maximum  $|Z|$  score change for different  $N$ . **c**, *resnorm* of different  $N = 2$  models. The set of outgroups used in this analysis in the terminology of Supplementary Information section 9 is 'World Foci 15 + Ancients'.



Extended Data Figure 3 | Modelling Europeans as mixtures of increasing complexity:  $N = 1$  (EN),  $N = 2$  (EN, WHG),  $N = 3$  (EN, WHG, Yamnaya),  $N = 4$  (EN, WHG, Yamnaya, Nganasan),  $N = 5$  (EN, WHG, Yamnaya,

Nganasan, BedouinB). The residual norm of the fitted model (Supplementary Information section 9) and its changes are indicated.



**Extended Data Figure 4 | Geographic distribution of archaeological cultures and graphic illustration of proposed population movements / turnovers discussed in the main text. a,** Proposed routes of migration by early farmers into Europe ~9,000–7,000 years ago. **b,** Resurgence of hunter-gatherer

ancestry during the Middle Neolithic 7,000–5,000 years ago. **c,** Arrival of steppe ancestry in central Europe during the Late Neolithic ~4,500 years ago. White arrows indicate the two possible scenarios of the arrival of Indo-European language groups. Symbols of samples are identical to those in Fig. 1.

**Extended Data Table 1 | Number of ancient Eurasian modern human samples screened in genome-wide studies to date**

First author	Description	No. samples at $\geq 0.05\times$ coverage (enough for Procrustes analysis)	No. samples at $>0.25\times$ coverage (enough to analyze in pairs)
Keller <sup>3</sup>	Tyrolean Iceman	1	1
Raghavan <sup>6</sup>	Upper Paleolithic Siberians	2	1
Olalde <sup>5</sup>	Mesolithic Iberian from LaBranca	1	1
Skoglund <sup>8</sup>	Farmers and hunter-gatherers from Sweden	5	2
Lazaridis <sup>4</sup>	Early European farmer from Germany & Mesolithic hunter-gatherers from Luxembourg and Sweden	7	4
Gamba <sup>2</sup>	Neolithic, Bronze Age, Iron Age Hungary	13	9
Fu <sup>1</sup>	Upper Paleolithic Siberian from Ust-Ishim	1	1
Seguin-Orlando <sup>7</sup>	Upper Paleolithic European from Kostenki	1	1
<i>Total before study</i>		<i>31</i>	<i>20</i>
This study	Hunter-gatherers and pastoralists from Russia, Mesolithic hunter-gatherers from Sweden, Early Neolithic from Germany, Hungary, and Spain, Middle Neolithic from Germany & Spain, Late Neolithic / Bronze Age from Germany	69	58

Only studies that produced at least one sample at  $\geq 0.05\times$  coverage are listed.

Extended Data Table 2 | Summary of the archaeological context for the 69 newly reported samples

Reich ID	Pop Label for Analysis	Culture	Group	Location and sample details (e.g. sample, grave and museum ID)	Date (lab no.)	Country	Sex	mt-hg	Y-hg	Autosomal SNPs
0061	Karelia_HG	Russian Mesolithic	EHG	Yuzhnyy Oleni Ostrov, Karelia, Russia, UzOO74, grave 142, MAE RAS 5773-74	5500-5000 BC	Russia	M	C1g (formerly C1f)	R1a1	341554
0124	Samara_HG	Russian Neolithic HG	EHG	Sok River, Samara, Russia, SVP44	5650-5555 cal BC (Beta - 392490)	Russia	M	U5a1d	R1b1	206748
0051	Motala_HG	Swedish Mesolithic	SHG	Motala, Sweden, Motala 1	5898-5531 cal BC	Sweden	F	U5a1	I2c2	229271
0012	Motala_HG	Swedish Mesolithic	SHG	Motala, Sweden, Motala 2	5898-5531 cal BC	Sweden	M	U2e1	I2c2	292853
0013	Motala_HG	Swedish Mesolithic	SHG	Motala, Sweden, Motala 3	5898-5531 cal BC	Sweden	M	U5a1	I2a1b	251108
0014	Motala_HG	Swedish Mesolithic	SHG	Motala, Sweden, Motala 4	5898-5531 cal BC	Sweden	F	U5a2d	I2a1b	311299
0015	Motala_HG	Swedish Mesolithic	SHG	Motala, Sweden, Motala 6	5898-5531 cal BC	Sweden	M	U5a2d	I2a1	285307
0016	Motala_HG	Swedish Mesolithic	SHG	Motala, Sweden, Motala 9	5898-5531 cal BC	Sweden	M	U5a2	I2a1	275233
0017	Motala_HG	Swedish Mesolithic	SHG	Motala, Sweden, Motala 12	5898-5531 cal BC	Sweden	M	U2e1	I2a1b	337794
0174	Starcevo_EN	Starcevo	EN	Aisóczy-Bálaszék, Mémóli telep, Hungary, BAM25a, feature 1532	5710-5550 cal BC (MAMS 11939)	Hungary	M	N1a1a1b	H2	101653
0176	LBKT_EN	LBKT	EN	Szemely-Hegyes, Hungary, SZEH4b, feature 1001	5210-4940 cal BC (Beta - 310038)	Hungary	M	N1a1a1a3		30718
0046	LBK_EN	LBK	EN	Halberstadt-Sonntagsfeld, Germany, HAL5, grave 2, feature 241.1	5206-5004 cal BC (MAMS 21479)	Germany	F	T2c1d'e'f		296764
0048	LBK_EN	LBK	EN	Halberstadt-Sonntagsfeld, Germany, HAL25, grave 28, feature 861	5206-5052 cal BC (MAMS 21482)	Germany	M	K1a	G2a2a	123628
0056	LBK_EN	LBK	EN	Halberstadt-Sonntagsfeld, Germany, HAL14, grave 15, feature 430	5206-5052 cal BC (MAMS 21480)	Germany	M	T2b(8)	G2a2a	136578
0057	LBK_EN	LBK	EN	Viesenhäuser Hof, Stuttgart-Mühlhausen, Germany, LBK1976	5207-5067 cal BC (MAMS 21483)	Germany	F	N1a1a1		55902
0100	LBK_EN	LBK	EN	Viesenhäuser Hof, Stuttgart-Mühlhausen, Germany, LBK1992	5032-4946 cal BC (KIA40341)	Germany	F	N1a1a1a		342342
0059	LBK_EN	LBK	EN	Halberstadt-Sonntagsfeld, Germany, HAL2, grave 35, feature 999	5073-4997 cal BC (KIA40350)	Germany	M	N1a1a	G2a2a1	191097
				5066-4979 cal BC (KIA30408)						
0821	LBK_EN	LBK	EN	Halberstadt-Sonntagsfeld, Germany, HAL24, grave 27, feature 867	5034-4942 cal BC (KIA40348)	Germany	M	Pre-K2d1	G2a2a1	55914
0795	LBK_EN	LBK	EN	Karsdorf, Germany, KAR6a, feature 170	5207-5070 cal BC (MAMS 22923)	Germany	M	H1	T1a	47804
0054	LBK_EN	LBK	EN	Oberwiederstedt-Unterviederstedt, UWSVA, Germany, grave 6, feature 1 14	5209-5070 cal BC (MAMS 21485)	Germany	F	J1c17		337625
0022	LBK_EN	LBK	EN	Viesenhäuser Hof, Stuttgart-Mühlhausen, Germany, LBK1976	5600-4800 BC	Germany	F	T2e		160852
0025	LBK_EN	LBK	EN	Viesenhäuser Hof, Stuttgart-Mühlhausen, Germany, LBK1992	5600-4800 BC	Germany	F	T2b		307686
0026	LBK_EN	LBK	EN	Viesenhäuser Hof, Stuttgart-Mühlhausen, Germany, LBK2155	5600-4800 BC	Germany	F	T2b		315484
0409	Spain_EN	Els_Trocs	EN	Els Trocs, Spain, Troc1	5311-5218 cal BC (MAMS 16159)	Spain	F	J1c3		172903
0410	Spain_EN	Els_Trocs	EN	Els Trocs, Spain, Troc3	5178-5086 cal BC (MAMS 16161)	Spain	M	pre-T2c1d2	R1b1	237595
0411	Spain_EN	Els_Trocs	EN	Els Trocs, Spain, Troc4	5177-5068 cal BC (MAMS 16162)	Spain	M	K1a2a	F*	31507
0412	Spain_EN	Els_Trocs	EN	Els Trocs, Spain, Troc5	5310-5206 cal BC (MAMS 16164)	Spain	M	N1a1a1	I2a1b1	333940
0413	Spain_EN	Els_Trocs	EN	Els Trocs, Spain, Troc7	5303-5204 cal BC (MAMS 16166)	Spain	F	V		295844
0405	Spain_MN	La_Mina	MN	La Mina, Spain, Mina9	3900-3600 BC	Spain	M	K1a1b1	I2a1a1/H27	123230
0406	Spain_MN	La_Mina	MN	La Mina, Spain, Mina9d	3900-3600 BC	Spain	M	H1	I2a2a1	324169
0407	Spain_MN	La_Mina	MN	La Mina, Spain, Mina9b	3900-3600 BC	Spain	F	K1b1a1		236225
0408	Spain_MN	La_Mina	MN	La Mina, Spain, Mina18a	3900-3600 BC	Spain	F	pre-U5b1f		321761
0172	Esperstedt_MN	Salzmünde/Bernburg	MN	Esperstedt, Germany, ESP24, feature 1841	3960-3686 cal BC (E18699)	Germany	M	T2b	I2a1b1a	279147
0559	Baalberge_MN	Baalberge	MN	Quedlinburg, Germany, QLB15d, feature 21033	3645-3537 cal BC (MAMS 22818)	Germany	M	HV6'17	R7*	64304
0560	Baalberge_MN	Baalberge	MN	Quedlinburg, Germany, QLB18a, feature 21039	3640-3510 cal BC (E17956)	Germany	F	T2e1		133305
0807	Baalberge_MN	Baalberge	MN	Esperstedt, Germany, ESP30, feature 6220	3687-3797 cal BC (E17784)	Germany	M	H1e1a	F*	33481
0231	Yamnaya	Yamnaya	EBA	Ekaterinovka, Southern Steppe, Samara, Russia, SVP3	2910-2875 cal BC (Beta 392487)	Russia	M	U4a1	R1b1a2a2	348142
0357	Yamnaya	Yamnaya	EBA	Lopatino I, Sok River, Samara, Russia, SVP5 same sample as SVP3	3090-2910 cal BC (Beta 392489)	Russia	F	W6		163845
0370	Yamnaya	Yamnaya	EBA	Ishtinovka I, Eastern Orenburg, Pre-Ural steppe, Samara, Russia, SVP10	3300-2700 BC	Russia	M	H13a1a1a	R1b1a2a2	199345
0429	Yamnaya	Yamnaya	EBA	Lopatino I, Sok River, Samara, Russia, SVP39	3339-2917 cal BC (AA47804)	Russia	M	T2c1a2	R1b1a2a2	217594
0438	Yamnaya	Yamnaya	EBA	Luzhki I, Samara River, Samara, Russia, SVP50	3021-2635 cal BC (AA47807)	Russia	M	U5a1a1	R1b1a2a2	213493
0439	Yamnaya	Yamnaya	EBA	Lopatino I, Sok River, Samara, Russia, SVP52	3305-2925 cal BC (Beta 392491)	Russia	M	U5a1a1	R1b1a	98900
0441	Yamnaya	Yamnaya	EBA	Kurmanavskii III, Buzuluk, Samara, Russia, SVP54	3010-2622 cal BC (AA47805)	Russia	F	H2b		51326
0443	Yamnaya	Yamnaya	EBA	Lopatino II, Sok River, Samara, Russia, SVP57	3300-2700 BC	Russia	M	W3a1a	R1b1a2a	343890
0444	Yamnaya	Yamnaya	EBA	Kutuluk I, Kutuluk River, Samara, Russia, SVP58	3300-2700 BC	Russia	M	H6a1b	R1b1a2a2	167126
0550	Karsdorf_LN	unknown	LN	Karsdorf, Germany, KAR22a, feature 191	2564-2475 cal BC (MAMS 22344)	Germany	F	T1a1		59907
10103	Corded_Ware_LN	Corded Ware	LN	Esperstedt, Germany, ESP16, feature 6236	2566-2477 cal BC (MAMS 21488)	Germany	F	W6a		336918
0049	Corded_Ware_LN	Corded Ware	LN	Esperstedt, Germany, ESP22, feature 6140	2454-2291 cal BC (MAMS 21489)	Germany	F	X2b4		167170
10106	Corded_Ware_LN	Corded Ware	LN	Esperstedt, Germany, ESP26, feature 6233.1	2454-2291 cal BC (MAMS 21490)	Germany	F	T2a1b1		69886
10104	Corded_Ware_LN	Corded Ware	LN	Esperstedt, Germany, ESP11, feature 6216	2473-2348 cal BC (MAMS 21487)	Germany	M	U4b1a1a1	R1a1a1	336637
0059	Benzingerode-Heimburg_LN	Bell Beaker?	LN	Benzingerode-Heimburg, Germany, BZH6, grave 2, feature/find 1287/1036	2286-2153 cal BC (MAMS 21486)	Germany	F	H1/H1b1ad		241081
0058	Benzingerode-Heimburg_LN	Bell Beaker	LN	Benzingerode-Heimburg, Germany, BZH4, grave 7, feature 4607	2283-2146 cal BC (MAMS 21491)	Germany	F	H1e		246728
01171	Benzingerode-Heimburg_LN	Bell Beaker?	LN	Benzingerode-Heimburg, Germany, BZH12, grave 3, feature 6256	2294-2136 cal BC (KIA27952)	Germany	F	U5a1a2a		96900
01112	Bell_Beaker_LN	Bell Beaker	LN	Quedlinburg XII, Germany, QUEX16, feature 6256	2340-2190 cal BC (Er7038)	Germany	F	H13a1a2		341003
01113	Bell_Beaker_LN	Bell Beaker	LN	Quedlinburg XII, Germany, QUEX14, feature 6255.1	2290-2130 cal BC (Er7283)	Germany	F	J1c14		190352
01108	Bell_Beaker_LN	Bell Beaker	LN	Rothenschimbach, Germany, ROT16, feature 10044	2497-2436 cal BC (Er8710)	Germany	F	H5a3		200528
01111	Bell_Beaker_LN	Bell Beaker	LN	Rothenschimbach, Germany, ROT14, feature 10142	2414-2333 cal BC (Er8712)	Germany	F	H5new		292556
0060	Bell_Beaker_LN	Bell Beaker	LN	Rothenschimbach, Germany, ROT3, feature 10011	2294-2206 cal BC (MAMS 22819)	Germany	F	K1a2c		47085
0806	Bell_Beaker_LN	Bell Beaker	LN	Quedlinburg VII 2, Germany, QLB28b, feature 19617	2296-2206 cal BC (MAMS 22820)	Germany	M	H1	R1b1a2a1a2	91757
01118	Alberstedt_LN	unknown	LN	Alberstedt, Germany, ALB3, feature 7144.2	2459-2345 cal BC (MAMS 21482)	Germany	F	HV6'17		349656
01114	Unetice_EBA_relative_of_01017	Unetice	EBA	Esperstedt, Germany, ESP2, feature 3340.1	2131-1979 cal BC (MAMS 21493)	Germany	M	I3a	I2a2	217031
01115	Unetice_EBA	Unetice	EBA	Esperstedt, Germany, ESP3, feature 1559.1	1931-1790 cal BC (MAMS 21494)	Germany	F	U5a1		123744
01116	Unetice_EBA	Unetice	EBA	Esperstedt, Germany, ESP4, feature 3322/3323	2118-1961 cal BC (MAMS 21495)	Germany	M	W3a1	I2c2	308158
01117	Unetice_EBA	Unetice	EBA	Esperstedt, Germany, ESP29, feature 3332/3333	2166-2064 cal BC (MAMS 21496)	Germany	F	I3a		279696
01164	Unetice_EBA	Unetice	EBA	Quedlinburg VIII, Germany, QUEVIII6, feature 3580	2012-1919 cal BC (MAMS 21497)	Germany	F	pre-U5b2a1b		323832
0803	Unetice_EBA	Unetice	EBA	Eulau, Germany, EUL41A, feature 882	2115-1966 cal BC (MAMS 22822)	Germany	F	H4a1a1		144186
0804	Unetice_EBA	Unetice	EBA	Eulau, Germany, EUL57b, feature1911.1	2131-1862 cal BC (MAMS 22821)	Germany	M	H3	I2	22689
0047	Unetice_EBA	Unetice	EBA	Halberstadt-Sonntagsfeld, Germany, HAL16, grave 19, feature 613.1	2022-1937 cal BC (MAMS 21481)	Germany	F	V		288353
0099	Halberstadt_LBA	Late Bronze Age	LBA	Halberstadt-Sonntagsfeld, Germany, HAL36c, grave 40, feature 1114	1113-1021 cal BC (MAMS 21484)	Germany	M	I23	R1a1a1b1a2	337566

Samples with direct radiocarbon dates are indicated by a calibrated date "cal bc" along with associated laboratory numbers. Dates that are estimated based on faunal elements associated with the samples are not indicated with 'cal' (although they are still calibrated, absolute dates).

