

# Early farmers from across Europe directly descended from Neolithic Aegeans

Zuzana Hofmanová, Susanne Kreutzer, Garrett Hellenthal, Christian Sell, Yoan Diekmann, David Díez-del-Molino, Lucy van Dorp, Saioa López, Athanasios Kousathanas, Vivian Link, Karola Kirsanow, Lara M. Cassidy, Rui Martiniano, Melanie Strobel, Amelie Scheu, Kostas Kotsakis, Paul Halstead, Sevi Triantaphyllou, Nina Kyparissi-Apostolika, Dushka Urem-Kotsou, Christina Ziota, Fotini Adaktylou, Shyamalika Gopalan, Dean M. Bobo, Laura Winkelbach, Jens Blöcher, Martina Unterländer, Christoph Leuenberger, Çiler Çilingiroğlu, Barbara Horejs, Fokke Gerritsen, Stephen J. Shennan, Daniel G. Bradley, Mathias Currat, Krishna R. Veeramah, Daniel Wegmann, Mark G. Thomas, Christina Papageorgopoulou, and Joachim Burger

## Supplementary information

SI1	Archaeological Background . . . . .	2
SI2	Sample Preparation . . . . .	11
SI3	Read Processing . . . . .	20
SI4	Analysis of Uniparental Markers and X Chromosome Contamination Estimates . . . . .	24
SI5	Genotype Calling for Ancient DNA . . . . .	31
SI6	PCA . . . . .	42
SI7	Using $f$ -statistics to Infer Genetic Relatedness and Admixture Amongst Ancient and Contemporary Populations . . . . .	49
SI8	Proportions of ancestral clusters in Neolithic populations of Europe . . . . .	68
SI9	Population continuity . . . . .	75
SI10	Comparing allele frequency patterns among samples using a mixture model . . . . .	81
SI11	Runs of Homozygosity . . . . .	100
SI12	Functional Markers . . . . .	102
	Data available online . . . . .	112

## **SI1 Archaeological Background**

### **Neolithic Archaeology in Northwest Anatolia, the Balkans and central Europe**

Çiler Çilingiroğlu, Fokke Gerritsen, Barbara Horejs, Joachim Burger & Stephen J. Shennan

#### **Introduction**

The regions of the Greek peninsula, - the Aegean including coastal western Anatolia, the southern Balkans and the Sea of Marmara in northwest Turkey - contain archaeological evidence of a complex transformation to a Neolithic way of life. Intensive archaeological investigations over the last several decades have provided insight into the potential trajectories of Neolithic spread over the European continent (e.g. [1, 2, 3, 4, 5, 6]). These data show varying local patterns of transformation, probably due to differing routes of Neolithization.

The coastal zones of the Aegean including Greece and western Anatolia appear to have been highly affected by seafaring connectivity between migratory groups, visible in the first pioneer coastal sites dating around 6700 cal BCE [7, 8, 9, 10, 11]. Northwest Anatolia and the Sea of Marmara, on the other hand, show clear affinities to the Central Anatolian Neolithic, which can be regarded as the origin of the earliest northwest Anatolian farmer-herders [5, 12, 13]. The spread of a Neolithic way of life to northwest Anatolia is most often described as a migration from inner Anatolia via terrestrial routes (e.g. [4, 14, 15, 16]).

Although current archaeological research has revealed various pathways of Neolithization in the first half of the 7th millennium BCE, questions still remain regarding how and where these trajectories overlapped and influenced each other in generating the complex emergence of agriculturalist lifestyles on the southeastern edge of Europe.

#### **The Northwest Anatolian Neolithic**

From 6600 cal BCE onwards, the first sedentary farming villages appear in northwest Anatolia, as revealed by excavations at Barcın, Menteşe and Aktopraklık [17, 18, 19]. The economy of the northwest Anatolian groups was based on farming and herding. In coastal sites, there is clear evidence of the exploitation of aquatic resources including pelagic fish and molluscs [20]. Production of dairy products is confirmed by lipid residue analyses on pottery from almost all excavated sites in the region [21, 22]. The rarity of clay stamps and figurines, the heavy use of timber as a building material and the use of round plans in architecture at some of the sites all indicate the development of a distinctly local culture, coined the "Fikirtepe Culture" by M. Özdoğan, which is easily distinguishable from Central, southwest and West Anatolia. This local trajectory should not, however, be taken to imply in any way that these groups were isolated. Both zooarchaeological studies and XRF studies on obsidian demonstrate that there was considerable mobility of animals and raw materials [20, 23].

## Neolithization of Southeast Europe (with Central Europe)

The archaeological evidence suggests that there are two, possibly three, major routes of Neolithic dispersal from northwest Anatolia into southeast Europe. The first route extending across Thrace is evidenced by the monochrome phase at Aşağı Pınar and the foundation of the first Karanovo I sites with painted pottery around 6000 BCE in western Thrace [24]. A second route was over the northern Aegean, penetrating the Struma Basin [25]. The third route may have followed the Black Sea littoral from northern Turkey to the lower Danube, from where it was possible for the pioneer farmers to enter Central Europe (see below). New archaeological findings from northern Greece suggest that the initial phase of Neolithization may have appeared almost simultaneously on both sides of the northern Aegean [26, 27, 28, 29, 30].

The early Neolithic sites in the Balkans were founded by farmer-herders who brought with them domesticated forms of cereals, pulses, and herd animals. Early Neolithic sites in Bulgaria have produced specimens of most of the founder crops that were first domesticated in southwest Asia in 9-8<sup>th</sup> millennia BCE [31, 32]. Similarly, domesticated sheep, goats, and cattle are present at the Early Neolithic in Bulgaria [33]. Recent investigations have confirmed that the Neolithic groups settling in northwest Anatolia and southeast Europe produced dairy products [21, 34]. The subsistence economies of Early Neolithic groups in southeastern Europe and northwest Anatolia shared many characteristics in common, both in terms of species composition and agropastoral practices. The appearance of all the aforementioned plant and animal species in a domesticated form supports the occurrence of extra-local Neolithic dispersals into the area ([35] but also see [30, 36, 37]). Therefore, an independent indigenous transition to Neolithic lifeways cannot be inferred from the available archaeological evidence.

In the following phase, dated roughly to 6000 BCE, the remains of Neolithic groups are better preserved and identified. Across the peninsula, farmer-herder groups with regional cultural characteristics appear and are named mainly according to the pottery they produced. Thus, in Bulgaria, white-on-red painted pottery and flint macro-blades are associated with early farmers of the “Karanovo I Culture”, while the northern Balkans is associated with the material culture of the “Starčevo-Criş-Körös Culture Group”.

It is now established that the Starčevo groups in the northern Balkans migrated to the Danube Valley and established the first Linearbandkeramik settlements in west Hungary around 5600/5500-5350 cal BCE [38]. The archaeobotanical, zoological and palaeogenetic evidence (e.g. [39, 40, 41]) confirms a non-indigenous appearance of Neolithic sites in Central Europe, with little admixture. The contribution of local foragers, however, was not insignificant as the lithic industries of early farmers show clear similarities to Mesolithic foragers toolkit ([38]; but see also [42]).

## Archaeological background of the early Neolithic individuals from Barcın (northwestern Anatolia)

Fokke Gerritsen

Barcın Höyük was established around 6600 cal BCE in the Yenişehir Valley east of the modern city of Bursa [43]. Since 2005, excavations by an international team headed by the Netherlands Institute in Turkey have investigated an uninterrupted occupational sequence spanning about six centuries, with occupation ending around 6000 cal BCE [17]. Other known Neolithic sites in the region, including sites associated with the Archaic and Classic Fikirtepe Horizon, appear to have started several centuries later. Other inland sites are known from about 6400 cal BCE onwards at Aktopraklık [19] and Menteşe [18], whereas sites in coastal locations, including Fikirtepe, Pendik and Yenikapı, appear around 6200 cal BCE [12]. The early date of Barcın Höyük underlines its importance for the study of when, how, and why farming became the main subsistence strategy in northwestern Anatolia. The cultural history of the region can be traced through the continuous and gradually evolving ceramic traditions throughout the Neolithic sequence at Barcın Höyük. The archeological data strongly suggest cultural continuity from a pre-Fikirtepe phase into the Fikirtepe Horizon [44].

The settlement was founded on a low natural elevation at the edge of a lake or marshland [45]. From the start of occupation at Barcın Höyük, the food economy was fully agrarian, based on cultivated cereals and pulses, as well as animal husbandry strategies relying largely on domestic cattle and sheep (Galik in [46]). Wild resources, including fish, molluscs, birds and game, as well as nuts, formed only a minor addition to the diet.

The absence of a transitional stage from foraging to farming is also clear from architectural traditions and habitation practices. The earliest architecture at the site, belonging to stratigraphic phase VIe, and discovered during the 2015 season of excavations, consists of timber buildings. Postholes indicate rectangular buildings aligned in a row, with open spaces in front of and behind the buildings. From phase VIId1 onwards, post-built walls were constructed over 30-50 cm deep foundation trenches. Mud mixed with straw was used to close the wall faces. The walls and centrally located upright posts would have carried the gabled roof. The floor surface of buildings varied between 12 and 25 m<sup>2</sup>.

The original layout of a row of buildings facing a central courtyard provides a settlement structure that is maintained during subsequent building phases throughout phases VIId1, VIId2 and VIc. This only changes in phase VIb, around 6200 cal BCE, when dispersed post-wall houses are built in the former courtyards [17]. Architectural remains of the final Neolithic phase of occupation, VIa, were not preserved well enough to establish the nature of the settlement structure.

The mounded site measures less than a hectare, and surface finds suggest that parts of the mound were not inhabited until well after the Neolithic period. The size of these communities is difficult to

ascertain based on current excavations, but is estimated to be on the order of dozens rather than hundreds of people.

The dead were buried in several locations within the settlement. Neonates and infants were buried within the houses, generally next to the walls, whereas juveniles and adults were buried in the central courtyard. All burials are primary inhumation graves in simple pits [47]. The corpse was typically placed on its side in a tightly flexed position, sometimes on its back with knees drawn up to the chest. In total, more than 100 burials have been excavated, of which about two thirds are comprised of neonates and infants.

Sample Bar8 (M10-106) was taken from a middle aged female with poor dental health [47]. The individual was buried in tightly flexed position in a simple pit among a small cluster of burials of adults and juveniles in the northern courtyard. It is dated stratigraphically to Barcm phase VIb, and has a  $^{14}\text{C}$  date of  $7238 \pm 38$  BP and a 95% calibrated range of 6212-6030 cal BCE (UBA-29837). The radiocarbon dating result was calibrated with OxCal v4.2.2 using INTCAL13.

Sample Bar31 (L11W-546) was taken from a genetically male adult buried in tightly flexed position in a simple pit. The grave was part of a cluster of adult burials in the central courtyard. It is dated stratigraphically to Barcm phase VIc or VIId, has a  $^{14}\text{C}$  date of  $7457 \pm 44$  BP, and has a 95% calibrated range of 6419-6238 cal BCE (UBA-29838). The radiocarbon dating result was calibrated with OxCal v4.2.2 using INTCAL13.

## The Mesolithic and Neolithic Period in Greece

Christina Papageorgopoulou, Kostas Kotsakis, Sevi Triantaphyllou, Dushka Urem-Kotsou, Nina Kyparissi-Apostolika, Fotini Adaktylou & Christina Ziota

The Greek Mesolithic is represented by archaeological and anthropological findings spanning the period 8600-6800/6700 BCE, most of which derive from cave sites (Franchthi, Theopetra, Cyclops) and open air sites (Sidari in Corfu, Maroulas in the island of Kythnos). There are no published data for the period 9300-8600 BCE, but archaeological artefacts dated before 9300 BCE have been reported from Palaeolithic sites (Klithi and Boila Rockshelter in Epirus) and from the Palaeolithic strata of Franchthi cave ([48]; for review see [49]). The Neolithic Period in Greece dates to 6700-3200 BCE and is divided into Early Neolithic (6700-5800/5600 BCE), Middle Neolithic (5800/5600-5400/5300 BCE), Late Neolithic (5400/5300-4500 BC), and Final Neolithic (4500-3300/3100 BC) [50, 51].

The lack of Mesolithic archaeological evidence and skeletal remains has created the perception that no significant cultural or economic development took place during the Mesolithic period [52]. Although this concept of the Greek Mesolithic is still accepted, especially among those stressing the allochthonous origin of the Neolithic way of life, new excavations and systematic analysis suggest the presence of complex Mesolithic groups both in the Aegean islands and mainland Greece. These groups followed a variety of subsistence strategies depending on their environments, and these differences make it challenging to identify specific unifying cultural elements for the Greek Mesolithic. Cultural and economic changes likely began earlier in the east, especially the Aegean Sea region, than in the west, namely in Epirus and the Ionian Sea [49].

Moreover, recent investigations, especially in northern Greece, have revealed a large number of early Neolithic sites in a region that was previously thought to be sparsely inhabited during that period. These findings have supported previous hypotheses that early sites were short-lived, only producing thin deposits that were subsequently covered by thick alluvial deposits (lowland sites) or underwent serious erosion during the dramatic geological changes of the Holocene (hillside sites) [53, 54]. The earliest Early Neolithic sites in Greece are found at Franchthi and Knossos (6700 cal BCE) in the south [55], in Thessaly (Argissa, Achilleion, Sesklo), dated at 6500 cal BCE, and recently in western and central Macedonia (Mauropigi, Paliambela Kolindrou, Revenia), where the earliest phases date to 6600 cal BCE [26, 27, 29].

There is a remarkable increase in the number of settlements from the Middle Neolithic onward, which could be interpreted either as a continuation of migration even after the beginning of the Neolithic, or an increase in population size; of course the two are not mutually exclusive.

The record of human skeletal remains is sparse for these periods. In the Mesolithic, human remains have been identified in the Theopetra cave [56], Franchthi cave [57] and the Mesolithic open site of Maroulas on the island of Kythnos [58]. Single burials predominate at Early Neolithic sites, and are generally placed in a flexed position inside the settlements (e.g. Argissa-Magula, Kefalovryssou,

Sesklo) [59]. Exceptions are the sites of Nea Nikomedeia with 23 single and 2 double burials and the 15 cremations of Soufli-Magoula, both found inside the settlements [60]. In the Middle Neolithic period, human remains are very sparse. Single skeletal elements are known from sites in southern Greece, such as the 30 burials from the Middle Neolithic phase of the Franchthi site (both the cave and the Paralia) found inside pits [57]. The mortuary practices changed to some extent during the Late and Final Neolithic. Cemeteries appear mainly outside the settlements and involve primary and secondary interments and cremations. Plateia-Magoula-Zarkou in Thessaly, dated to 5300-4800 BCE, consists of an organized cemetery [61]. In the cave of Alepotrypa, on the Diros Bay in South Peloponnese, occupied from 5000 to 3200 BCE, inhumations alongside secondary burials and cremations have been found [62, 63]. In the Aegean islands, organized cemeteries have been discovered on the small island of Yiali, on the island of Euboea outside the Late Neolithic site of Tharrounia, and on the island of Kea in the cemetery of Kephala (Angel, in [64]). In northern Greece, disarticulated and fragmented human remains as well as a few articulated inhumations disposed of in ditches located at the margins of the settlements appear to be the normal burial practice during the Late Neolithic period (e.g. LN Makriyalos, LN Toumba Kremastis Koiladas, Paliambela Kolindrou) [65].

## Archaeological background of the Greek samples

### Mesolithic individuals

#### Theopetra Cave, Thessaly

The cave of Theopetra is situated in the Thessalian plain on the north side of a limestone formation almost 100 m above the plain and 300 m above sea level. The cave lies between the edge of the plain and the foothills of the East Pindus mountains and it is the westernmost prehistoric site of the Thessalian plain. The cave has a quadrilateral shape and measures approximately 500 m<sup>2</sup>. Excavations started in 1987 and, after fourteen excavation periods, ended in 2002 [66]. A second phase of excavations was conducted from 2005 to 2008. The archaeological remains are dated from the Middle Palaeolithic ( $46330 \pm 1590$  BP) to the Late Neolithic [67, 68], revealing a long sequence of deposits that extend across the Pleistocene-Holocene boundary, with the oldest dates ranging from 110-140 ka [69, 70]. The importance of the site lies in its Mesolithic layers [56] that bridge the gap between the numerous Neolithic settlements of the Thessalian plain and the open-air Palaeolithic findings. The samples analyzed (Theo1 and Theo5) belong to Mesolithic burials found *in situ* in the cave. Theo1 belongs to a subadult female [71] but the anthropological study of Theo5 has not yet been completed. The sample Theo1 was dated during the excavation period ([72] - see table 1: H6, Human skeleton, burial *in situ*). The estimated age BP is  $8070 \pm 60$ . This radiocarbon dating result was calibrated with OxCal v4.2.2 using INTCAL13 to a corresponding age of 7288-6771 cal BCE. The second sample Theo5 was dated at the Research Laboratory for Archaeology at the University of

Oxford to an age of  $8549 \pm 40$  BP and calibrated with OxCal v4.2.2 using INTCAL13 at 7605-7529 cal BCE.

## **Early Neolithic Period**

### **Revenia, Macedonia**

The flat-extended settlement “Revenia” Korinou lies in a small valley in Pieria, central Macedonia, northern Greece, and covers at least 4 ha, as the surface finds suggest. Rescue excavations conducted at the settlement during 2001-2004, revealed an extraordinarily dense series of pits and other features cut into bedrock [73] and provisionally dated, on the basis of ceramic typology, to a period spanning the Early Neolithic and the beginning of the Middle Neolithic [74]. The pits vary significantly in their shape and dimensions, with some identifiable as semi-subterranean dwellings on the basis of associated postholes and (occasionally burnt) clay remains of superstructure. The fills of individual pits vary significantly, with some being particularly rich in ceramics, others in chipped stone or ground stone tools, and others in animal bone or marine shell. Rescue excavations also revealed the foundations (postholes and wall ditches) of at least one above-ground rectangular building. These well-defined contexts have yielded unusually rich Early Neolithic and Early Middle Neolithic assemblages of several categories of artefacts and ecofacts. Six articulated burials in a flexed position and other five ‘deviant’ burials were found within the investigated area [75].

The sample Revenia 5 (Rev5) was taken from the petrous bone of burial number 2 recovered in a trial trench. The skeleton belongs to a female who was 30-40 years of age. The sample was dated at the Curt-Engelhorn-Zentrum of Archaeometry (Mannheim, Germany). This radiocarbon dating result was calibrated with OxCal v4.2.2 using INTCAL13. For sample Rev5 (Lab number MAMS23036) a  $^{14}\text{C}$  age of  $7505 \pm 25$  BP was estimated resulting in a 95% calibrated range of 6438-6264 cal BCE.

## **Late/Final Neolithic Period**

### **Paliambela, Macedonia**

The Neolithic settlement at Paliambela is situated in the rolling landscape of the coastal lowlands of Pieria in Central Macedonia. An ongoing excavation of the site started in 2000 as a joint project between the Universities of Thessaloniki (Greece) and Sheffield (UK). The site is a low mound with Neolithic deposits exceeding 3 m, spanning the period from the Early to the Final Neolithic. The total excavated area exceeds well over 500 m<sup>2</sup>.



In the earlier phases of the Neolithic, it seems that Paliambela was a flat, extended site, while in the later phases the site took the form of a mound. It does not, however, represent a classic tell village of densely packed houses [29]. On the basis of the first preliminary reports, at least one deep ditch dug in natural bedrock, and a series of pits containing pottery, chipped stone and bone tools and other findings, belong to the Early Neolithic Period (6600-5900 BCE) [76]. During the Middle Neolithic (5900-5400 BCE) the settlement was also encircled by at least one deep, wide ditch. ‘Domestic’ architecture in this phase was comprised of closely-set rectangular buildings separated by cobbled yards. There is some evidence that habitation shifted across the site during the Middle Neolithic. Pottery from this phase links the settlement with both the southern Balkans and Thessaly [77].

Late Neolithic [LN] (5400-4700 BCE) deposits found on top of the hill have been heavily eroded and disturbed by ploughing. It seems, however, that during the LN the settlement was encircled by a pair of stone enclosures [78]. Based on the presence of pottery, the site was inhabited during the whole LN period up to the Final Neolithic (4700-3300 BC). Apart from black burnished and black topped vessels with characteristics of the LN period, carinated and conical shapes, and painted pottery of the ‘Dimini’ style is also found [79]. Other findings include a large number of ground, polished, and chipped stone and bone tools, and figurines.

Most of the human skeletal material belongs to scattered postcranial and a few cranial remains [65]. The sample Paliambela 7 (Pal7) was taken from the petrous bone and belongs to a 7-10 year old individual; morphological sex determination was not possible due to a lack of necessary anatomical elements. The sample was dated at the Curt-Engelhorn-Zentrum of Archaeometry (Mannheim, Germany). This radiocarbon dating result was calibrated with OxCal v4.2.2 using INTCAL13. For Paliambela (Lab number MAMS 23037) a  $^{14}\text{C}$  age of  $5559 \pm 29$  BP was estimated resulting in a 95% calibrated range of 4452-4350 cal BCE.

## **Final Neolithic Period**

### **Kleitos**

The archaeological site of Kleitos is situated in western Macedonia, in northern Greece. In 2006, an extensive rescue excavation took place covering an area of 7.5 hectares. This brought to light two neighboring Neolithic settlements (Kleitos 1 and 2) and findings of the Bronze Age, the Hellenistic, Late Roman and Early Byzantine Period.

Kleitos 1 is a flat Neolithic settlement which was inhabited during the early phases of the Late Neolithic period (second half of the 6<sup>th</sup> and early 5<sup>th</sup> millennium BCE) and is one of the very few sites in the Balkan region that has been excavated throughout its entire area [80]. Ten quadrangular ground-floor buildings made of a wooden framework and covered with clay have been identified. Inside the buildings were structures designed for food preparation and storage, along with clay vessels, tools, and pots. The buildings had been destroyed by fire and none of them preserved more

than three construction phases. Among and around the building complex, a variety of workshops, storage areas, and refuse pits were uncovered. The settlement covers an area of approximately 2 ha and is bounded by a system of ditches and wooden enclosures. Inside Building C, amongst and outside the houses, 16 Neolithic graves were found. The burials were single inhumations in contracted positions and there was one cremation in an urn [80].

The sample Kleitos 10 (Klei10) was taken from a petrous bone of a 15 year old individual interred in grave 9. Morphological sex determination was not possible due to a lack of reliable criteria. The sample was dated at the Curt-Engelhorn-Zentrum of Archaeometry (Mannheim, Germany). The results were calibrated as described above. For Klei10 (Lab number MAMS 23038) a  $^{14}\text{C}$  age of  $5559 \pm 22$  BP was estimated resulting in a 95% calibrated range of 4230-3995 cal BCE.

## SI2 Sample Preparation

Susanne Kreutzer, Zuzana Hofmanová, Melanie Strobel, Laura Winkelbach & Amelie Scheu

The seven Mesolithic and Neolithic Northern Aegean samples (Table S1) were analyzed in the dedicated ancient DNA facilities of the Palaeogenetics Group, Institute of Anthropology, Mainz. Details on the decontamination procedures and sample preparation of bones and teeth are described in [35].

As shown previously, the inner core of the petrous bone is likely to contain a high amount of endogenous DNA [81]. Therefore, the petrous bone was separated from the temporal bone at the base of petrous part (where it fuses with the squamous and mastoid part). The outer surface of the petrous bone was removed with a saw (Electer Emax IH-300, MAFRA) in order to identify the densest parts of this bone fragment. All parallel canals, fossa, sinuses and canaliculi were cleaned of dirt by sandblasting (P-G 400, Harnisch & Rieth, Winterbach, Germany). The densest inner parts of the petrous bone were sawn into small cubes and UV-irradiated for 30 minutes per side before being milled into fine powder (MM200, Retsch).

Milling controls as described in [35] were processed in parallel to control for the decontamination procedure of the devices used. These controls were treated as samples for all subsequent steps, including extraction, library preparation and quantification.

**Table S1:** *Sample summary (n.d. - not determined, \* anthropological sex determination not feasible due to the lack of morphological criteria)*

sample	age (cal BCE)	site, geograph. region, country	skeletal element	anthropological information
Theo 1	7288-6771	Theopetra, Thessaly, Greece	tibia	subadult female
Theo 5	7605-7529	Theopetra, Thessaly, Greece	tibia	n.d.
Rev 5	6438-6264	Revenia, Northern Greece	petrous bone	30-40 year old female
Bar 8	6212-6030	Barcin, Western Anatolia	petrous bone/ tooth	middle aged female
Bar 31	6419-6238	Barcin, Western Anatolia	petrous bone	male adult
Pal 7	4452-4350	Paliambela, Northern Greece	petrous bone	7-10 year old individual*
Klei 10	4230-3995	Kleitos, Northern Greece	petrous bone	15 year old individual*

### Extraction

Extraction was performed similarly to [35]. For the lysis step, EDTA (10ml-14ml, 0.5M, pH8; Ambion/Applied Biosystems, Life technologies, Darmstadt, Germany), N-laurylsarcosine (250 $\mu$ l, 0.5%; Merck Millipore, Darmstadt, Germany) and proteinase K (30 $\mu$ l, 18U $\mu$ l; Roche, Mannheim, Germany) were added to the powdered sample. The EDTA volume was adjusted by the amount and density of the bone/tooth powder being extracted (200-500 mg).

The lysis solution was incubated on rocking shakers at 37°C until the powder was dissolved and the DNA was released from the bone matter. The DNA was isolated via phenol/chloroform/isoamyl

alcohol (25:24:1, Roth, Karlsruhe, Germany) extraction, then desalted by stepwise washes with HPLC-water (2-12 ml) and concentrated to approximately 200 $\mu$ l using Amicon Ultra-15 Centrifugal Filter Units (Merck Millipore, Darmstadt, Germany).

For deeper shotgun sequencing approaches, additional extractions were performed, resulting in 2-4 extractions (I-IV, see Table S3 for details) per sample. Blank controls were processed with every extraction.

## Library preparation

All libraries were prepared according to [82] with slight modifications. USER<sup>TM</sup> (NEB) treatment of the DNA extract was performed prior to library preparation for 18 out of 51 libraries, but for initial screenings were left untreated, and the full damage patterns were used to authenticate the age of the sample (see Table S3 for details). We used hybridized adapters P5 and P7 (IDT, Leuven, Belgium) at a concentration of 1,25 $\mu$ M. Amplifications of all libraries were performed with AmpliTaq Gold<sup>®</sup> DNA Polymerase (Applied Biosystems) in at least three PCR parallels and 10-16 cycles. Indices on both sides of the library molecule were added simultaneously. Double indexing was performed according to [82], including additional index sequences from the Nextera<sup>XT</sup> index Kit v2 (Illumina).

Blank controls were processed with every library step and each PCR. Purification during library preparation was conducted using the MinElute PCR Purification Kit (Qiagen, Hilden, Germany) and amplified libraries were purified with MSB<sup>®</sup> Spin PCRapace (Invitex, Stratec Molecular, Berlin, Germany). Library concentrations were measured by Qubit<sup>®</sup> Fluorometric quantitation (dsDNA HS assay, Invitrogen) and fragment length distributions of libraries were estimated on the Agilent 2100 Bioanalyzer System (HS, Agilent Technologies) following the manufacturer's protocols. Occasional primer dimers of <100 bp length were removed prior to sequencing by additional purification with Agencourt<sup>®</sup> AMPure<sup>®</sup> XP beads (Beckmann Coulter).

## Quality assessment of samples

The endogenous DNA content of an ancient sample depends on variety of factors, including the type of skeletal element, age, climate, post-burial taphonomic processes and storage conditions after excavation. Not all samples contain enough endogenous DNA to generate sufficient genome-wide coverage after shotgun sequencing for subsequent population genetic inference. For this reason, DNA quality and quantity was gauged using a combination of shallow shotgun sequencing and quantitative real-time PCR.

## Quantitative real-time PCR

A direct measurement of the number of molecules in a library can be obtained by quantitative real-time PCR (qPCR) on library fill-in products as described in [83]. To detect unique molecules, not biased by PCR duplicates, a qPCR measurement was performed with KAPA Sybr Fast Universal Mastermix (PeqLab, VWR International) on a Step One Plus<sup>TM</sup> Real-Time PCR system (Applied Biosystems, Thermo Fisher Scientific) with primer pair IS7/IS8. A synthetic standard (artificial library molecule, 89bp) was run with every measurement setup as described in [84]. Resulting molecule numbers were corrected for fragment lengths of each library, as determined by Bioanalyzer measurement. Measurements of libraries prepared from the milling, extraction and blank controls allowed us to determine the quantity of contaminating molecules incorporated during each step of the protocol. A theoretical contamination level was calculated for each sample by comparing each library to its respective blank controls. Samples with high unique copy numbers are assumed to be influenced by contaminating molecules to a lesser extent. For deep shotgun sequencing, the molecule number allowed to estimate the number of libraries and PCR parallels needed to reach the desired coverage without depleting library complexity (i.e. sequencing duplicates; see Figure S3 in [85] for detail).

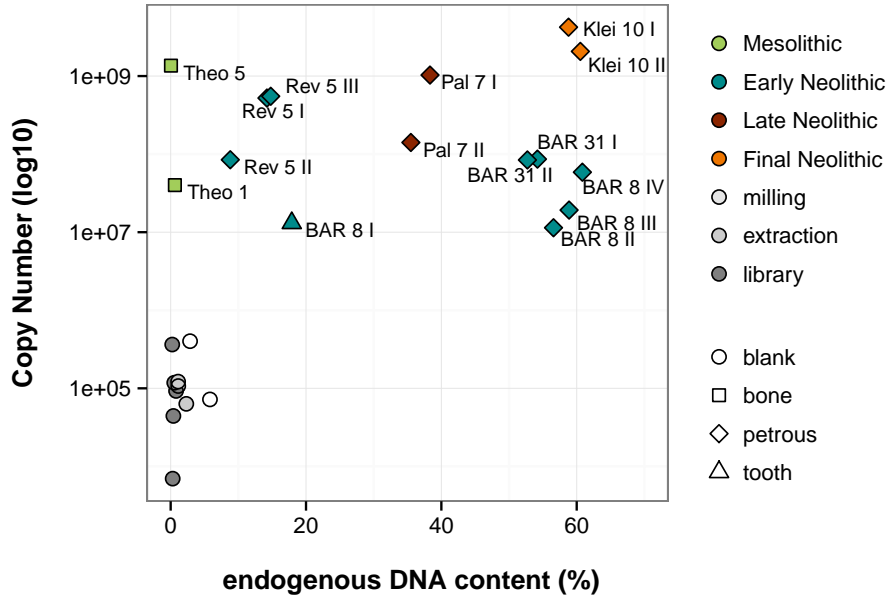
## MiSeq screening

Prior to deeper shotgun sequencing at least one library per sample per extract was sequenced on Illumina's MiSeq with 50bp read length (single end) resulting in 1-2 million reads per sample. This quantity of sequencing data was sufficient to evaluate sample authenticity (by identifying the expected damage patterns) and sample quality (by determining the endogenous DNA content).

MiSeq sequencing data were analysed with the pipeline described in SI3. The endogenous DNA content of a sample was calculated after duplicate removal.

## Results sample quality

The combination of endogenous DNA content and the quantity of unique molecules per library is highly informative of sample complexity. Figure S1 and Table S2 provide a summary of the quality and contamination assessment of all samples. The extracts of the two Mesolithic samples Theo1 and Theo5 contain only between 0.05 and 0.62% endogenous DNA. While they are not suitable for whole genome shotgun sequencing, they were selected for targeted enrichment of the mitochondrial genome. The five Neolithic samples (Bar8, Bar31, Klei10 and Pal7 and Rev5) show endogenous DNA contents between 8.80 and 60.83 %. On the basis of these results, we developed a library and PCR pooling strategy as outlined in Table S3 to reach genome-wide coverage for each sample between 1 and 7 x.



**Figure S1: Quality Assessment.** *Log(10) transformed copy numbers per  $\mu\text{l}$  library of each extract screened and corresponding blank controls (measured after fill-in reaction during library preparation) are plotted against their endogenous DNA content.*

## Blank controls

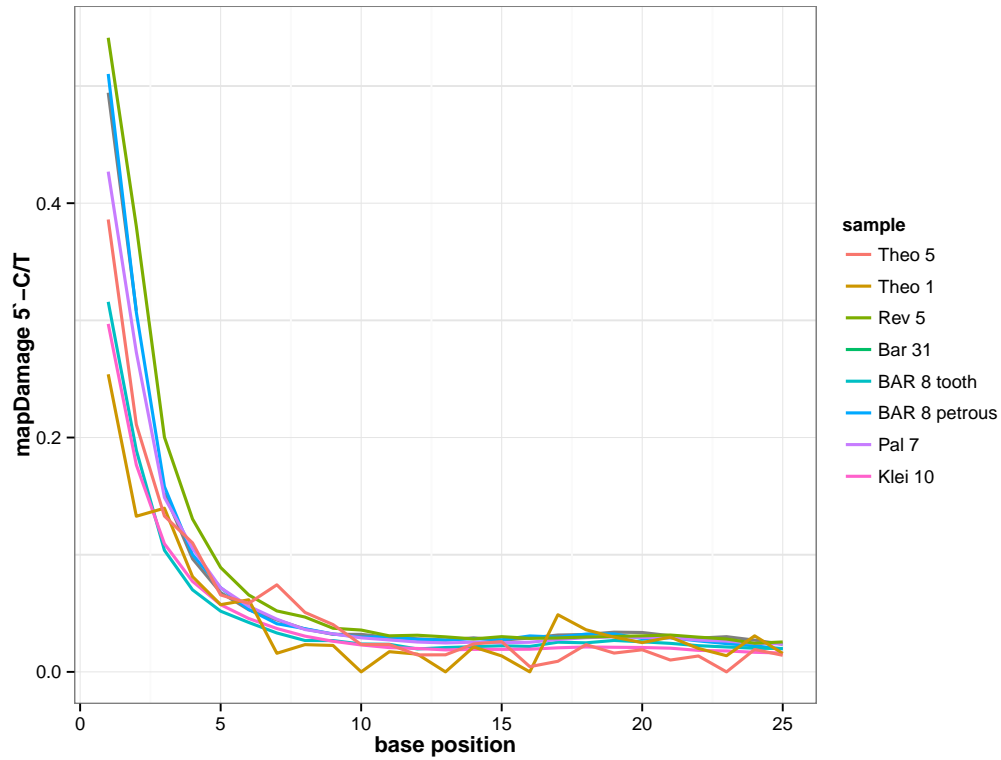
We determined the theoretical contamination level influencing our samples from all stages of lab work. We calculated the number of endogenous molecules per  $\mu\text{l}$  of library after fill-in step in samples and controls by multiplying the unique copy number by the endogenous DNA proportion (see “column unique 2” in Table S2). By comparing molecule numbers of the samples to their corresponding milling controls (the first control taken during sample preparation) we obtain a maximum contamination level of 0.53% for Theo 1 (lowest sample complexity) and a minimum contamination level of 0.00005% for Klei 10 I (highest sample complexity). By inspecting the blank controls along each step we find only a slight accumulation of molecules mapping the human genome. Library controls contain a range of 100-1000 human molecules per  $\mu\text{l}$ , extraction controls show approximately 1500 molecules per  $\mu\text{l}$  and milling controls contain a range of 1000-10000 molecules per  $\mu\text{l}$ . In the samples, we find endogenous molecule counts per  $\mu\text{l}$  library between  $2.51 \cdot 10^5$  for Theo 1 and  $2.40 \cdot 10^9$  for Klei 10.

**Table S2:** *Quality assessment of samples and contamination level of blank controls (unique 1 = unique molecules per  $\mu\text{l}$  library, unique 2 = endogenous molecule count per  $\mu\text{l}$  library, mD5 = deamination rate of first position at 5' molecule end)*

sample	skeletal element	unique 1	endogenous DNA content (%)	mD5	unique 2
Rev 5 I	petrous	5.23E+08	14.23	0.541	74416316
Rev 5 II	petrous	8.47E+07	8.81	0.537	7461209
Rev 5 III	petrous	5.51E+08	14.79	0.510	81471468
Klei 10 I	petrous	4.23E+09	58.79	0.296	2486983763
Klei 10 II	petrous	2.07E+09	60.54	0.262	1253117319
Pal 7 I	petrous	1.03E+09	38.32	0.429	394716616
Pal 7 II	petrous	1.41E+08	35.49	0.386	50046282
Theo 1	bone	4.01E+07	0.62	0.386	250714
Theo 5	bone	1.36E+09	0.05	0.253	618845
BAR 31 I	petrous	8.63E+07	54.20	0.494	46772311
BAR 31 II	petrous	8.42E+07	52.71	0.447	44384667
BAR 8 I	tooth	1.31E+07	17.90	0.315	2344287
BAR 8 II	petrous	1.14E+07	56.57	0.510	6449094
BAR 8 III	petrous	1.92E+07	58.87	0.500	11303622
BAR 8 IV	petrous	5.87E+07	60.84	0.507	35711588
milling	blank	1.07E+05	1.12	-	1203
milling	blank	7.21E+04	5.81	-	4189
milling	blank	4.02E+05	2.87	-	11517
extraction	blank	1.22E+05	1.10	-	1340
extraction	blank	6.33E+04	2.31	-	1462
library	blank	6.97E+04	0.30	-	21
library	blank	4.43E+03	0.40	-	177
library	blank	9.26E+04	0.80	-	741
library	blank	3.64E+05	0.23	-	837
library	blank	1.18E+05	0.53	-	625

## Damage patterns

The damage patterns identified by Mapdamage2.0 [86] are visualized in Figure S2. All samples show damage patterns typical of ancient DNA. C to T transitions occur with a frequency of 25-54% at the first base of the 5' end.



**Figure S2:** Damage patterns at 5'end (C to T transition) from shallow 50bp MiSeq single-end sequencing.



## Deep shotgun sequencing

Additional libraries were prepared for deeper shotgun sequencing with adjustments made according to the previously estimated sample quality. For low copy number samples, we decreased adapter ( $0.75\mu\text{M}$ ) and primer ( $0.5\mu\text{M}$ ) concentrations. We lowered the cycle number during PCR amplification and increased the number of parallels from 3 to 4-6. Portions of the extracts were treated with USER<sup>TM</sup> enzyme before library preparation as described elsewhere [87] (detailed information in Table S3). For Rev5, 4 out of 11 libraries, for Bar31, 2 out of 8 libraries and for Bar8, 12 out of 23 libraries were prepared from USER<sup>TM</sup> treated extract.

Selected samples were sequenced on Illumina Hiseq2500 (either a 100bp paired-end or single-end run) and Illumina NextSeq (75bp paired end).

Bar8 was sequenced on 6 lanes (one lane was paired-end) and Bar31 on 3 lanes (all single-end). Due to the low quantity of unique copy numbers in Barcin samples, we increased the number of libraries for sequencing. A total of 23 separately indexed libraries produced in 123 independent PCR reactions were pooled for Bar8. A total of 8 separately indexed libraries produced in 45 independent PCR reactions were pooled for Bar 31.

Rev5 was sequenced on NextSeq (low output) and three lanes of HighSeq (two lanes paired-end, one lane single-end). Pal7 and Klei10 were sequenced each on one lane of HighSeq (paired-end).

For additional details regarding library preparation and sequencing see Table S3.

**Table S3:** Sample preparation and shotgun sequencing strategy. # Extr. = number of extractions per sample, # Lib. = number of libraries per sample/extraction, %endo = percentage of reads mapped to hg19, mD5 (C/T) = deamination rate at the 5'-end, coloring scheme visualizes different sequencing runs/machines/modes

sample	material	# extr.	# lib.	USER treatment	PCR cycles (parallels)	% endo	mD5 (C/T)	shotgun sequencing I	shotgun sequencing II	shotgun sequencing III
Rev 5	petrous	I	1	n	16 (4)	14,23%	0.541	Next Seq (Low Output, PE)	HiSeq (one lane, PE)	
Rev 5	petrous	I	2	n	14 (4)					
Rev 5	petrous	I	3,4	n	14 (4)					
Rev 5	petrous	II	5,6	n	14 (4)	8,81%	0.537			
Rev 5	petrous	III	7	n	14 (4)	14,79%	0.510	sequenced on MiSeq only		
Rev 5	petrous	III	8-11	y	14 (4)				HiSeq (one lane, SE)	
Klei 10	petrous	I	1	n	16 (4)	58,79%	0.296			
Klei 10	petrous	I	2	n	14 (4)			HiSeq (one lane, PE)		
Klei 10	petrous	II	3	n	14 (4)	60,54%	0.262			
Klei 10	petrous	II	4	n	14 (4)					
Pal 7	petrous	I	1	n	16 (4)	38,32%	0.429			
Pal 7	petrous	I	2	n	14 (4)			HiSeq (one lane, PE)		
Pal 7	petrous	II	3	n	14 (4)	35,49%	0.386			
Pal 7	petrous	II	4	n	14 (4)					
BAR 31	petrous	I	1	n	12 (6)	54,20%	0.494			
BAR 31	petrous	I	2	n	12 (6)			HiSeq (three lanes, SE)		
BAR 31	petrous	I	3	n	12 (6)					
BAR 31	petrous	I	4	n	12 (6)					
BAR 31	petrous	I	5	n	12 (6)					
BAR 31	petrous	II	6	n	12 (3)	52,71%	0.447			
BAR 31	petrous	II	7	y	10 (6)					
BAR 31	petrous	II	8	y	10 (6)					
BAR 8	tooth	I	1	n	22 (6)	17,90%	0.315			
BAR 8	tooth	I	2	n	18 (3)					
BAR 8	tooth	I	3	n	18 (3)					
BAR 8	tooth	I	4	n	18 (3)					
BAR 8	petrous	II	5	n	12 (6)	56,57%	0.510	HiSeq (one lane, PE)		
BAR 8	petrous	II	6	n	12 (6)					
BAR 8	petrous	II	7	n	12 (6)					
BAR 8	petrous	II	8	n	12 (6)					
BAR 8	petrous	II	9	n	12 (6)					
BAR 8	petrous	III	10	n	12 (3)	58,87%	0.500			
BAR 8	petrous	III	11	y	12 (6)					
BAR 8	petrous	III	12	y	12 (6)					
BAR 8	petrous	III	13	y	12 (6)					
BAR 8	petrous	III	14	y	12 (6)					
BAR 8	petrous	III	15	y	12 (6)					
BAR 8	petrous	III	16	y	12 (6)					
BAR 8	petrous	IV	17	n	12 (3)	60,84%	0.507			
BAR 8	petrous	IV	18	y	10 (6)					
BAR 8	petrous	IV	19	y	10 (6)					
BAR 8	petrous	IV	20	y	10 (6)					
BAR 8	petrous	IV	21	y	10 (6)					
BAR 8	petrous	IV	22	y	10 (6)					
BAR 8	petrous	IV	23	y	10 (6)					

## Mitochondrial capture

The Mesolithic samples from Theopetra Cave contained low amounts of endogenous DNA (Theo 1: 0.6%, Theo 5: 0.05%). Therefore, we performed whole mitochondrial genome enrichment of these samples with Agilent's SureSelect<sup>XT</sup> in solution target enrichment kit (custom design) [88]. 120bp RNA baits covering the whole mitochondrial genome were built with 8-fold tiling from the main mitochondrial haplogroups given by Phylotree 8 (mtDNA tree built14; 5. April 2012, [89]). To ensure good coverage in the control region additional baits (10-fold tiling) for this region (15900-16569 & 1-600) were designed on 18 different haplogroups and additional baits were introduced manually to close the circular structure of the mitochondrial genome. To correct for GC bias, extra baits for GC-low regions of the target sequence were introduced.

Capture was performed according to a modified version of the manufacturer's protocol. Modifications included:

- (i) dilution of RNA-baits set in the hybridization reaction,
- (ii) in house preparation of hybridization and washing buffers [90],
- (iii) design of blocking adapter sequences around index sequences for the reverse strand only,
- (iv) setting washing temperature to 57°C [91] and
- (v) amplification in three PCR parallels after capture.

Purification after PCR was performed with MSB<sup>®</sup> Spin PCRapace (Invitex, Stratec Molecular, Berlin, Germany).

In order to obtain adequate coverage to allow for accurate SNP calling and contamination estimates of the mitochondrial DNA, a double capture of pooled libraries of each sample was performed. For both Theo 1 and Theo 5, the complete product from the first reaction was used in the second reaction. Only the cycle number differed between the two rounds of target enrichment.

The samples were then sequenced on Illumina machines (MiSeq 50bp single end, HiSeq2500 100bp paired end) yielding 1-2 million sequence reads.

## SI3 Read Processing

Christian Sell, Susanne Kreutzer, Zuzana Hofmanová & Jens Blöcher

Illumina HiSeq runs were carried out at the sequencing facilities of the University of Mainz (Institute of Molecular Genetics, Mainz, Germany). Screening runs were subjected to the Illumina MiSeq outsourced to StarSEQ GmbH (Mainz, Germany). FASTQ files were generated by the sequencing facility. Samples were demultiplexed with a threshold of 0 mismatches per index read.

Adapters were trimmed at the 3' end of each sequence if the adapter sequence and the read were at least 90% identical, while allowing a minimum adapter length of 1 bp [92]. Reads showing a base quality score  $\leq 15$  in more than 5% of the bases of a sequencing read were removed from the dataset [92]. For paired end reads, a custom python script was applied to order reads in pairs by their names. Subsequently, the ea-utils package [93] was used to merge overlapping read pairs, with the default parameters of  $\geq 6$ bp overlap and 92% sequence identity in the overlapping region.

Reads were aligned using BWA aln [94] to the human reference build GRCh37/hg19 or the revised Cambridge Reference Sequence (rCRS) for mitochondrial capture data with the default parameters. The “MarkDuplicates” command from the Picard tools package (picardtools, <http://broadinstitute.github.io/picard>) was used to remove duplicate reads. To sort and index the alignments, the Samtools package [95] was used. Additionally, all sequences with a length  $< 30$ bp were removed from the alignment with NGSUtils [96].

Local realignment was performed using GATK v. 3.3.0 [97] according to the recommendations of the GATK development team using “IndelRealigner”. We estimated 5'- and 3'-deamination patterns in aligned sequence reads using MapDamage 2.0 [86]. Read groups were set with the “AddOrReplaceReadGroups” command from the Picard tools package in order to differentiate sequencing machines/modes and various post-mortem damage (PMD) patterns, since we were using combined sequencing data of a sample including USER<sup>TM</sup> treated and untreated extracts. Genotyping of the Aegean samples is described in SI5.

Different sequencing runs of the same sample were merged using the Samtools package. Mean coverage and standard deviation of each ancient genome were estimated in non-overlapping sliding windows of 1Mbp.

In addition to the variant calls described in SI5, a variant call with the GATK “HaplotypeCaller” was performed. Only regions with a  $> 2x$  coverage were selected for calling using GATK’s “CallableLoci”. To allow for sites being called as the reference allele to be displayed, the emit reference confidence option (-ERC) was used with BP\_RESOLUTION in GATK’s HaplotypeCaller.

## Results

Table S4 and S5 provide an overview of each sample by number of reads sequenced, sequencing run, number of reads mapped to the human reference genome, coverage, damage patterns and mean fragment lengths. Fragment length distributions and corresponding deamination rates are listed in Table S5 and parts of the dataset are displayed in Figure S3.

**Table S4:** Overview sequencing results for the different sequencing runs, *HS* = *HiSeq* (100bp single end), *HP* = *HiSeq* (100bp paired end), *NS* = *Next Seq* (75bp paired end), *PE* = paired end

sample	seq. mode	reads total	joined reads (PE)	mapped hg19	coverage	% covered hg19 1x (2x)
Rev 5	NS	136240888	49660203	6668229		
Rev 5	HP	671583640	246033810	36568311	1.16 ± 0.73	58.46 (32.19)
Rev 5	HS	232833397	-	32815612		
Bar 31	HS	358863810	-	150475345		
Bar 31	HS	116924532	-	47032756	3.66 ± 2.04	82.32 (71.86)
Bar 8	HP	88093160	37265273	6493549		
Bar 8	HP	164981292	68408997	37149128		
Bar 8	HS	50824229	-	6946472	7.13 ± 4.56	86.65 (83.86)
Bar 8	HS	170028128	-	85709753		
Bar 8	HS	496840553	-	238736411		
Pal 7	HP	420078186	147638025	60080727	1.28 ± 1.01	65.18 (37.08)
Klei 10	HP	366685662	139638436	87116507	2.01 ± 2.20	76.33 (55.84)

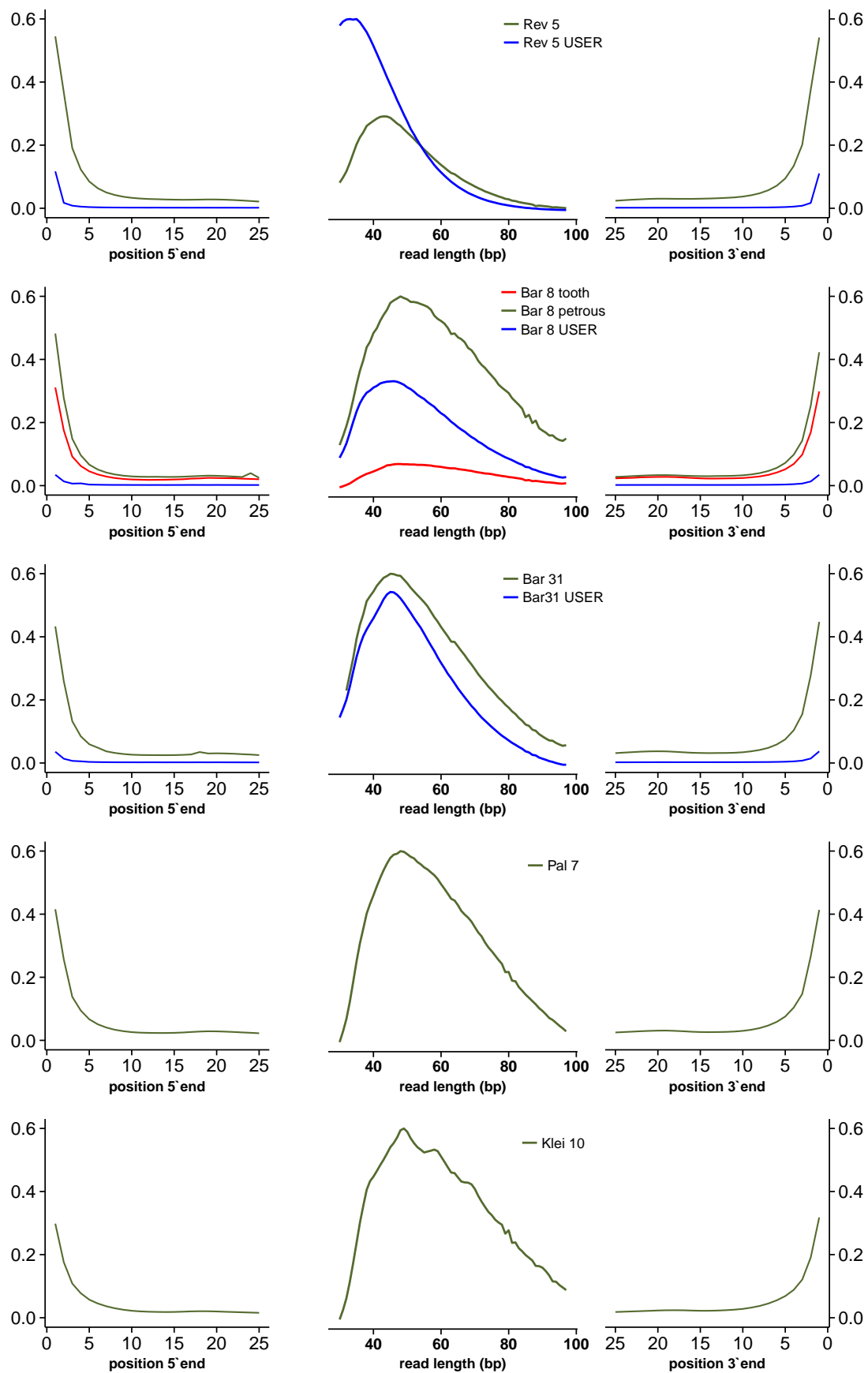
**Table S5:** *mapDamage* results for the different sequencing runs, *HS* = *HiSeq* (100bp single end), *HP* = *HiSeq* (100bp paired end) *NS* = *Next Seq* (75bp paired end), *mD5/mD3* = deamination rate of first position at particular molecule end

sample	seq. mode	USER treatment	skeletal element	mD5	mD3	read length $\phi$
Rev 5	NS	-	petrous	0.546	0.545	49.41 ± 15.03
Rev 5	HP	-	petrous	0.555	0.541	51.33 ± 14.68
Rev 5	HS	+	petrous	0.130	0.120	43.74 ± 11.64
Bar 31	HS	-	petrous	0.453	0.443	59.48 ± 19.43
Bar 31	HS	+	petrous	0.033	0.036	54.32 ± 16.48
Bar 8	HP	-	tooth	0.304	0.320	69.21 ± 28.26
Bar 8	HP	-	petrous	0.495	0.505	66.31 ± 25.94
Bar 8	HS	-	tooth	0.305	0.292	66.16 ± 21.37
Bar 8	HS	-	petrous	0.489	0.434	62.94 ± 20.20
Bar 8	HS	+	petrous	0.057	0.054	57.03 ± 18.38
Pal 7	HP	-	petrous	0.425	0.436	66.62 ± 25.37
Klei 10	HP	-	petrous	0.280	0.298	72.22 ± 28.35

While small fragment lengths of around 50 bp and strong deamination patterns are expected for ancient DNA, we find particularly elevated deamination rates of up to 56% at the 5' ends of sample Rev5. USER treatment still left 13% damage while reducing the mean fragment length from 51 to 44 bp. Around 20% of all reads were smaller than 30 bp, and hence filtered out. During downstream

analyses, we accounted for the pronounced damage patterns by developing and applying a novel SNP caller that recalibrates the quality scores of damaged bases as described in SI5, thereby allowing us to analyze sequence data from both USER and non-USER-treated extract together.

For sample Bar8, we also observe shorter fragments (around 5bp less) and elevated 5' end deamination rates (around 20% more) from the petrous bone compared to the tooth. Despite the relatively high endogenous DNA content of petrous bone material, which we can confirm in the present study, they appear to be more susceptible to PMD. The reason for this observation is unclear, and will require further investigation of how taphonomic processes act on different skeletal elements.



**Figure S3:** Deamination patterns and fragment length of different sequencing runs per sample. Right: Deamination pattern 5'-end, middle: Sequenced fragment length extracted from alignments, left: Deamination pattern 3'-end.

## SI4 Analysis of Uniparental Markers and X Chromosome Contamination Estimates

Susanne Kreutzer, Rui Martiniano, Zuzana Hofmanová & Christian Sell

Read processing was performed as in SI 3. Local realignment, base quality recalibration, and SNP calling were implemented using the GATK program v. 3.3.0 [97]. Indel realignment and base recalibration were performed according to the recommendations of the GATK development team using the IndelRealigner and BaseRecalibrator tools using the GATK resource bundle callset.

Final variant calling was performed using the GATK UnifiedGenotyper [97] with ploidy set to 1. Variants were then filtered for  $>5x$  coverage and a quality of Phred quality  $>50$ .

Haplogroups were estimated using Haplofind [98]. To generate the consensus sequence of the complete mitochondrial genome SAMtools mpileup/bcftools/vcfutils [95] was applied to compile FASTQ data with subsequent transformation to FASTA format. Sites that were not assayed are filled with N.

### Contamination estimates

To estimate mitochondrial contamination, a likelihood-based method described in [99] was used (Table S6). Additionally, we applied the genotype caller GATK HaplotypeCaller [97] in combination with a custom python script to detect positions that could not be determined reliably (i.e. heterozygous genotype was called) and manually examined these positions in order to identify possible contamination.

### Mitochondrial shotgun data

Mitochondrial regions were extracted from alignments to GRCh37/hg19 and realigned against the revised Cambridge Reference Sequence (rCRS). Contamination estimates, SNP calling and haplogroup determination were conducted as described above (Table S7). The distribution of mitochondrial genome coverage was highly uneven. Regional coverage varied between extremely high coverage (up to 746x) and no coverage at all. Because we extracted the mitochondrial data from the full shotgun alignment to the complete genome, this result could be due to mitochondrial reads mapping with higher or equal quality to nuclear regions known as nuclear mitochondrial sequences (numts).



**Table S6: Mitochondrial Capture results of Mesolithic samples***(mD5: deamination rate of first position at 5', Contamination estimates were subtracted from 1, HG: haplogroup)*

Sample	Material	reads mapped rCRS	coverage	mD5	Contamination estimate [%]	authentic data [%]	HG
Theo 1	bone	7359	21.7	0.322	1.84 - 6.71	96.4	K1c
Theo 5	bone	12600	62.74	0.344	0.05 - 3.8	99.3	K1c

**Table S7: Mitochondrial genomes extracted from Shotgun data***(Contamination estimates were subtracted from 1, read lengths were extracted from BAM files, HG: haplogroup)*

Samples	reads mapped rCRS	coverage	not covered	authentic data [%]	Contamination estimate [%]	read length	HG
Rev5	11850	34.71	1465	99.96	0.006 - 0.628	48.47	X2b
Pal7	12243	52.60	425	99.96	0.006 - 0.772	70.98	J1c1
Klei10	24522	119.01	129	99.10	0.363 - 1.772	80.13	K1a2
Bar8	60460	220.00	319	99.86	0.744 - 1.619	60.12	K1a2
Bar31	25323	88.90	722	99.96	0.006 - 0.628	58.04	X2m

## Mitochondrial haplogroups

Mitochondrial haplogroups analysed in this study are consistent with previously described lineages from Balkan [100] and Central European [101, 102] first farmer populations. Samples Bar8 and Klei10 are both typed as K1a2, but differ at seven sites among each other. The Mesolithic samples share a similar lineage differing in three positions between each other, one of them in hypervariable-region I.

Interestingly, the two Mesolithic individuals from Theopetra Cave (Thessaly, Greece) display a K1c haplogroup, a lineage which has never been observed in the European Mesolithic context. The <sup>14</sup>C datings of these samples indicate they belong to a period that early Neolithic sites had been absent in the Aegeans and Greece.

The geographically and temporally closest individual from a Mesolithic context analyzed so far is from the Adriatic site Vela Spila (Croatia) dating to 6210-6000 cal. BCE [100] and showing a typical European Mesolithic lineage (U5b2a5). This indicates that the Mesolithic populations of the Balkans share a common ancestry with corresponding populations from Central Europe. However, this scenario is altered in the Greek Mesolithic samples, revealing Central-Anatolian/Near Eastern affinities. Based on the data presented here, the current hypothesis of autochthonous formation of farming cultures in Greece cannot be rejected. Taking into account the sparse existing Mesolithic findings together with the presence of trans-Aegean networks in Mesolithic times, an alternative interpretation could be that a small-scaled migration from Central-Anatolia or the Near East passed through the Aegean in pre-Neolithic times. Genomic data from Mesolithic Aegean individuals would be helpful to address this question. Due to the low amount of endogenous DNA content of samples in this study, additional sampling will be required.

## Y-chromosomal lineage determination

We determined Y chromosomal lineages in ancient male samples with `clean_tree` [103]. This software requires BAM format files as input, calling alleles with SAMtools `mpileup` at given SNP positions (Table S8). We used the same markers that were included with the `clean_tree` software (539 SNPs) for haplogroup determination, which are based on ISOGG 2013 (International Society of Genetic Genealogy; <http://www.isogg.org/>). Table S9 and Table S10 show the Y-chromosomal polymorphisms that define haplogroup G2 and sub-lineages. Although Bar31 presents the derived allele for mutation G2a2b2a1a1-L78, this is inconsistent with the ancestral alleles found at upstream markers G2a2b2-L141.1 and G2a2b2a-P303. Nevertheless we can confidently assign this individual to haplogroup G2a2b because of derived alleles at 2 distinct markers L32 and L30. Regarding Klei10, this sample can be included in the G2a2a1b-L91 sub lineage of the Y-chromosome, which was also identified in the Tyrolean Iceman [104]. G2a derived lineages are common in the Neolithic and have been found in Germany [102, 105], in Southern France [106] and Spain [107]. In present-day populations G2a reaches its highest frequencies in the Georgian and Balkar populations from the Caucasus but exist at relatively low frequencies in Europe [108]. The age of haplogroup G-M201 has been estimated to be 17.000 years [109], while G2 was given earlier date of 12.500 years [110]. Overall, the abundance of G2a lineages in the Neolithic and their absence in samples from preceding time periods suggests a primarily Neolithic expansion into modern Europe.

**Table S8:** *Number of Y-chromosome reads and SNPs for the ancient male samples in this study.*

Sample	chrY reads	chrY SNPs
Bar31	432616	408
Klei10	232245	350

**Table S9:** Y-chromosome polymorphisms for markers downstream haplogroup G2 for sample Bar31  
(A = ancestral, D = derived, Cov = coverage).

Bar31								
Position	Marker	Haplogroup	Mutation	A	D	Cov	Bases	State
17174741	L156	G2	A->T	A	T	2	TT	D
14028148	L31	G2a	C->A	C	A	2	AA	D
23244026	P15	G2a	C->T	C	T	4	TTTT	D
14692227	L32	G2a2b	T->C	T	C	6	CCCCCC	D
15604899	L30	G2a2b	C->T	C	T	1	T	D
14871976	L78	G2a2b2a1a1	C->T	C	T	1	T	D
9985022	L293	G2a1	G->C	G	C	3	GGG	A
22741799	M286	G2a2a1a	G->A	G	A	2	GG	A
21645555	L91	G2a2a1b	G->C	G	C	1	G	A
23989884	L166	G2a2a1b1a	C->A	C	A	1	C	A
2749995	M406	G2a2b1	T->G	T	G	4	TTTT	A
21628300	L90	G2a2b1a	G->A	G	A	2	GG	A
22917995	L14	G2a2b1a	C->T	C	T	2	CC	A
2888608	L141.1	G2a2b2	del->A	T	A	3	TTT	A
21645348	P303	G2a2b2a	T->C	T	C	1	T	A
6738741	CTS417	G2a2b2a1a1a	A->G	A	G	1	A	A
16903025	L1266	G2a2b2a1a2	T->G	T	C	2	TT	A
14231229	L1265	G2a2b2a1a2a	A->G	A	G	5	AAAAA	A
7644368	L1264	G2a2b2a1a2a	A->G	A	G	1	A	A
17423320	L497	G2a2b2a1b	C->T	C	T	1	C	A
17937365	L43	G2a2b2a1b1a1	A->G	A	G	1	A	A
16660759	L42	G2a2b2a1b1a1a	C->A	C	A	1	C	A
17111777	CTS6796	G2a2b2a1b1a2	T->C	T	C	7	TTTTTTT	A
6835545	Z724	G2a2b2a1c1	C->T	C	T	2	CC	A
16596946	CTS5990	G2a2b2a1c1a	A->G	A	G	2	AA	A
18393688	L640	G2a2b2a1c1a1	A->G	A	G	1	A	A
8051187	L1263	G2a3b1a1a1a1a	G->A	G	A	1	G	A
8467136	L183	G2b1	G->C	G	C	2	GG	A
15031385	M283	G2b1a	A->G	A	G	3	AAA	A

**Table S10:** *Y-chromosome polymorphisms for markers downstream haplogroup G2 for sample Klei10 (A = ancestral, D = derived, Cov = coverage).*

Klei10								
Position	Marker	Haplogroup	Mutation	A	D	Cov	Bases	State
17174741	L156	G2	A->T	A	T	2	TT	D
14028148	L31	G2a	C->A	C	A	1	A	D
23244026	P15	G2a	C->T	C	T	1	T	D
21645555	L91	G2a2a1b	G->C	G	C	1	C	D
9985022	L293	G2a1	G->C	G	C	3	GGG	A
22741799	M286	G2a2a1a	G->A	G	A	1	G	A
14692227	L32	G2a2b	T->C	T	C	1	T	A
15604899	L30	G2a2b	C->T	C	T	1	C	A
21628300	L90	G2a2b1a	G->A	G	A	3	GGG	A
22917995	L14	G2a2b1a	C->T	C	T	2	CC	A
21645348	P303	G2a2b2a	T->C	T	C	4	TTTT	A
14871976	L78	G2a2b2a1a1	C->T	C	T	1	C	A
16903025	L1266	G2a2b2a1a2	T->G	T	C	1	T	A
14231229	L1265	G2a2b2a1a2a	A->G	A	G	1	A	A
7644368	L1264	G2a2b2a1a2a	A->G	A	G	2	AA	A
17937365	L43	G2a2b2a1b1a1	A->G	A	G	1	A	A
16660759	L42	G2a2b2a1b1a1a	C->A	C	A	3	CCC	A
6835545	Z724	G2a2b2a1c1	C->T	C	T	3	CCC	A
16596946	CTS5990	G2a2b2a1c1a	A->G	A	G	2	AA	A
18393688	L640	G2a2b2a1c1a1	A->G	A	G	3	AAA	A
17937308	L662	G2a2b2a1c2a	C->T	C	T	1	C	A
15027433	M377	G2b1	A->G	A	G	4	AAAA	A
8467136	L183	G2b1	G->C	G	C	1	G	A

## X-chromosome contamination in ancient male samples

We used ANGSD v.0.614 [111] to determine X-chromosome contamination in male samples [112]. Because male samples only have one copy of the X-chromosome, it is expected that only one allele is observed at polymorphic sites. When a second allele is observed, this implies that it has arisen due to contamination or sequencing error. ANGSD v.0.614 calculates the rate of heterozygous alleles at SNP sites and compares this value to the rate of mismatches at adjacent monomorphic sites. We started by generating a binary count file for each sample, followed by the script “contamination.R”, which performs a Fisher’s exact test and jackknife to estimate contamination. This analysis is restricted to unique regions of the X-chromosome (“RES/ChrX.unique.gz”) and to known HapMap polymorphic sites (“RES/HapMapChrX.gz”). Very low contamination values obtained for both Bar31 and Klei10 (<2.2%) support the validity of our results (Table S11).

**Table S11:** X-chromosome based contamination estimates in ancient male samples.

Sample	Method1	SE	p-value	Method 2	SE	p-value
Bar31	2.069	0.127	2.20E-16	2.133	0.002	2.20E-14
Klei10	1.255	0.133	2.20E-16	1.502	0.214	5.01E-14

**Table S12:** Number of major and minor bases detected in SNP sites and adjacent positions in the X-chromosome for Methods 1 and 2.

		Method1									
Sample	Base	-4	-3	-2	-1	SNP site	1	2	3	4	
Bar31	Minor base	853	807	836	954	1644	993	831	852	856	
	major base	97178	97335	97277	97019	96046	96924	97228	97317	97236	
	Minor base	332	345	341	355	562	363	325	309	344	
Klei10	Major base	52573	52626	52679	52632	52441	52701	52736	52773	52719	
	Method 2										
	Bar31	Minor base	256	249	256	264	474	300	251	265	256
major base		29285	29292	29285	29277	29067	29241	29290	29276	29285	
Klei10		Minor base	121	132	119	125	219	123	106	111	127
	major base	20232	20221	20234	20228	20134	20230	20247	20242	20226	
	Major base										

## Authenticity of sequencing results

All samples reported in this study were extracted independently in at least two separated reactions. We estimated the damage pattern and plotted the read length distribution to display the fragmentation pattern (see SI3 Figure S3). We found sample dependent deamination pattern varying between 30 - 55% in combination with fragment length distributions below 60bp, indicating for a prehistoric origin of our samples. Blank controls were conducted and quantified in all laboratory steps (Table S2). Contamination estimates were determined from merged datasets, including different extractions of the same sample, showing combined estimates (Table S6 & Table S7) of 0.004 - 1.619 % for shotgun data and 0.05 - 6.71 % for capture data. Additionally, X-chromosomal estimates for the male individuals Bar31 and Klei10 show low levels of contamination (<2.2%).

# SI5 Genotype Calling for Ancient DNA

Athanasios Kousathanas, Vivian Link, Christoph Leuenberger, Daniel Wegmann

## Calling Diploid Genotypes from Ancient DNA

### Basic Likelihood of Sequencing Data

Our approach is a direct extension of current approaches to genotyping for modern DNA and our model follows most closely the basic model introduced by Li [113].

The observed data  $d_i$  at site  $i$  corresponds to what is typically obtained when individual reads  $j$  of next generation sequencing were aligned by mapping against a reference genome. Here we will assume that all sequencing reads were accurately mapped and hence that reads with low mapping qualities have been filtered out. The data  $d_i$  obtained at site  $i$  thus consists of a list of  $n_i$  observed bases  $d_i = \{d_{i,1}, \dots, d_{i,n_i}\}$ ,  $d_{i,j} = A, C, G, T$ .

Let us denote the hidden genotype at site  $i$  by  $g_i$  where  $g_i$  consists of a pair of nucleotides  $rs$  with  $r, s = A, G, C, T$ . The basic likelihood function for genotype calling relates the observed bases at site  $i$  to the hidden genotype  $g_i =$  at this location, while also accounting for sequencing errors. Let us denote by  $\epsilon_{i,j}$  the probability of a sequencing error at the base of read  $j$  covering site  $i$ . The likelihood of the full data is thus given by

$$\mathbb{P}(d_i | g_i = rs, \epsilon_i) = \prod_{k=1}^{n_i} \mathbb{P}(d_{i,j} | g_i = rs, \epsilon_{i,j}),$$

where  $\epsilon_i = \{\epsilon_{i,1}, \dots, \epsilon_{i,n_i}\}$ .

Under the assumption that a sequencing read is equally likely to cover any of the two alleles of an individual, the probability of observing a base  $d_{i,j} = T$  given the underlying genotype  $g_i = CT$  is then given by the probability of sequencing allele  $T$  without making a sequencing error with probability  $\frac{1}{2}(1 - \epsilon_{i,j})$  or sequencing allele  $C$  and making a sequencing error with probability  $\frac{1}{2}\frac{\epsilon}{3}$ . We thus have

$$\mathbb{P}(d_{i,j} | g_i = kl, \epsilon_{i,j}) = \begin{cases} 1 - \epsilon_{i,j} & \text{if } r = s = d_{i,j} \\ \frac{\epsilon}{3} & \text{if } r \neq d_{i,j}, s \neq d_{i,j} \\ \frac{1}{2} - \frac{\epsilon_{i,j}}{3} & \text{if } r \neq s, r = d_{i,j} \text{ or } s = d_{i,j} \end{cases}$$

Here we wish to extend this basic model to incorporate particular features of ancient DNA. For that, let us denote for each individual observed base  $d_{i,j}$  a vector  $q_{i,j}$  of external information such

as the quality score reported by the sequencing machine, the position within the read, the distance from the 5' and 3' ends or the nucleotide context (i.e. the previous base read).

## Post-Mortem Damage

A characteristic feature of ancient DNA is post-mortem damage (PMD). The most common form of post-mortem damage is  $C$  deamination, which leads to a  $C \rightarrow T$  transition on the affected strand and a  $G \rightarrow A$  transition on the complementary strand [e.g. 114]. These deaminations do not occur randomly along the whole read, but instead are observed much more frequently at the beginning of a read. This is due to fragment ends being more often single-stranded and thus subject to a much higher rate of deamination.

Here we follow [115] in assuming that the probability of observing such a deamination  $D_m, m = C \rightarrow T, G \rightarrow A$  decays exponentially with distance  $p$  from the end of the read. We chose to model this as

$$D_m(q_{i,j}) = A_m \cdot \exp(-p_{i,j}^{(m)} \cdot B_m) + C_m \quad (1)$$

where  $p_{i,j}^{(m)} = p_{i,j}^{(5')}$  and  $p_{i,j}^{(m)} = p_{i,j}^{(3')}$  are the distances of the observed base at site  $i$  from the 5' or 3' end of read  $j$  for  $m = C \rightarrow T$  and for  $m = G \rightarrow A$ , respectively. Further,  $A_m, B_m$  and  $C_m$  are assumed to be known constants (i.e. learned from the data *a priori*).

We now seek to develop a model for the emission probabilities  $\mathbb{P}(d_{i,j}|g_i, \epsilon_{i,j}, q_{i,j})$  that takes both sequencing errors and post-mortem damage patterns into account. Under the assumption that sequencing errors and post-mortem damage are occurring independently among reads, the emission probability is given by

$$\mathbb{P}(d_i|g_i, \epsilon_i, q_i) = \prod_{j=1}^{n_i} \mathbb{P}(d_{i,j}|g_i, \epsilon_{i,j}, q_{i,j}), \quad (2)$$

where  $q_i = \{q_{i,j}, \dots, q_{i,n_i}\}$ .

While we report the emission probabilities  $\mathbb{P}(d_{i,j}|g_i, \epsilon_{i,j}, q_{i,j})$  for all combinations of  $d_{i,j} = A, C, G, T$  and  $g_i = kl$  with  $k, l = A, C, G, T$ , we illustrate here our reasoning to derive those probabilities for the case of observing  $d_{i,j} = T$  given the underlying genotype  $g_i = CT$ . There are three possible ways to obtain a  $T$ : either by sequencing allele  $T$  without error, sequencing allele  $C$  that was deaminated (post-mortem damage) without sequencing error, or sequencing allele  $C$  that was not deaminated with error. We thus have



$$\begin{aligned}\mathbb{P}(d_{i,j} = T | g_i = CT, \epsilon_{i,j}, q_{i,j}) &= \frac{1}{2} \left( (1 - \epsilon_{i,j}) + D(q_{i,j})(1 - \epsilon_{i,j}) + (1 - D(q_{i,j})) \frac{\epsilon_{i,j}}{3} \right) \\ &= \frac{1}{2} \left( (1 - \epsilon_{i,j})(1 + D(q_{i,j})) + (1 - D(q_{i,j})) \frac{\epsilon_{i,j}}{3} \right).\end{aligned}$$

Here,  $D(q_{i,j}) = D_{C \rightarrow T}$  is given by eq. 1 based on the distance from the 5' end of the read. Also note that under the assumption of random sequencing errors, only one out of three will result in the actual base being observed.

## Genotyping

The basic caller we implement here is a maximum likelihood (ML) caller based on the likelihood function described above. Specifically, and independently for each site, we calculate the likelihood  $\mathbb{P}(d_i | g_i, \epsilon_i, q_i)$  according to eq. 2 for each genotype. The genotype with the highest likelihood will then be called for this position. Following GATK [116], we calculate the Phred-scaled quality of the call as the likelihood difference  $\Delta LL$  between the ML and the second most likely genotype.

## Recalibration

While all genotyping algorithms rely on the sequencing error rates, these rates are usually not known. The goal is thus to estimate sequencing error rates by leveraging external information provided for each observation by the sequencing machine. Typically, the quality reported by the sequencing machine, the position in the read, and the base context are considered [97]. The goal is then to fit a statistical model that accurately predicts sequencing error rates from these covariates.

Here we adopt a similar strategy as implemented in the Base Quality Score Recalibration (BQSR) tool available in GATK [116], but extend it to ancient DNA by incorporating post-mortem damage. We assume that the error rates of a specific observation  $\epsilon_{i,j}$  can be decomposed as

$$\epsilon_{i,j} = \epsilon_q(q_{i,j}) \prod_{k=1}^K \alpha_k(q_{i,j}),$$

where  $\epsilon_q(q_{i,j})$  is the mean error rate for all observations with a given quality score at observation  $j$  at site  $i$  and all  $\alpha_k$  are scaling factors for additional covariates. Here, the covariates we consider are position in the read and di-nucleotide context.

To learn all  $\epsilon_q$  and  $\alpha_k$  from the data, we follow the implementation in GATK as described in [97] and will focus on a subset of the data for which it can be assumed that all differences from the reference genome are sequencing errors. Following standard recommendations, we achieve this by

focusing only on autosomes and the X chromosome and masking i) all sites known to be polymorphic in humans as predicted by 1000G, HapMap and dbSNP and ii) repetitive, telomeric and centromeric regions retrieved from the UCSC Table Browser [117] by using track *Repeatmasker* in group *Repeats* and track *Gap* in group *Mapping and Sequencing*.

### Learning all $\epsilon_q$

Under the assumption that the genotypes at all remaining sites are known to be homozygous reference, the relevant log-likelihood function to learn the mean error rate for all sites with quality  $q_{i,j} = q$  is then given by

$$\log \mathbb{P}(\mathbf{d}|\mathbf{g}, \epsilon_q) = \sum_{i=1}^I \sum_{j=1}^{n_i} \mathbb{I}(q_{i,j} = q) \log \mathbb{P}(d_{i,j}|g_i, \epsilon_q, q_{i,j}),$$

where  $\mathbb{I}(q_{i,j} = q)$  is an indicator function that equals one if the quality score reported by the sequencing machine for the observation in read  $j$  at site  $i$  is equal to  $q$  and zero otherwise. The emission probabilities  $\mathbb{P}(d_{i,j}|g_i, \epsilon_q, q_{i,j})$  are calculated as indicated in Table S13 for the genotypes  $g_i = rr, r = A, C, G, T$  (top four rows). To generalize the notation, we will denote these probabilities by

$$\mathbb{P}(d_{i,j}|g_i, \epsilon_q, q_{i,j}) = (1 - D_{i,j,g_i}) \frac{\epsilon_q}{3} + D_{i,j,g_i} (1 - \epsilon_q),$$

where

$$D_{i,j,g_i} = \begin{cases} 1 - D_{C \rightarrow T}(q_{i,j}) & \text{if } g_i = CC, d_{i,j} = C \\ D_{C \rightarrow T}(q_{i,j}) & \text{if } g_i = CC, d_{i,j} = T \\ D_{G \rightarrow A}(q_{i,j}) & \text{if } g_i = GG, d_{i,j} = A \\ 1 - D_{G \rightarrow A}(q_{i,j}) & \text{if } g_i = GG, d_{i,j} = G \\ 1 & \text{if } g_i = AA, d_{i,j} = A \text{ or } g_i = TT, d_{i,j} = T \\ 0 & \text{otherwise.} \end{cases}$$

Since this function can not be maximized for  $\epsilon_q$  analytically, we adopt a Newton-Ralphson scheme to find the maximum likelihood estimate of  $\epsilon_q$  starting at the value indicated by the quality score  $q$  provided by the machine. The relevant derivatives are given by

$$\frac{\partial}{\partial \epsilon_q} \mathbb{P}(\mathbf{d}|\mathbf{g}, \epsilon_q) = \sum_{i=1}^I \sum_{j=1}^{n_i} \mathbb{I}(q_{i,j} = q) \frac{1 - 4D_{i,j,g_i}}{D_{i,j,g_i}(3 - 4\epsilon_q) + \epsilon_q}$$

and

$$\frac{\partial^2}{\partial \epsilon_q^2} \mathbb{P}(\mathbf{d}|\mathbf{g}, \epsilon_q) = - \sum_{i=1}^I \sum_{j=1}^{n_i} \mathbb{I}(q_{i,j} = q) \left( \frac{(1 - 4D_{i,j,g_i})}{(D_{i,j,g_i}(3 - 4\epsilon_q) + \epsilon_q)} \right)^2.$$

### Learning scaling factors $\alpha_k$ of covariates

To learn the scaling factors of covariates sequentially, we adopt a very similar strategy, but always conditioning on the currently best estimate of  $\epsilon_{i,j}$  given  $\epsilon_q$  and all previously estimated scaling factors. The log likelihood function to learn the scaling factor of covariate  $l$  is then given by

$$\log \mathbb{P}(\mathbf{d}|\mathbf{g}, \epsilon_q) = \sum_{i=1}^I \sum_{j=1}^{n_i} \mathbb{I}(q_{i,j}, a) \log \mathbb{P}(d_{i,j} | g_i, \hat{\epsilon}_{i,j}, \alpha_{l,a}, q_{i,j}),$$

where  $a$  is a particular value of this covariate (e.g. a particular position in the read or a particular context),  $\mathbb{I}(q_{i,j}, a)$  is an indicator function that is one if the covariate for the observation in read  $j$  at site  $i$  is equal to the  $a$  and zero otherwise, and

$$\hat{\epsilon}_{i,j} = \epsilon_q(q_{i,j}) \prod_{k=1}^{l-1} \alpha_k(q_{i,j}).$$

Again, we find the maximum likelihood scheme using a Newton-Raphson scheme starting at  $\alpha_{l,a} = 1$ . The relevant derivations are

$$\frac{\partial}{\partial \alpha_{l,a}} \mathbb{P}(\mathbf{d}|\mathbf{g}, \epsilon_q) = \sum_{i=1}^I \sum_{j=1}^{n_i} \mathbb{I}(q_{i,j}, a) \frac{\hat{\epsilon}_{i,j}(1 - 4D_{i,j,g_i})}{D_{i,j,g_i}(3 - 4\hat{\epsilon}_{i,j}\alpha_{l,a}) + \hat{\epsilon}_{i,j}\alpha_{l,a}}$$

and

$$\frac{\partial^2}{\partial \alpha_{l,a}^2} \mathbb{P}(\mathbf{d}|\mathbf{g}, \epsilon_q) = - \sum_{i=1}^I \sum_{j=1}^{n_i} \mathbb{I}(q_{i,j}, a) \left( \frac{\hat{\epsilon}_{i,j}(1 - 4D_{i,j,g_i})}{(D_{i,j,g_i}(3 - 4\hat{\epsilon}_{i,j}\alpha_{l,a}) + \hat{\epsilon}_{i,j}\alpha_{l,a})} \right)^2.$$

### Calling Allele Presence

Since coverage of many ancient samples does not allow for accurate diploid genotype calling, we develop here a strategy to make a haploid calls at each position that takes observation-specific error rates  $\epsilon_{i,j}$  as well as post-mortem damage into account. We do this by calculating at each position the probability for each possible base to be present. Specifically, we are interested in calculating

$$\mathbb{P}(g_i = r \cdot | d_i, \epsilon_i, q_i) = \mathbb{P}(g_i = rr | d_i, \epsilon_i, q_i) + 2 \sum_{s \neq r} \mathbb{P}(g_i = rs | d_i, \epsilon_i, q_i),$$

where  $g_i = r \cdot$  denotes any genotype that contains at least one allele  $r$  (e.g.  $AA$ ,  $AC$ ,  $AG$ ,  $AT$ ,  $CA$ ,  $GA$  or  $TA$  for  $r = A$ ). Note that, due to symmetry,  $\mathbb{P}(g_i = rs|d_i, \epsilon_i, q_i) = \mathbb{P}(g_i = sr|d_i, \epsilon_i, q_i)$ .

To calculate this probability, we adopt a Bayesian approach and assume that the prior probability of each genotype  $g_i$  is given by the base frequencies  $\boldsymbol{\pi} = \{\pi_A, \pi_C, \pi_G, \pi_T\}$  in the region and the heterozygosity in the genome. We employ the classic substitution model by Felsenstein [118], which expresses the probability of each genotype  $g_i$  as a function of  $\boldsymbol{\pi}$  and the substitution rate  $\theta = 2T\mu$  along the genealogy connecting the two alleles of an individual, where  $T$  corresponds to the time to the most recent common ancestor of the two lineages and  $\mu$  the mutation rate. Under this model, the probability of observing a specific genotype  $g_i = rs$  is given by

$$\mathbb{P}(g_i = rs|\theta, \boldsymbol{\pi}) = \begin{cases} \pi_r q_{rr}(2T) = \pi_r(e^{-\theta} + \pi_r(1 - e^{-\theta})) & \text{if } r = s, \\ \pi_r q_{rs}(2T) = \pi_r \pi_s(1 - e^{-\theta}) & \text{if } r \neq s, \end{cases}$$

where  $q_{rs}(2T)$  denote the substitution rate as a function of time  $2T$ .

This then allows us to calculate all

$$\mathbb{P}(g_i = rs|d_i, \epsilon_i, q_i, \boldsymbol{\pi}, \theta) \propto \mathbb{P}(d_i|g_i = rs, \epsilon_i, q_i)\mathbb{P}(g_i|\boldsymbol{\pi}, \theta),$$

and hence the posterior probabilities  $\mathbb{P}(g_i = r \cdot |d_i, \epsilon_i, q_i, \boldsymbol{\pi}, \theta)$  for each base  $r = A, C, G, T$  at each position. We can then accurately call the most likely base present even for low coverage genomes, and report quality scores as the Phred-scaled probability  $1 - \mathbb{P}(g_i = r \cdot |d_i, \epsilon_i, q_i, \boldsymbol{\pi}, \theta)$  of the most likely allele  $r$  to be present.

To use this approach in practice, we assume  $\theta = 10^{-3}$  in the genome and that the base frequencies  $\boldsymbol{\pi}$  can be learned from the observed base frequencies within the  $1Mb$  window containing site  $i$ .

## Benefit of taking Post Mortem Damage into account

To test if the approaches introduced here lead to a higher accuracy in genotype calling, we simulated 10 Mb of data for a diploid individual matching the common characteristics of ancient samples as follows:

1. We simulated a reference genome with equal base frequencies (25% each).
2. The genotypes of the individual were then simulated from the reference by adding mutations. Since heterozygous sites are more difficult to call, we simulated a heterozygosity of 1% (instead of  $10^{-3}$ , as is observed for humans).
3. We then simulated reads of 100bp with starting positions chosen randomly within the sequence up to a coverage of 20X.

4. Following [115], we added post mortem damage by simulating  $C \rightarrow T$  and  $G \rightarrow A$  transitions with probability

$$D = (1 - \lambda)^{p_i - 1} p_i + C,$$

where  $\lambda = 0.3$ ,  $C = 0.01$  and  $p_i$  is the relevant position of base  $i$  within the read (from the 3' and 5' end for  $C \rightarrow T$  and  $G \rightarrow A$  transitions, respectively).

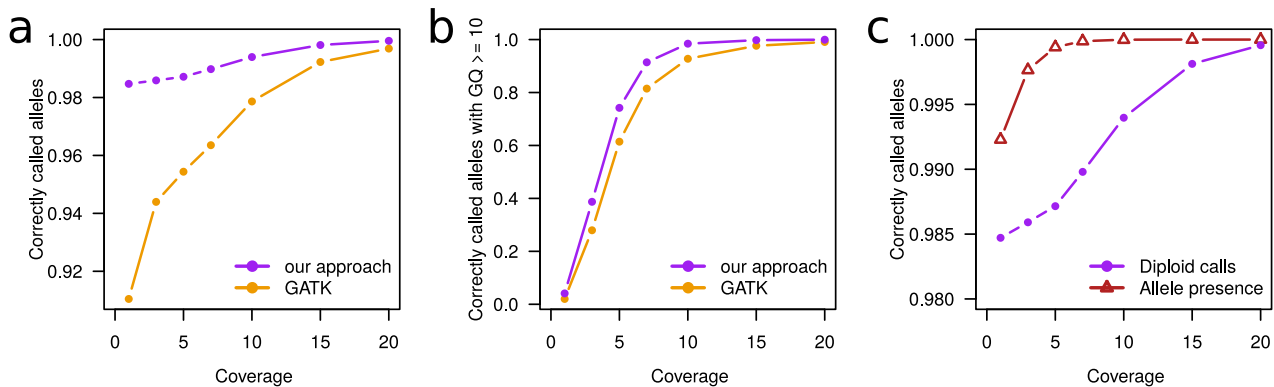
5. For each observed base we simulated fake quality scores  $q_i$  that were then transformed into true error rates  $\epsilon_i$  to test our BQSR approach. We started by simulating fake qualities from a normal distribution with mean 25 and sd 7.5. To simulate a quality effect for recalibration, we then transformed those qualities as  $q'_i = 10 \log(q)$ . To simulate a position effect, we first transformed the phred-scaled  $q'$  into error rates  $\epsilon'$ . These error rates were then scaled as  $\epsilon_i = \epsilon'_i * 0.5 + \frac{1.5}{99}(p_i - 1)$ , where  $p_i$  is the position of base  $i$  within the read. Finally, we simulated a context effect by simulating errors with unequal base probabilities. We chose to have errors resulting in a  $T$  in 70% of all cases, but in a  $A, C$  or  $G$  in only 5%, 15% and 10%, respectively.
6. The fake qualities together with the simulated bases were written to a BAM file, while the simulated reference was written to a FASTA file.

The data simulated in this way was then downsampled to obtain a range of coverages and then pushed through both GATK as well as our pipeline, consisting of BQSR and genotype calling in both cases. For the BQSR step we masked all truly polymorphic sites for both approaches. In addition, we provided the true PMD probability distributions to our pipeline. For variant calling with GATK, we used the recommended Haplotype Caller and forced GATK to print all sites using the options `-ERC BP_RESOLUTION` and `-mmq 0` but otherwise default settings.

We compared the accuracy of the genotype calls between the two approaches by calculating the fraction of correctly called alleles. For a single site with true genotype  $AA$  and a genotype call  $AA$ , both allele were considered correctly called, but for a call of  $AC$ ,  $CC$  or  $CT$  only 1, 0 and 0 alleles were considered correctly called. For sites without a call we counted a distance of 2.

We note that the Haplotype Caller in GATK has, by design, a bias towards the reference allele. At truly heterozygous sites, for instance, GATK calls more frequently homozygous reference than homozygous alternative genotypes. Since we simulated only 1% of all sites to be heterozygous, this bias is in favor of GATK over our caller that does not take information about the reference into account.

Despite this, we found our caller to outperform GATK for all coverages tested (Figure S4a), but with the largest benefit for relatively low coverages, where the amount of errors is effectively halved. This increased accuracy was even more pronounced when filtering for a genotype quality of at least 10 (Figure S4b), which is more frequently obtained using our pipeline, most likely because the BQSR in GATK overestimates general error rates due to the presence of PMD.



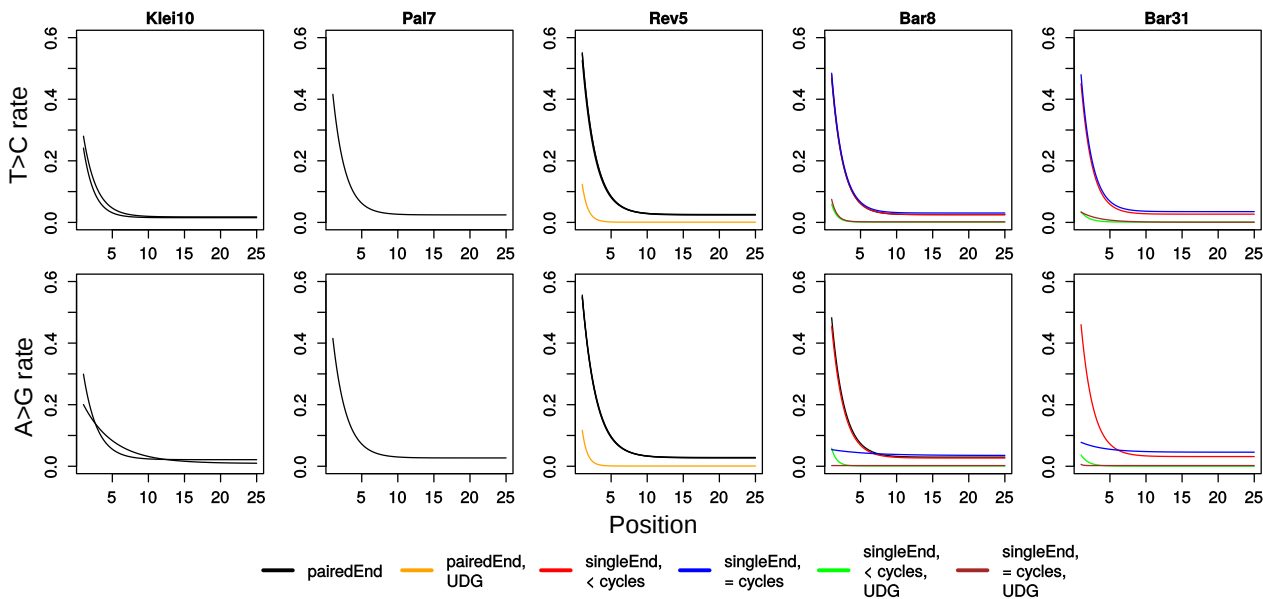
**Figure S4:** Accuracy of genotype calling approaches. We assessed the accuracy of our genotype calling approach using 10Mb of data simulated to match the characteristics of ancient DNA data (see text). This data was then called using our approach, as well as GATK for comparison. In both cases, we used the implemented BQSR approach to recalibrate quality scores. We then assessed the fraction of correctly called alleles for all sites with a coverage of at least one at different coverages. **a)** Shown is the fraction of correctly called alleles for GATK and the approach developed here. **b)** Same as a) but for correctly called alleles with a genotype quality  $\geq 10$ . This is almost impossible to obtain at very low coverages. **c)** Comparison between our approaches to call diploid individuals (same as in a)) and allele presence.

Taken together, these results suggest that modeling PMD is critical for genotype calling in ancient DNA, particularly when coverage is low. However, they also illustrate the general difficulty of genotype calling at low coverages. Indeed, the majority of the wrongly inferred alleles are at heterozygous sites, even if we simulated on 1% of all sites to be heterozygous.

One way to diminish these errors is to abstain from diploid calls at low coverage. Here we present a natural way to infer the most likely allele present at each site by means of Bayesian approach. When applying this approach to the simulated data, we indeed observe much higher accuracy at low coverages than when calling diploid genotypes (Figure S4c). However, this obviously comes at the cost of only calling half of the data.

## Genotyping of Aegean samples

We called both MLE genotypes as well as Bayesian allele presence for the three Greek (Klei10, Pal7 and Rev5) and two Anatolian (Bar8 and Bar31) samples. To do so, we first split the full data into different read-groups based on the extraction protocol used and the sequencing run. Specifically, we made sure to treat protocols with and without the application of uracil-DNA glycosylase (UDG) differently, as well as those resulting in single-end and paired-end reads. While we merged all paired-end reads prior to the analysis, single-end reads were split by their length into one group containing all reads that were shorter than the number of sequencing cycles, and another group containing all reads as long as the number of sequencing cycles. This was done to account for the different  $G \rightarrow A$



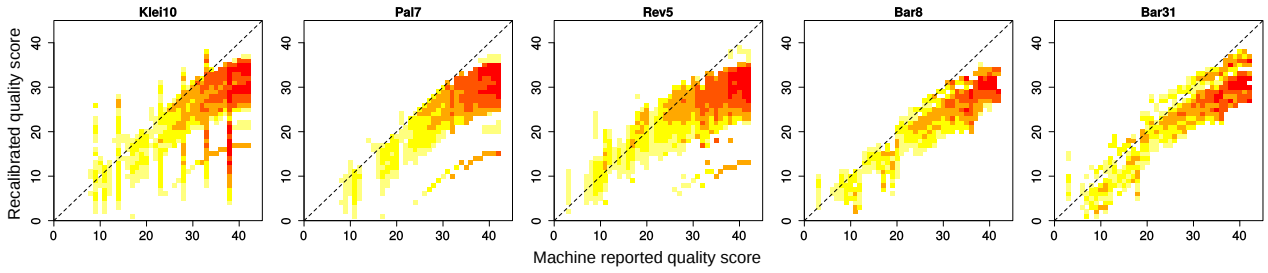
**Figure S5:** Post-mortem damage pattern profiles of Aegean samples. Plotted is the frequency of reads showing  $T$  at all sites where the reference is  $C$  (top panels) and the frequency of reads showing  $A$  at all sites where the reference is  $G$  (bottom panels) as a function of the distance from the 5' and 3' end of the read, respectively, for each individual and read group. The patterns are colored according to the sequencing scheme (paired-end or single-end) and whether or not they were treated with uracil-DNA glycosylase (UDG) before sequencing to remove damaged bases. For single end sequencing runs, damage patterns are further shown individually for reads that are shorter ( $<$  cycles) or as long ( $=$  cycles) as the number of sequencing cycles used.

damage patterns of those groups that result from the fact that all reads from long fragments are expected to show little damage at their 3' end.

Next, we inferred damage patterns for each read group individually using `MapDamage` [86] and fitted the parameters of eq. 1 using ordinary least squares (Figure S5). In order to account for the expected variation between the reference and our samples, we fitted the probability of  $C \rightarrow T$  damage patterns to the frequency of  $C \rightarrow T$  changes at a specific position in the read, minus the frequency of  $T \rightarrow C$  at the same position. Reassuringly, observed damage patterns closely matched those expected given the protocols used.

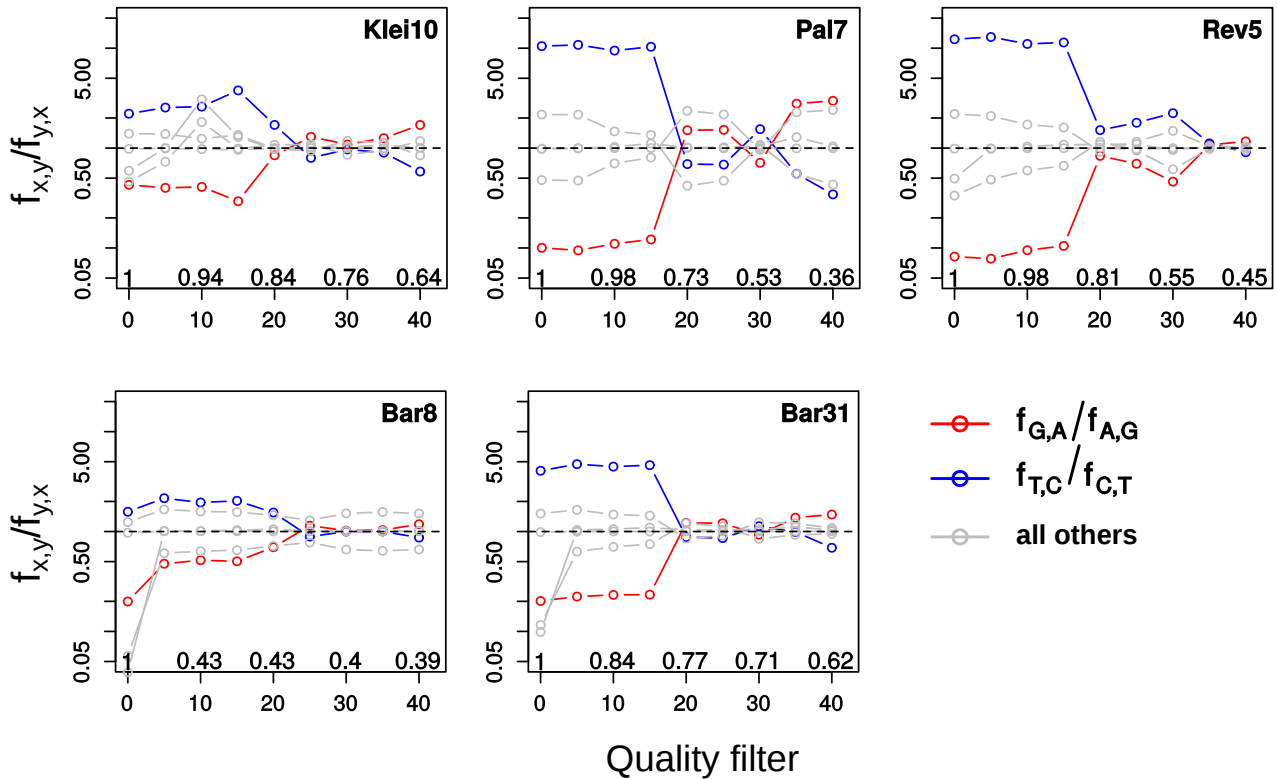
We then performed BQSR recalibration as described above (subsection Recalibration) by considering the original quality scores as well as the position in read and di-nucleotide context. As shown in Figure S6, this recalibration step was crucial, as it resulted in substantially lower qualities than reported by the sequencing machines, despite fully accounting for post-mortem damage.

The allele calls were then obtained using the recalibrated quality scores and by taking post-mortem damage into account as described above. To check if our allele presence calls are properly accounting for post-mortem damage, we determined the proportion  $f_{x,y}$  of sites at which we called allele  $y = A, C, G, T$  while the reference was  $x = A, C, G, T$ . Given the symmetry of the mutational process,



**Figure S6:** Effect of Recalibration on Quality Scores. Shown are the density distributions of the quality transformations as a result of the applied quality recalibration for each sample. It appears that the machine-reported qualities were overall too high, and in particular for the qualities  $> 30$ .

we expect the ratio of  $f_{x,y}/f_{y,x}$  to approach 1. In Figure S7 we plot these ratios for all pairs  $x \neq y$  as a function of applying increasingly strict quality filters. We found that while there was an excess of  $C \rightarrow T$  and  $G \rightarrow A$  calls in the raw data, allele presence calls with a quality  $\leq 20$  are very robust to post-mortem damage.



**Figure S7:** Effect of filtering on allele presence calls. To quantify the effect of post-mortem damage on our allele presence calls, we determined the proportion  $f_{x,y}$  of sites at which we called allele  $y = A, C, G, T$  while the reference was  $x = A, C, G, T$ . Plotted here is the ratio of  $f_{x,y}/f_{y,x}$  for all pairs  $x \neq y$  as a function of filtering, with pairs affected by post-mortem damage in color ( $G/A$  in red and  $T/C$  in blue). For calls associated with qualities  $> 20$ , these patterns become indistinguishable from those of pairs of bases not affected by post-mortem damage. The numbers at the bottom of each plot indicate the fraction of the genome passing the respective filters.



**Table S13: Emission probabilities**

The probability of observing a specific base depending on the underlying genotype (first column), the observation-specific error probability and the probability of post-mortem damage. For the sake of brevity, we denote the observation-specific error probability by  $\epsilon$  instead of the full term  $\epsilon_{i,j}$  and use  $D_{C \rightarrow T} = D_{C \rightarrow T}(q_{i,j})$  and  $D_{G \rightarrow A} = D_{G \rightarrow A}(q_{i,j})$ .

	<b>A</b>	<b>C</b>	<b>G</b>	<b>T</b>
<b>AA</b>	$(1 - \epsilon)$	$\frac{\epsilon}{3}$	$\frac{\epsilon}{3}$	$\frac{\epsilon}{3}$
<b>CC</b>	$\frac{\epsilon}{3}$	$(1 - D_{C \rightarrow T})(1 - \epsilon) + D_{C \rightarrow T}\frac{\epsilon}{3}$	$\frac{\epsilon}{3}$	$(1 - D_{C \rightarrow T})\frac{\epsilon}{3} + D_{C \rightarrow T}(1 - \epsilon)$
<b>GG</b>	$(1 - D_{G \rightarrow A})\frac{\epsilon}{3} + D_{G \rightarrow A}(1 - \epsilon)$	$\frac{\epsilon}{3}$	$(1 - D_{G \rightarrow A})(1 - \epsilon) + D_{G \rightarrow A}\frac{\epsilon}{3}$	$\frac{\epsilon}{3}$
<b>TT</b>	$\frac{\epsilon}{3}$	$\frac{\epsilon}{3}$	$\frac{\epsilon}{3}$	$(1 - \epsilon)$
<b>AC</b>	$\frac{1}{2} - \frac{\epsilon}{3}$	$\frac{(1 - D_{C \rightarrow T})(1 - \epsilon)}{2} + \frac{(1 + D_{C \rightarrow T})\epsilon}{6}$	$\frac{\epsilon}{3}$	$\frac{D_{C \rightarrow T}(1 - \epsilon)}{2} + \frac{(2 - D_{C \rightarrow T})\epsilon}{6}$
<b>AG</b>	$\frac{(1 + D_{G \rightarrow A})(1 - \epsilon)}{2} + \frac{(1 - D_{G \rightarrow A})\epsilon}{6}$	$\frac{\epsilon}{3}$	$\frac{(1 - D_{G \rightarrow A})(1 - \epsilon)}{2} + \frac{(1 + D_{G \rightarrow A})\epsilon}{6}$	$\frac{\epsilon}{3}$
<b>AT</b>	$\frac{1}{2} - \frac{\epsilon}{3}$	$\frac{\epsilon}{3}$	$\frac{\epsilon}{3}$	$\frac{1}{2} - \frac{\epsilon}{3}$
<b>CG</b>	$\frac{D_{G \rightarrow A}(1 - \epsilon)}{2} + \frac{(2 - D_{G \rightarrow A})\epsilon}{6}$	$\frac{(1 - D_{C \rightarrow T})(1 - \epsilon)}{2} + \frac{(1 + D_{C \rightarrow T})\epsilon}{6}$	$\frac{(1 - D_{G \rightarrow A})(1 - \epsilon)}{2} + \frac{(1 + D_{G \rightarrow A})\epsilon}{6}$	$\frac{D_{C \rightarrow T}(1 - \epsilon)}{2} + \frac{(2 - D_{C \rightarrow T})\epsilon}{6}$
<b>CT</b>	$\frac{\epsilon}{3}$	$\frac{(1 - D_{C \rightarrow T})(1 - \epsilon)}{2} + \frac{(1 + D_{C \rightarrow T})\epsilon}{6}$	$\frac{\epsilon}{3}$	$\frac{(1 + D_{C \rightarrow T})(1 - \epsilon)}{2} + \frac{(1 - D_{C \rightarrow T})\epsilon}{6}$
<b>GT</b>	$\frac{D_{G \rightarrow A}(1 - \epsilon)}{2} + \frac{(2 - D_{G \rightarrow A})\epsilon}{6}$	$\frac{\epsilon}{3}$	$\frac{(1 - D_{G \rightarrow A})(1 - \epsilon)}{2} + \frac{(1 + D_{G \rightarrow A})\epsilon}{6}$	$\frac{1}{2} - \frac{\epsilon}{3}$

## SI6 PCA

Yoan Diekmann, Mark Thomas

### Methods

The PCA plots shown in Figure 2 in the main text are generated using LASER version 2.02 [119] (see 3D-figure S4 at [https://figshare.com/articles/Hofmanova\\_et\\_al.3D\\_figure\\_S4/3188767](https://figshare.com/articles/Hofmanova_et_al.3D_figure_S4/3188767) for an interactive version in three dimensions). The coordinates are computed in two steps, for modern and ancient samples separately (datasets used here are described in the sections below).

First, a reference space is generated by standard PCA on genotype data of modern individuals. Missing entries are imputed by averaging the genotypes encoded as  $\{0, 1, 2\}$  over all individuals.

Second, ancient samples are mapped into the reference space. In this step, the input for LASER is sequence reads in form of BAM files. Each sample is placed by first simulating sequencing data for each reference individual that matches the coverage pattern of the ancient sample. Then, a PCA is generated for the simulated reference data together with the ancient individual. Finally, Procrustes analysis of the latter PCA with respect to the reference from step 1 allows to project the sample into the reference space [120].

The advantages of this approach are twofold and especially relevant in the context of ancient DNA. First, the use of BAM files for the ancient samples circumvents the need for genotype calling, which can be challenging for ancient DNA due to contamination, post-mortem damage and low coverage. However, the modern reference is still based on genotype data, which is abundant and of high quality for present-day populations. Second, Procrustes analysis allows for robust placement of samples despite low coverage, which is particularly valuable in the case of ancient DNA data.

### Modern reference data

We use a genotyping dataset with samples from present-day European and Middle Eastern populations (see legend of Figure 2 in the main text) from a merged dataset published as part of Hellenthal *et al.* [121] and Busby *et al.* [122] based on Illumina’s genotyping platform (Illumina 550, 610, 660W arrays). We remove duplicate individuals (identified by  $\hat{\pi} = 1$  running PLINK version 1.9 [123, 124] with option ‘--genome’) and a visual PCA outlier (sample Jordan444). After lifting the coordinates to human genome version hg19 (using PyLiftover [125, 126]) we obtain a total of 1051 individuals and 510,811 autosomal polymorphic loci.

## Ancient data

All ancient samples shown in Figure 2 in the main text are projected into the reference PCA space by LASER as explained above. Please see Supplementary Table S14 for details on the ancient samples, their IDs and the corresponding references.

The accuracy of the projection depends on the number of covered SNP positions and their read depth. LASER provides a sample-specific Procrustes similarity score  $t$  ranging from 0 and 1 that quantifies the confidence in the projection, with values closer to 0 indicating higher uncertainty (see Wang *et al.* [120] for details on the statistic). Table S16 lists number of sites covered by at least one read, the average read depth over all reference sites and the  $t$ -statistic for all samples projected here.

**Table S14:** List with all genomes discussed in the text giving relevant information in chronological order.

Site	Sample IDs	Coverage	mtDNA HG	Country	abbr.	Culture (arch. context)	abbr.	Time cal BCE	Label PCA	Reference
<b>45000 - 10000 cal BCE</b>										
Ust-Ishim	Ust-Ishim	42	R*	Siberia	SIB	Upper Paleolithic	UP	44930 - 41260	Ust-Ishim-SIB-UP	[127]
Kostenki	K14	2.42 (-UDG), 2.84 (+UDG)	U2	Russia	RU	Upper Paleolithic	UP	36734 - 34312	Kostenki-RU-UP	[128]
Bichon	Bichon	9.5	U5b1h	Switzerland	CH	Upper Paleolithic	UP	11820 - 11610	Bichon-CH-UP	[129]
Satsurblia	Satsurblia	1.44	K3	Georgia	GE	Upper Paleolithic	UP	11430 - 11882	Satsurblia-GE-UP	[129]
Kotias	Kotias	15.38	H13c	Georgia	GE	Mesolithic	M	7945 - 7579	Kotias-GE-M	[129]
<b>7000 - 6000 cal BCE</b>										
Barcın	see Table S15			Turkey	TR	Neolithic	N	6500 - 6200	Barcın-TR-N	[130]
Mentese	I0723, I0724, I0726, I0727	see Table S15		Turkey	TR	Neolithic	N	6400 - 5600	Mentese-TR-N	[130]
Barcın	Bar8, Bar31	7.21, 3.71	K1a2, X2m	Turkey	TR	Early Neolithic	N	6419 - 6030	Bar8-TR-N, Bar31-TR-N	This study
Revenia	Rev5	1.17	X2b	Greece	GR	Early Neolithic	N	6438 - 6264	Rev5-GR-N	This study
Loschbour	Loschbour	22	U5b1a	Luxembourg	LU	Mesolithic	M	6220 - 5990	Loschbur-LU-M	[41]
Motala	Motala1, 2, 3, 4, 6, 9, 12	0.18, 0.15, 0.55, 0.07, 0.02, 0.01, 2.4	U5a1, U2e1, U5a1, U5a2d, U5a2d, U5a2, U2e1	Sweden	SW	Mesolithic	M	6361 - 5516	Motala-SW-M	[41]
<b>6000 - 5000 cal BCE</b>										
La Braña	La Braña1	3.4	U5b2c1	Spain	ES	Mesolithic	M	5990 - 5740	La Braña-ES-M	[131]
Tiszaszölös-Domaháza	KO1	1.24	R3	Hungary	HUN	Körös/Mesolithic	KÖR HG	5780-5650	Tisz.-Doma.-HUN-KÖR-HG	[81]
Berettyóújfalú-Morotva-liget	KO2	0.13	K1	Hungary	HUN	Körös	N	5710 - 5570	Ber.-Moro.-HUN-N	[81]
Stora Förvar	StoraFörvar11	0.09	U5a1	Sweden	SW	Mesolithic	M	5550 - 5300	Stora Förvar-SW-M	[132]
Cova Bonica	CB13	1.10	K1a2a	Spain	ES	Cardial	N	5470 - 5360	Els-Trocs-ES-N	[131]
Els Trocs	Troc1,3,5,7	0.68, 1.30, 13.78, 1.57 <sup>1</sup>	J1c3, T2c1d, N1a1a1, V	Spain	ES	Epicardial	N	5310 - 5066	Els-Trocs-ES-N	[130]
Polgár-Ferenci-hát	NE1, NE4	22.12, 0.10	U5b2c, J1c	Hungary	HUN	Alföld LBK (ALP)	N	5310 - 5070	Polg.-Fer.-HUN-N	[81]
Debrecen Tócsópart Erdoalja	NE2	0.19	H	Hungary	HUN	Alföld LBK (ALP)	N	5290 - 5060	Debr.-Tócsó.-HUN-N	[81]
Garadna	NE3	0.13	X2b	Hungary	HUN	Alföld LBK (ALP)	N	5210 - 5010	Garadna-HUN-N	[81]
Kompolt-Kigyósér	NE5	1.04	J1c1	Hungary	HUN	Alföld LBK (ALP)	N	5210 - 4990	Komp.-Kig.-HUN-N	[81]
Apc-Berekalja I.	NE6	1.18	K1a3a3	Hungary	HUN	LBK	N	5300 - 4950	Apc-Berekalja-HUN-N	[81]
Viesenhäuser Hof, Stuttgart-Mühlhausen	Stuttgart	19	T2c1d1	Germany	GER	LBK	N	5100 - 4800	Stuttgart-GER-N	[41]
<b>5000 - 4000 cal BCE</b>										
Kumtepe	Kumtepe6	0.13	H2a	Turkey	TR	Chalcolithic	CHALC	4846 - 4618	Kumtepe-TR-CHALC	[133]
Apc-Berekalja I.	NE7	1.14	N1a1a1a	Hungary	HUN	Lengyel culture	LN	4490 - 4360	Apc-Berekalja-HUN-LN	[81]

<sup>1</sup>coverage calculated on Human Origins array ([102], [130])

Site	Sample IDs	Coverage	mtDNA HG	Country	abbr.	Culture (arch. context)	abbr.	Time cal BCE	Label PCA	Reference
Paliambela	Pal7	1.29	J1c1	Greece	GR	Late Neolithic	LN	4452 - 4350	Pal7-GR-LN	This study
Kleitios	Klei10	2.01	K1a2	Greece	GR	Final Neolithic	FN	4230 - 3995	Klei10-GR-FN	This study
<b>4000 - 1000 cal BCE</b>										
Gökhem	Gökhem2, 4, 5, 7	1.33, 0.04, 0.02, 0.01	H1c, H, K1e, H24	Sweden	SW	TRB	N	3330 - 2800	Gökhem-SW-N	[132]
Apc-Berekalja I.	CO1	1.13	H	Hungary	HUN	Late Copper Age	CA	2900 - 2700	Apc-Berekalja-HUN-CA	[81]
Ajvide	Ajvide52, 53, 58, 59, 70	0.09, 0.03, 2.22, 0.01, 0.16	V, U4d, U4d, U	Sweden	SW	Pitted Ware Culture	PWC	2950 - 2650	Ajvide-SW-PWC	[132]
Ire	Ire8	0.04	U4d	Sweden	SW	Pitted Ware Culture	PWC	3150 - 2200	Ire-SW-PWC	[132]
Kompolt-Kigyósér	BR1	0.81	K1c1	Hungary	HUN	Early Bronze Age	BA	2190 - 1980	Komp.-Kig.-HUN-BA	[81]
Ludas-Varjú-dűlű	BR2	21.25	K1a1a	Hungary	HUN	Late Bronze Age	BA	1270 - 1110	Lud.-Var.-HUN-BA	[81]

**Table S15:** *Barcm/Menteşe IDs and corresponding labels PCA. Coverage calculated on Human Origins array ([102], [130])*

Unique ID ([130])	Renaming PCA (this study)	Coverage	SNPs	mtDNA HG	Other IDs
I0707	BarM1	9,431	1.026.916	K1a4	BAR2 / L11-213
I0708	BarM2	6,948	1.007.876	N1b1a	BAR6 / L11-439
I0709	BarM3	9,765	1.015.958	U3	BAR20/ M13-170
I0723	MenM1	0,472	431.024	X2m2	T1, M229 / UH
I0724	MenM2	0,048	56.787	K1a4	T2 / UP
I0726	MenM3	0,221	236.641	H or H5-C16192T	M15, M15.2, M15.2 / UF
I0727	MenM4	0,039	46.069	K1a2	M24 / UA JK 16
I0736	BarM4	2,248	825.825	N1a1a1a	L11-216
I0744	BarM5	2,39	907.009	J1c11	M10-275
I0745	BarM6	7,785	1.025.390	U8b1b1	M11-363
I0746	BarM7	8,465	1.033.308	K1a or K1a1	L11-322
I1096	BarM8	2,865	780.073	N1a1a1	BAR26 / M10-76
I1097	BarM9	2,125	776.163	W1-T119C	BAR271 / M10-271
I1098	BarM10	2,994	809.388	X2d2	BAR99 / M10-352
I1099	BarM11	0,827	556.176	T2b	L11-S-488
I1100	BarM12	0,331	304.802	K1a or K1a6	M11-351
I1101	BarM13	1,56	680.133	T2b	M11-352a
I1102	BarM14	0,55	423.429	K1a3a	M11-354
I1103	BarM15	1,135	622.029	K1b1b1	M11-S-350
I1579	BarM16	2,507	847.361	K1a-C150T	M13-72
I1580	BarM17	2,903	934.855	H5	L12-393
I1581	BarM18	2,803	853.125	U3	L12-502
I1583	BarM19	9,396	1.012.449	K1a2	L14-200
I1585	BarM20	3,098	852.587	J1 or J1c	M11-59

## Results

The modern populations shown in Figure 2 in the main text roughly follow their geographic distribution. Near Eastern populations are located close to Western Asians (Turks, Kurds, Armenians), Iranians and Caucasians along the second principal component (PC2, 0.36% variance explained). Sardinians form a clear outlier compared to other Mediterranean populations which are flanked by Cypriots and Basques in direction of the first PC (PC1, 0.76% variance explained). French bridge the Mediterranean to Central and Eastern European populations, with samples from the British Isles and Scandinavia partially overlapping with Germans and Slavic populations.

The location of the ancient samples in reference to modern populations is consistent with previous reports [41, 102]. Holocene hunter-gatherers differentiate from the bulk of modern European and Near-Eastern populations along PC1, whereas Neolithic farmers differ from all modern populations except Sardinians along PC2. Middle Neolithic farmers from Scandinavia moved away from the

Early Neolithic samples and populate a distinct area. Figure 2 also contains older pleistocene hunter-gatherers, that fall close to modern Caucasians and Turks but are separated from modern populations by the third PC (0.23% variance explained).

The Greek and Anatolian genomes reported here locate close to Early Neolithic samples. They fall in direct proximity of KO2 and other Neolithic individuals from the Great Hungarian Plane [81], as well as Stuttgart [41], a genome from the Linearbandkeramik (LBK) culture in southwestern Germany. If the movement of Early- to Mid-Neolithic samples along PC1 is interpreted as a signature of genetic drift and admixture, then the position of the Greek and Anatolian genomes reported here suggests that they are representatives of early migrants that have not yet picked up the genetic signature resulting from admixture.

Within the five individuals newly reported here, two clusters can be observed. An early and late Greek genome fall nearly on top of each other (Rev5 and Pal7), whereas the to Anatolian (Bar8 and Bar31) and the latest of the Greek genomes (Klei10) cluster together.

**Table S16:** For each sample projected here, the number of sites from the reference data covered by at least one read, the average read depth over all reference sites and the  $t$ -statistic are given.

sample	sites covered by $\geq 1$ read	avg. depth	$t$ -stat.
Ust'Ishim	510857	39.1623	0.999918
Kostenki	440881	2.34704	0.998352
BR1	258878	0.714649	0.994931
BR2	510942	20.1961	0.99992
CO1	191583	0.595747	0.989068
IR2	291611	0.981857	0.994726
KO1	110479	0.368503	0.975656
KO2	44347	0.0918224	0.969886
NE1	510812	20.1356	0.999914
NE2	76499	0.162697	0.981807
NE3	46177	0.0951098	0.970114
NE4	38606	0.0789232	0.966231
NE5	298846	0.891962	0.995725
NE6	335769	1.08388	0.996464
NE7	329049	1.04291	0.996162
La Brana	464447	3.54143	0.998917
Loschbour	510852	19.2679	0.999896
Stgrt. LBK	510779	17.3774	0.999891
Ajivide52	44633	0.0962819	0.970217
Ajivide53	12998	0.0268746	0.908246
Ajivide58	450797	2.34339	0.998269
Ajivide70	76450	0.16948	0.981133
Gökhem2	307497	1.42543	0.996821
Gökhem4	16990	0.0367837	0.930021
Gökhem5	9270	0.0188772	0.879518
Gökhem7	5548	0.013995	0.829325
Ire8	19990	0.0413489	0.936403
StoraFörvar11	44153	0.090709	0.969098
Motala1	84522	0.181222	0.983445

sample	sites covered by $\geq 1$ read	avg. depth	<i>t</i> -stat.
Motala2	73153	0.154561	0.980785
Motala3	225500	0.584104	0.993969
Motala4	35535	0.0720764	0.96128
Motala6	12785	0.0253365	0.898789
Motala9	5031	0.00990136	0.795668
Motala12	460688	2.37039	0.998415
Kotias	510284	12.956	0.99981
Satsurbliia	359487	1.26943	0.996742
Bichon	508882	8.66574	0.999688
CB13	323372	1.03757	0.996285
Kum6	56563	0.117758	0.974658
I0707	420444	8.2573	0.999026
I0708	406577	6.15252	0.998853
I0709	408007	8.44978	0.998954
I0723	181287	0.489996	0.992495
I0724	22385	0.0458965	0.942461
I0725	19635	0.0409282	0.934697
I0726	98182	0.225759	0.985908
I0727	19613	0.0400927	0.935508
I0736	329398	2.05583	0.997601
I0744	366865	2.30579	0.997968
I0745	412760	6.8549	0.998943
I0746	416482	7.45434	0.999023
I0854	310137	1.7018	0.997088
I1096	319755	2.98223	0.99795
I1097	319313	2.21417	0.997574
I1098	335732	3.1158	0.998018
I1099	232253	0.855552	0.995115
I1100	128831	0.341039	0.990039
I1101	281646	1.62425	0.996995
I1102	178203	0.57035	0.993305
I1103	258893	1.17854	0.995919
I1579	352679	2.58461	0.99801
I1580	387425	2.99033	0.998289
I1581	354926	2.8851	0.998149
I1583	411113	9.78011	0.999097
I1585	354183	3.18912	0.998143
I0409	211729	0.66879	0.994465
I0410	214874	1.16888	0.99545
I0405	128393	0.355465	0.989912
I0407	218062	1.2101	0.995531
I0412	411864	12.3254	0.999094
I0408	386587	12.2416	0.998933
I0406	299127	3.62348	0.997701
I0413	279001	1.60645	0.996762
Klei10	455268	2.3389	0.998344
Pal7	377606	1.39845	0.997177
Rev5	335068	1.20947	0.996563
Bar8	507629	7.60681	0.999649
Bar31	476974	3.65432	0.998979



# SI7 Using $f$ -statistics to Infer Genetic Relatedness and Admixture Amongst Ancient and Contemporary Populations

Zuzana Hofmanová, Krishna R. Veeramah

## Introduction

We used a suite of methods developed by Patterson *et al.* [134] that utilize the expected correlation in allele frequencies along drift paths in a population tree in order to make inferences about the population history of our ancient Aegean samples within the context of other existing ancient and modern Eurasian genomic data. In particular we used these methods to infer a) the degree of shared ancestry (i.e. drift) between our two populations and b) the extent of admixture amongst three populations. The advantage of these methods are that they are highly robust to ascertainment bias, and, in the case of the  $f_4$ -statistic, can be applied to even very low coverage ancient genomes in order to make inferences about admixture.

## Methods

In order to make our results comparable to other recently published paleogenomic studies [40, 41, 81, 102, 104, 135], we limited our analyses to the  $\sim 300\text{K}$  SNP positions utilized in Haak *et al.* [102]. In addition, to ensure our results were robust to errors in genotype calling and comparable to the reference dataset, the ancient Greek and Anatolian genomes were analyzed using pseudo-haploid calls (using positions with a minimum genotype quality of Q30).

To account for the possible effects of low coverage and post-mortem DNA damage on our inferences, all analysis was repeated excluding transitions (C $\leftrightarrow$ T and G $\leftrightarrow$ A sites are the most likely to be affected by DNA damage). However, both Z-scores ( $\rho=0.91$ ) and the actual  $f$ -statistics ( $\rho=0.85$ ) were highly correlated (Spearman's rank correlation coefficient was calculated in R) between the full and filtered datasets indicating that our results should be relatively robust. It should be noted that as the ancient genomes vary widely in their coverage, the analysis of individual samples was based on different numbers of SNP positions after quality filtering, though this variation was accounted for in the construction of Z-scores to assess significance (see below).

Plink 1.9 [123, 124] was used for data formatting. Programs contained within the ADMIXTOOLS package [134] were used to calculate the  $f$ -statistics and associated Z-scores via a block jackknife resampling procedure using the software's default options.

## Relationship of Greek and Anatolian samples to other ancient and modern populations

The relative genetic similarity of Greek and Anatolian Neolithic samples to other ancient populations was estimated using an *outgroup*  $f_3$ -statistic [136]. An  $f_3$ -statistic is commonly notated as  $F_3(C : A, B)$ , where population C is examined as a *target* population for evidence of admixture with populations A and B, where a negative value indicates that C possesses ancestry from both populations. However, if C is chosen to be an *outgroup* population that has not experienced any post divergence gene flow with either A or B, then the value of the  $f_3$ -statistic will be a positive value proportional to the length of the shared drift path of populations A and B with C. Therefore, by fixing population A (or B) but substituting different populations for B (or A), it is possible to infer which populations (tested as B) are genetically more similar to population A based on the relative magnitudes of the  $f_3$ -statistics. Under a simple three-population tree model with no post divergence admixture, the relative  $f_3$ -statistics for different B populations will be proportional to the relative population divergence time of A and B and will be robust to differences in genetic drift that occurred after these populations diverged. However, if this simple tree model is violated, the  $f_3$ -statistics will reflect a more complex demographic history of different timing of population divergence, proportion of admixture and population-specific drift (for example changes in  $N_e$ ) that occurred during the period between divergence and the time of admixture.

Our form of the  $f_3$ -statistic was  $f_3(\neq\text{Khomani}; \text{TEST}, \text{Greek}/\text{Anatolian})$  where *TEST* was one of the ancient populations from the reference datasets (see Figure S8). The grouping of individuals to these ancient populations was kept as defined in Haak *et al.* [102] for comparability (detailed results of these tests in Dataset S2) with occasional shortening of the group names, whereas samples from additional studies are tested individually.

The  $\neq$ Khomani San were selected as an outgroup as they are considered to be the most diverged extant human population, having diverged from all other modern humans at least 100kya, and are highly unlikely to have experienced substantial Eurasian admixture [137, 138]. However, we also tested the Mbuti and Yoruba in the place of the  $\neq$ Khomani San and we note that our results were robust to the choice of sub-Saharan African population utilized. The results are also consistent for all the levels of filtering and for the dataset without transitions.

The greatest amount of genetic similarity as reflected by the largest  $f_3$ -statistics is generally found between the Greek and Anatolian Neolithic genomes generated in this study and the Chalcolithic Anatolian sample, Kumtepe6 (see Figure S8). Other populations demonstrating high  $f_3$  values with the Greek/Anatolian population (considered either separately or together) are other European Early and Middle Neolithic populations. Especially high amounts of shared drift can be seen with Spanish Neolithic farmers, LBK and Starcevo. This suggests a common ancestry component for populations found throughout Europe during this era, confirming the patterns we observed in the first two principle components of our PCA (see SI6). Interestingly Pal7 and Klei10 demonstrate relatively

high levels of shared drift with the KK and SATP genomes from Jones *et al.* [129] compared to the other Aegean samples, potentially indicating some admixture between late Neolithic Greeks and incoming Caucasus hunter-gatherers (CHG).

We also used the outgroup  $f_3$ -statistics to examine the level of genetic similarity between our Greek/Anatolian Neolithic samples and contemporary humans populations (see Figure S9 and detailed results in Dataset S2). The geographical distributions of the values can be seen in Figure S10. The modern populations with the highest  $f_3$ -statistics are those located in the Mediterranean area (Italians, Sardinians, Greeks etc.) as well as Basques. All the highest values were obtained for modern Sardinians, a population previously noted for its genetic similarity to early farmers [41], possibly because of their relative geographic isolation from mainland Europe. However, we did not observe particularly high genetic similarity between the ancient samples excavated in Anatolia and the geographically closest modern Turkish populations. This pattern is also supported by our PCA, mixture model and simulation-based continuity analysis (see SI6, SI9 and SI10).

## Genetic Structure of Aegean Neolithic populations

$f_4$ -statistics are commonly notated as  $F_4(A, B : C, D)$  and provide a more model-based framework with which to investigate population similarity than *outgroup*  $f_3$ -statistics, though their interpretation is also based on the extent to which populations share drift paths. If population D is chosen to be an outgroup, it provides a three-way test of population genetic similarity through the sum of genealogical topologies across loci. An  $f_4$  value of zero can be obtained in two ways. In the first, genealogical topologies are always consistent with population A being closer to population B (i.e. population A and B form a bifurcating clade). In the second there are balanced numbers of topologies where A is closer to C and B is closer to C. However, positive and negative  $f_4$  values can only be obtained if the most common topologies are (A,C)(B,D) and (B,C)(A,D) respectively. Thus, the  $f_4$ -statistic with an outgroup can be used to test whether population C is more similar to population A or B. An additional advantage of the  $f_4$ -statistic is that it can be applied on a per sample basis (i.e. there is no need to pool multiple samples to represent a population to obtain allele frequencies).

Both the PCA and *outgroup*  $f_3$ -statistics (see Figure S8) indicate that the Greek and Anatolian samples are highly similar compared to all other populations. However, the  $f_4$ -statistic allowed us to more explicitly examine the level of population genetic structure amongst these samples with regard to geography and chronology.

We first examined the  $f_4$ -statistics of the form  $f_4(\textit{Anatolian}, \textit{Greek}, \textit{Early\_farmer}, \neq \textit{Khomani})$  (see Table S17 and Dataset S2) in order to examine whether there are differences in the level of non-Aegean early farmer ancestry in Greek versus Anatolian samples. While not all pairwise comparisons were significant, there was a slight general trend of negative  $f_4$  values indicating that Greeks were more genetically similar to other early farming populations from Spain and Central Europe (perhaps

**Table S17:**  $f_4(\text{Anatolian}, \text{Greek}, \text{Early\_farmer}, \neq \text{Khomani})$ , values for  $|Z| > 2$  shown.

Bar31	Klei10	Spain_EN	$\neq$ Khomani	-0.0129	-2,017
Bar8	Pal7	LBK_EN	$\neq$ Khomani	-0.0103	-2,033
Bar31	Pal7	Spain_EN	$\neq$ Khomani	-0.0112	-2,130
Bar31	Pal7	LBKT_EN	$\neq$ Khomani	-0.0440	-2,163
Bar8	Pal7	Stuttgart	$\neq$ Khomani	-0.0199	-2,169
Bar8	Pal7	Stuttgart	$\neq$ Khomani	-0.0199	-2,169
Bar31	Pal7	Alberstedt_LN	$\neq$ Khomani	-0.0173	-2,193
Bar8	Pal7	Spain_MN	$\neq$ Khomani	-0.0154	-2,198
Bar8	Rev5	LBK_EN	$\neq$ Khomani	-0.0101	-2,223
Bar8	Klei10	LBK_EN	$\neq$ Khomani	-0.0133	-2,227
Bar8	Pal7	Alberstedt_LN	$\neq$ Khomani	-0.0204	-2,621
Bar8	Pal7	Spain_EN	$\neq$ Khomani	-0.0178	-2,820
Bar8	Klei10	Spain_MN	$\neq$ Khomani	-0.0230	-3,322
Bar8	Klei10	Spain_EN	$\neq$ Khomani	-0.0242	-3,518
Bar8	Klei10	Esperstedt_MN	$\neq$ Khomani	-0.0319	-3,757

indicative of a movement of Neolithic farmers across the Aegean sea from Anatolia into the rest of Europe, though this also may simply be an isolation by distance pattern). However, we note that no significant comparisons were observed when using only transition mutations (i.e. removing potential post-mortem damage but lowering the number of positions analyzed). If population structure did exist between the Anatolian and Greek Neolithic farmers it was likely relatively subtle.

Given that there is a substantial time gap of  $\sim 2,000$  years between Early Neolithic (Rev5, Bar8, Bar31) and Middle Neolithic (Klei10, Pal7) samples, we calculated  $f_4$ -statistics additionally in the forms  $f_4(\text{Greek1}, \text{Greek2}, \text{Early\_farmer}, \neq \text{Khomani})$  and  $f_4(\text{Bar8}, \text{Bar31}, \text{Early\_farmer}, \neq \text{Khomani})$ . We found no significant pairwise comparisons using these chronological groupings (see Dataset S2). Assuming the Aegean as the source of European Neolithic ancestry, this would indicate that once early European farmers diverged from this source, the Aegean populations remained relatively isolated from later European farmers (i.e. there were no major episodes of gene flow back into Aegean farming populations from the west).

When Western and Eastern hunter-gatherers were included in our analysis using  $f_4$ -statistics of the form  $f_4(\text{Aegean}, \text{Aegean}, \text{HG}, \neq \text{Khomani})$  (see Dataset S2), we obtained no significantly positive values for  $f_4(\text{Greek}, \text{Anatolian}, \text{HG}, \neq \text{Khomani})$  comparisons, suggesting that Aegean populations also formed a clade with respect to HG and we did not observe significant shared drift violating this tree.

**Table S18:**  $f_4(\text{Early\_farmer}, \text{Iceman}, \text{Aegean}, \neq \text{Khomani})$ , values for  $|Z| > 3$  shown.

HungaryGamba_CA	Iceman	Pal7	$\neq$ Khomani	-0.0373	-5,409
SwedenSkoglund_MN	Iceman	Bar31	$\neq$ Khomani	-0.0358	-4,062
HungaryGamba_EN	Iceman	Pal7	$\neq$ Khomani	-0.0207	-3,801
SwedenSkoglund_MN	Iceman	Pal7	$\neq$ Khomani	-0.0305	-3,542
HungaryGamba_EN	Iceman	Bar31	$\neq$ Khomani	-0.0193	-3,084
HungaryGamba_EN	Iceman	Rev5	$\neq$ Khomani	-0.0228	-3,021

## Aegean as the source for early farming populations in Europe

We next estimated  $f_4$  values using the form  $f_4(\text{Early\_farmer}, \text{post-Neolithic}, \text{Neolithic Greek/Anatolian}, \neq \text{Khomani})$  in order to formally examine whether the Greek and Anatolian Neolithic samples were genetically closer to other Early and Middle Neolithic European farmers compared to other post Neolithic European and Middle Eastern populations, as indicated by the first two principle components of our PCA analysis. As expected, almost all  $f_4$  values were significantly positive (with  $|Z| \gg 3$ , see Dataset S2) (for the only exception, see Iceman below), consistent with the Aegean Neolithic populations being more genetically similar, as defined by greater levels of shared drift paths, to other early European farmers than with any other tested populations from more recent eras. Similarly,  $f_4$  tests of the form  $f_4(\text{Early\_farmer}, \text{HG}, \text{Neolithic Greek/Anatolian}, \neq \text{Khomani})$  also clearly demonstrated that the Neolithic Aegeans (see Dataset S2) were genetically more related to early farmers than any hunter-gatherer populations.

Interestingly, comparisons of our Neolithic Greek/Anatolian samples with the Late Neolithic/Early Bronze age Iceman [104] resulted in significantly negative values when compared to Neolithic farmers (see Table S18 and Dataset S2). If our Aegean populations are assumed to be the source of Neolithic genetic ancestry, it is thus possible that Ötzi and his ancestors did not admix with local populations after an initial spread from the Aegean to the same extent as other Middle and Late Neolithic farmer populations. This unique drift shared between Iceman and Aegeans might also suggest that the ancestors of this individual either shared substantial exchange with the Aegean farmer core area after the original spread, or they migrated from the core area later.

Based on the genetic similarity between the Early and Middle Neolithic populations and the archaeological context of the samples, it is reasonable to assume that this genetic ancestry arose in and around Anatolia and spread out to the rest of Europe. Given the presence of Early Neolithic farmers stretching from Anatolia all the way to Spain, might this spread have arisen via a serial range expansion moving westwards? If this was the case, then some non-Aegean early farmers further along this route might be expected to share unique drift compared to the original source Aegean farmers. However,  $f_4$ -statistics of the form  $f_4(\text{Early\_farmer2}, \text{Aegean}, \text{Early\_farmer1}, \neq \text{Khomani})$  using all pairwise combinations of non-Aegean Early Neolithic farmers from different geographic locations (Spain, Hungary/Central Europe) demonstrated negative or non-significant values rather

than positive ones with the exception of  $f_4(LBK\_EN, Bar8, SPAIN\_EN, \neq Khomani)$  (see Table S19).

These results suggest that either the initial Neolithic expansion from the Aegean region to the rest of Europe involved multiple independent migrating groups from the same source or was very rapid such that there was insufficient time for genetic differentiation. Interestingly, positive  $f_4$  values were obtained when comparing pairs of non-Aegean early farmers from the same region but different time periods, suggesting some geographic-specific drift once these populations were established and diverged from Aegeans. Some positive values were also obtained when comparing Middle Neolithic Spanish to Middle Neolithic Esperstedt, which may be the result of the proposed resurgence of hunter-gather ancestry during this era and the Late Neolithic. It is of interest to note that the positive significant values were observed almost exclusively with Anatolian samples (see in Table S20), which is consistent with the previous observation from the  $f_4(Anatolian, Greek, Early\_farmer, \neq Khomani)$  test that the Greek Aegean samples are genetically closer to Early Neolithic farmers.

Under a scenario of a rapid migration from Central Europe and then to Spain, we would assume that non-Aegean farmers would form a clade to the exclusion of Aegeans. However, when performing a test of the form  $f_4(Early\_farmer1, Early\_farmer2, Aegean, \neq Khomani)$ , we observed unique drift between Aegeans and Spanish farmers (see Table S21). This points to a gene flow event through the Mediterranean between Greece and Spain that did not include Central Europe. Given the previously discussed result (see S17) and the archaeological record, the most likely scenario would be an independent migration of Aegean farmers to Spain distinct from an initial migration to Central Europe (though migration from Spain back to the Aegean would also fit the data). Due to the observation of unique drift between LBK and Spanish early farmers after their split from Bar8 (one significant positive value for  $f_4(LBK\_EN, Bar8, SPAIN\_EN, \neq Khomani)$ , see Table S19), we can speculate that it happened chronologically after this individual lived.

## The relationship between Neolithic Aegeans and Chalcolithic Anatolians

Given their geographic proximity, the Aegean population characterized by the genomes sequenced in this study could potentially be the source population for both the Anatolian Kumtepe [133] that is dated to Chalcolithic as well as European Neolithic farmers. Interestingly,  $f_4$  tests of the form  $f_4(Aegean, Kumtepe, Early\_farmer, \neq Khomani)$  were often significantly positive (see Table S22 and Dataset S2), suggesting that Aegeans share ancestry with Neolithic European farmers (especially with LBK, Starcevo and Early Hungarian Neolithic farmers) not present in Kumtepe samples. Thus, the Kumtepe likely descend from an Aegean population that was somewhat differentiated from the one that expanded from Anatolia into the rest of Europe.

We also examined whether Kumtepe shared more unique drift with Anatolian samples from Barcin or later Greek samples by performing  $f_4$  tests of the form  $f_4(Greek, Kumtepe, Neo\_Anatolian, \neq Khomani)$  and  $f_4(Neo\_Anatolian, Kumtepe, Greek, \neq Khomani)$  (Table S23). Kumtepe6 demon-

**Table S19:**  $f_4(\text{Early\_farmer2}, \text{Aegean}, \text{Early\_farmer1}, \neq \text{Khomani})$ , values for  $Z > 3$  shown.

Esperstedt_MN	Bar8	Baalberge_MN	$\neq$ Khomani	0.0248	3,158
LBK_EN	Bar8	Stuttgart	$\neq$ Khomani	0.0153	3,290
LBK_EN	Bar8	Stuttgart	$\neq$ Khomani	0.0153	3,290
Starcevo_EN	Rev5	HungaryGamba_EN	$\neq$ Khomani	0.0277	3,316
<b>Spain_EN</b>	<b>Bar8</b>	<b>LBK_EN</b>	<b><math>\neq</math>Khomani</b>	<b>0.0130</b>	<b>3,352</b>
Starcevo_EN	Bar8	HungaryGamba_EN	$\neq$ Khomani	0.0203	3,360
Spain_EN	Bar31	Spain_MN	$\neq$ Khomani	0.0161	3,505
Starcevo_EN	Bar8	Esperstedt_MN	$\neq$ Khomani	0.0438	3,532
Spain_MN	Rev5	Esperstedt_MN	$\neq$ Khomani	0.0221	3,603
Spain_EN	Bar31	Esperstedt_MN	$\neq$ Khomani	0.0249	3,614
Baalberge_MN	Bar8	Esperstedt_MN	$\neq$ Khomani	0.0304	3,732
Spain_MN	Bar31	Esperstedt_MN	$\neq$ Khomani	0.0289	3,816
Stuttgart	Bar8	Esperstedt_MN	$\neq$ Khomani	0.0218	4,004
Stuttgart	Bar8	Esperstedt_MN	$\neq$ Khomani	0.0218	4,004
Starcevo_EN	Bar8	LBK_EN	$\neq$ Khomani	0.0269	4,057
Esperstedt_MN	Bar31	LBK_EN	$\neq$ Khomani	0.0188	4,075
Esperstedt_MN	Bar31	Spain_MN	$\neq$ Khomani	0.0220	4,379
Esperstedt_MN	Bar8	Spain_EN	$\neq$ Khomani	0.0245	4,647
Spain_MN	Bar8	Spain_EN	$\neq$ Khomani	0.0234	4,696
LBK_EN	Bar31	Esperstedt_MN	$\neq$ Khomani	0.0294	4,755
Esperstedt_MN	Bar8	Spain_MN	$\neq$ Khomani	0.0329	5,514
Spain_EN	Bar8	Spain_MN	$\neq$ Khomani	0.0276	5,667
Esperstedt_MN	Bar8	LBK_EN	$\neq$ Khomani	0.0299	6,215
Spain_EN	Bar8	Esperstedt_MN	$\neq$ Khomani	0.0371	6,441
LBK_EN	Bar8	Esperstedt_MN	$\neq$ Khomani	0.0394	7,725
Spain_MN	Bar8	Esperstedt_MN	$\neq$ Khomani	0.0417	8,020

**Table S20:**  $f_4(\text{Early\_farmer2}, \text{Aegean}, \text{Early\_farmer1}, \neq \text{Khomani})$ , values for  $Z < -3$  shown.

HungaryGamba_EN	Pal7	Spain_MN	$\neq$ Khomani	-0.0314	-7,625
HungaryGamba_EN	Klei10	Spain_MN	$\neq$ Khomani	-0.0322	-7,047
HungaryGamba_EN	Pal7	Spain_EN	$\neq$ Khomani	-0.0282	-6,825
Stuttgart	Pal7	Spain_MN	$\neq$ Khomani	-0.0277	-5,806
Stuttgart	Pal7	Spain_MN	$\neq$ Khomani	-0.0277	-5,806
HungaryGamba_EN	Klei10	Spain_EN	$\neq$ Khomani	-0.0279	-5,332
Stuttgart	Klei10	Spain_MN	$\neq$ Khomani	-0.0288	-5,055
Stuttgart	Klei10	Spain_MN	$\neq$ Khomani	-0.0288	-5,055
Stuttgart	Bar31	Spain_MN	$\neq$ Khomani	-0.0206	-4,642
Stuttgart	Bar31	Spain_MN	$\neq$ Khomani	-0.0206	-4,642
LBK_EN	Pal7	Spain_MN	$\neq$ Khomani	-0.0155	-4,292
HungaryGamba_EN	Bar31	Spain_MN	$\neq$ Khomani	-0.0194	-4,268
LBK_EN	Pal7	Spain_EN	$\neq$ Khomani	-0.0139	-4,119
Spain_MN	Pal7	Stuttgart	$\neq$ Khomani	-0.0295	-3,950
Spain_MN	Pal7	Stuttgart	$\neq$ Khomani	-0.0295	-3,950
Stuttgart	Pal7	Spain_EN	$\neq$ Khomani	-0.0244	-3,875
Stuttgart	Pal7	Spain_EN	$\neq$ Khomani	-0.0244	-3,875
HungaryGamba_EN	Rev5	Spain_MN	$\neq$ Khomani	-0.0206	-3,711
HungaryGamba_EN	Rev5	Spain_EN	$\neq$ Khomani	-0.0179	-3,610
HungaryGamba_EN	Pal7	LBK_EN	$\neq$ Khomani	-0.0124	-3,560
HungaryGamba_EN	Bar31	Spain_EN	$\neq$ Khomani	-0.0141	-3,460
Stuttgart	Klei10	Spain_EN	$\neq$ Khomani	-0.0237	-3,447
Stuttgart	Klei10	Spain_EN	$\neq$ Khomani	-0.0237	-3,447
Spain_MN	Bar31	Stuttgart	$\neq$ Khomani	-0.0187	-3,193
Spain_MN	Bar31	Stuttgart	$\neq$ Khomani	-0.0187	-3,193
Spain_MN	Pal7	HungaryGamba_EN	$\neq$ Khomani	-0.0143	-3,180
Stuttgart	Rev5	Spain_MN	$\neq$ Khomani	-0.0232	-3,113
Stuttgart	Rev5	Spain_MN	$\neq$ Khomani	-0.0232	-3,113
SwedenSkoglund_MN	Klei10	LBK_EN	$\neq$ Khomani	-0.0165	-3,046



*Table S21:  $f_4(\text{Early\_farmer1}, \text{Early\_farmer2}, \text{Aegean}, \neq \text{Khomani})$ , values for  $Z < -2$  shown.*

HungaryGamba_EN	Spain_EN	Pal7	$\neq$ Khomani	-0.0236	-5.379
HungaryGamba_EN	Spain_EN	Klei10	$\neq$ Khomani	-0.0255	-5.201
HungaryGamba_EN	LBK_EN	Rev5	$\neq$ Khomani	-0.0142	-4.352
HungaryGamba_EN	Spain_MN	Klei10	$\neq$ Khomani	-0.0194	-4.26
HungaryGamba_EN	Spain_EN	Rev5	$\neq$ Khomani	-0.0172	-3.859
Stuttgart	Starcevo_EN	Klei10	$\neq$ Khomani	-0.0377	-3.387
Stuttgart	Starcevo_EN	Klei10	$\neq$ Khomani	-0.0377	-3.387
HungaryGamba_EN	Starcevo_EN	Klei10	$\neq$ Khomani	-0.029	-3.369
HungaryGamba_EN	LBK_EN	Bar31	$\neq$ Khomani	-0.0092	-3.334
HungaryGamba_EN	LBK_EN	Klei10	$\neq$ Khomani	-0.0125	-3.264
HungaryGamba_EN	Spain_MN	Pal7	$\neq$ Khomani	-0.0171	-3.245
Stuttgart	Spain_EN	Klei10	$\neq$ Khomani	-0.0201	-3.133
Stuttgart	Spain_EN	Klei10	$\neq$ Khomani	-0.0201	-3.133
HungaryGamba_EN	LBK_EN	Pal7	$\neq$ Khomani	-0.0118	-3.107
Stuttgart	Starcevo_EN	Bar8	$\neq$ Khomani	-0.0264	-3.095
Stuttgart	Starcevo_EN	Bar8	$\neq$ Khomani	-0.0264	-3.095
LBK_EN	Starcevo_EN	Bar8	$\neq$ Khomani	-0.0194	-3.031
HungaryGamba_EN	Starcevo_EN	Rev5	$\neq$ Khomani	-0.0269	-2.996
LBK_EN	Spain_EN	Klei10	$\neq$ Khomani	-0.013	-2.871
HungaryGamba_EN	Spain_EN	Bar31	$\neq$ Khomani	-0.0128	-2.811
LBK_EN	Spain_EN	Pal7	$\neq$ Khomani	-0.0121	-2.8
Stuttgart	Spain_MN	Klei10	$\neq$ Khomani	-0.0154	-2.548
Stuttgart	Spain_MN	Klei10	$\neq$ Khomani	-0.0154	-2.548
HungaryGamba_EN	Starcevo_EN	Bar8	$\neq$ Khomani	-0.0175	-2.453
HungaryGamba_EN	Stuttgart	Pal7	$\neq$ Khomani	-0.018	-2.25
HungaryGamba_EN	Stuttgart	Pal7	$\neq$ Khomani	-0.018	-2.25
HungaryGamba_EN	Spain_MN	Rev5	$\neq$ Khomani	-0.0111	-2.19
LBK_EN	Starcevo_EN	Klei10	$\neq$ Khomani	-0.0185	-2.107
HungaryGamba_EN	Starcevo_EN	Pal7	$\neq$ Khomani	-0.0203	-2.10

**Table S22:**  $f_4(\text{Aegean}, \text{Kumtepe}, \text{Early\_farmer}, \neq \text{Khomani})$ , values for  $|Z_3| > 3$  shown.

Bar8	Kumtepe6	HungaryGamba_EN	$\neq$ Khomani	0.0706	4.526
Bar31	Kumtepe6	LBK_EN	$\neq$ Khomani	0.0571	4.321
Klei10	Kumtepe6	HungaryGamba_BA	$\neq$ Khomani	0.083	4.101
Bar8	Kumtepe6	LBK_EN	$\neq$ Khomani	0.0514	4.002
Klei10	Kumtepe6	HungaryGamba_EN	$\neq$ Khomani	0.0616	3.71
Klei10	Kumtepe6	Stuttgart	$\neq$ Khomani	0.0851	3.63
Klei10	Kumtepe6	Stuttgart	$\neq$ Khomani	0.0851	3.63
Bar31	Kumtepe6	HungaryGamba_BA	$\neq$ Khomani	0.0682	3.588
Pal7	Kumtepe4	Esperstedt_MN	$\neq$ Khomani	0.2761	3.502
Bar8	Kumtepe6	Corded_Ware_LN	$\neq$ Khomani	0.0604	3.476
Klei10	Kumtepe6	Bell_Beaker_LN	$\neq$ Khomani	0.0633	3.439
Bar8	Kumtepe6	HungaryGamba_BA	$\neq$ Khomani	0.0662	3.436
Bar31	Kumtepe6	Bell_Beaker_LN	$\neq$ Khomani	0.0493	3.404
Bar31	Kumtepe6	HungaryGamba_EN	$\neq$ Khomani	0.0519	3.331
Klei10	Kumtepe6	LBK_EN	$\neq$ Khomani	0.0576	3.33
Klei10	Kumtepe4	Esperstedt_MN	$\neq$ Khomani	0.2394	3.268
Rev5	Kumtepe6	HungaryGamba_EN	$\neq$ Khomani	0.0649	3.188
Pal7	Kumtepe4	Spain_MN	$\neq$ Khomani	0.1862	3.164
Bar31	Kumtepe6	HungaryGamba_CA	$\neq$ Khomani	0.0833	3.058
Klei10	Kumtepe6	Starcevo_EN	$\neq$ Khomani	0.1247	3.028
Pal7	Kumtepe6	HungaryGamba_CA	$\neq$ Khomani	0.1179	3.018
Klei10	Kumtepe6	Esperstedt_MN	$\neq$ Khomani	0.0863	3.013

strated unique drift with Neolithic Greeks, especially Late Neolithic ones (Klei10, Pal7), which could be explained by gene flow that was maintained over the Aegean throughout the Neolithic. Results for Kumtepe4 showed indications of shared ancestry in the opposite direction (i.e. greater affinity with Barcm), but this result was barely significant, perhaps as a consequence of the much lower coverage of this genome.

Finally,  $f_4$ -statistics of the form  $f_4(\text{Aegean}, \text{Kumtepe}, \text{CHG}, \neq \text{Khomani})$  showed that CHG populations shared unique drift with Kumtepe6 when compared to both Greek and Anatolian Aegeans (Table S25). Though little is known about hunter-gatherers in Anatolia, this suggests that towards the end of, or directly following, the Neolithic expansion there was gene flow from the Caucasus and neighboring regions to Anatolia. If there was continued gene flow across the Aegean at this time between Greece and Anatolia, this would also be compatible with the  $f_3$  *outgroup* results which show the later Greek samples to be closer to CHG than the Rev5 and two early Neolithic Anatolian samples.

**Table S23:**  $f_4(\text{Greek}, \text{Kumtepe}, \text{Neo\_Anatolian}, \neq\text{Khomani})$  and  $f_4(\text{Neo\_Anatolian}, \text{Kumtepe}, \text{Greek}, \neq\text{Khomani})$ , all values shown. *Kumtepe6* shows negative, *Kumtepe4* positive values.

Bar8	Kumtepe6	Pal7	$\neq$ Khomani	-0.1438	-6.821
Bar31	Kumtepe6	Pal7	$\neq$ Khomani	-0.1183	-4.571
Bar8	Kumtepe6	Klei10	$\neq$ Khomani	-0.0672	-3.097
Bar8	Kumtepe6	Rev5	$\neq$ Khomani	-0.0652	-2.739
Pal7	Kumtepe6	Bar31	$\neq$ Khomani	-0.0465	-1.935
Rev5	Kumtepe6	Bar8	$\neq$ Khomani	-0.0395	-1.557
Rev5	Kumtepe6	Bar31	$\neq$ Khomani	-0.0547	-1.544
Bar31	Kumtepe6	Klei10	$\neq$ Khomani	-0.0296	-1.018
Pal7	Kumtepe6	Bar8	$\neq$ Khomani	-0.0255	-1.004
Bar31	Kumtepe6	Rev5	$\neq$ Khomani	-0.0114	-0.415
Klei10	Kumtepe6	Bar8	$\neq$ Khomani	-0.0072	-0.291
Klei10	Kumtepe6	Bar31	$\neq$ Khomani	-0.0017	-0.055
Bar31	Kumtepe4	Klei10	$\neq$ Khomani	0.0473	0.599
Bar8	Kumtepe4	Rev5	$\neq$ Khomani	0.1316	1.47
Klei10	Kumtepe4	Bar8	$\neq$ Khomani	0.1158	1.493
Klei10	Kumtepe4	Bar31	$\neq$ Khomani	0.1029	1.513
Pal7	Kumtepe4	Bar31	$\neq$ Khomani	0.1203	1.725
Bar8	Kumtepe4	Klei10	$\neq$ Khomani	0.1132	1.787
Bar31	Kumtepe4	Rev5	$\neq$ Khomani	0.1647	1.857
Rev5	Kumtepe4	Bar8	$\neq$ Khomani	0.1453	2.026
Bar31	Kumtepe4	Pal7	$\neq$ Khomani	0.1825	2.243
Bar8	Kumtepe4	Pal7	$\neq$ Khomani	0.1655	2.283
Pal7	Kumtepe4	Bar8	$\neq$ Khomani	0.1892	2.327
Rev5	Kumtepe4	Bar31	$\neq$ Khomani	0.2766	3.529

**Table S24:**  $f_4(\text{Aegeans}, \text{Kumtepe}, \text{WHG}, \neq \text{Khomani})$ , all values shown.

Pal7	Kumtepe4	LaBranal	$\neq$ Khomani	0.0846	1.143
Rev5	Kumtepe4	LaBranal	$\neq$ Khomani	0.0805	1.082
Klei10	Kumtepe4	LaBranal	$\neq$ Khomani	0.0731	1.015
Pal7	Kumtepe6	LaBranal	$\neq$ Khomani	0.0026	0.088
Klei10	Kumtepe6	LaBranal	$\neq$ Khomani	0.0007	0.025
Rev5	Kumtepe6	LaBranal	$\neq$ Khomani	-0.0003	-0.008
Bar31	Kumtepe6	LaBranal	$\neq$ Khomani	-0.0006	-0.025
Bar8	Kumtepe6	LaBranal	$\neq$ Khomani	-0.002	-0.086
Bar31	Kumtepe4	LaBranal	$\neq$ Khomani	-0.0538	-0.725
Bar8	Kumtepe4	LaBranal	$\neq$ Khomani	-0.0756	-1.203
Rev5	Kumtepe6	Loschbour	$\neq$ Khomani	0.0771	2.791
Bar8	Kumtepe6	Loschbour	$\neq$ Khomani	0.0521	2.24
Klei10	Kumtepe6	Loschbour	$\neq$ Khomani	0.0609	2.12
Bar31	Kumtepe6	Loschbour	$\neq$ Khomani	0.0418	1.438
Rev5	Kumtepe4	Loschbour	$\neq$ Khomani	0.0933	1.21
Pal7	Kumtepe6	Loschbour	$\neq$ Khomani	0.0298	0.92
Pal7	Kumtepe4	Loschbour	$\neq$ Khomani	0.0091	0.119
Klei10	Kumtepe4	Loschbour	$\neq$ Khomani	-0.0139	-0.211
Bar8	Kumtepe4	Loschbour	$\neq$ Khomani	-0.0167	-0.248
Bar31	Kumtepe4	Loschbour	$\neq$ Khomani	-0.1371	-1.797

**Table S25:**  $f_4(\text{Aegean}, \text{Kumtepe}, \text{CHG}, \neq \text{Khomani})$ , all values shown.

Bar8	Kumtepe6	SATP	$\neq$ Khomani	-0.1584	-6.017
Bar8	Kumtepe6	KK1	$\neq$ Khomani	-0.0855	-4.505
Rev5	Kumtepe6	KK1	$\neq$ Khomani	-0.099	-4.116
Bar31	Kumtepe6	KK1	$\neq$ Khomani	-0.0836	-3.639
Rev5	Kumtepe6	SATP	$\neq$ Khomani	-0.1119	-3.218
Bar31	Kumtepe6	SATP	$\neq$ Khomani	-0.0917	-2.954
Klei10	Kumtepe6	SATP	$\neq$ Khomani	-0.1049	-2.93
Pal7	Kumtepe6	SATP	$\neq$ Khomani	-0.0975	-2.82
Pal7	Kumtepe6	KK1	$\neq$ Khomani	-0.0728	-2.81
Klei10	Kumtepe6	KK1	$\neq$ Khomani	-0.0684	-2.793
Bar8	Kumtepe4	SATP	$\neq$ Khomani	0.0699	0.813
Bar31	Kumtepe4	KK1	$\neq$ Khomani	0.0733	0.945
Bar31	Kumtepe4	SATP	$\neq$ Khomani	0.111	1.248
Bar8	Kumtepe4	KK1	$\neq$ Khomani	0.0913	1.609
Pal7	Kumtepe4	KK1	$\neq$ Khomani	0.1294	1.858
Klei10	Kumtepe4	SATP	$\neq$ Khomani	0.21	2.156
Pal7	Kumtepe4	SATP	$\neq$ Khomani	0.2387	2.876
Rev5	Kumtepe4	SATP	$\neq$ Khomani	0.2859	2.877
Klei10	Kumtepe4	KK1	$\neq$ Khomani	0.2282	3.118
Rev5	Kumtepe4	KK1	$\neq$ Khomani	0.2555	3.833

## Hunter-gather contributions to farming societies

Recent studies have shown that European Neolithic populations likely experienced some level of western hunter-gatherer (WHG) admixture. In particular, Haak *et al.* [102] have suggested there was a resurgence of hunter-gatherer ancestry in Middle and Later Neolithic European farmers. Our Anatolian farmers likely possess genetic ancestry that is most representative of the ancestral Neolithic component, thus presenting an opportunity for us to refine our understanding of degree of Neolithic vs. WHG admixture in Europe.

Again we used the  $f_4$ -statistics, this time of the form  $f_4(\text{Neolithic\_farmer}, \text{Anatolian}, \text{HG}, \neq \text{Khomani})$  (for the results see Dataset S2). The observation of a positive value under this test would indicate admixture between the Neolithic farmer and hunter-gatherer populations. The results shown (see Table S26) utilize Loschbour to represent *HG* (other HG in Dataset S2).

Amongst Early Neolithic populations, only Neolithic Spain and Hungarian early farmer show significant positive values ( $|Z| > 3$ ; see Table S26) and therefore evidence of hunter-gatherer gene flow. As previously noted by Haak *et al.* [102], there is, however, evidence of the resurgence of the hunter-gatherer admixture component in Middle and Late Neolithic samples from Spain, Hungary and Central Europe (see Dataset S2).

Regarding the differential affinities of Kumtepe and Neolithic Aegeans to hunter-gatherers, we did not observe significant drift with WHG for any of Neolithic or Chalcolithic Aegeans studied. If the Final and Chalcolithic Aegeans samples are representative of their respective populations, we can therefore conclude that the WHG resurgence did not happen in the Aegean. However, it should be noted that for  $f_4(\text{Aegeans}, \text{Kumtepe6}, \text{Loschbour}, \neq \text{Khomani})$ , the values are positive and some Z-scores are above 2 for Bar8 and Rev5 (see Table S24) showing rather an opposite trend (the decrease of WHG-like drift over time in the Aegean).

When we examined an  $f_4$ -statistic of the form  $f_4(\text{Neolithic\_farmer}, \text{Aegeans}, \text{CHG}, \neq \text{Khomani})$ , we obtained almost exclusively negative results, consistent with CHG admixture with the Aegean (see Dataset S2). Again, consistent with the results described the  $f_3$  tests, an  $f_4$  test of the form of  $f_4(\text{Aegean}, \text{Aegean}, \text{CHG}, \neq \text{Khomani})$  (Table S27) demonstrated greater shared drift between CHG and Late Neolithic Greeks.

*Table S26:  $f_4(\text{Early\_farmer}, \text{Greek/Anatolian}, \text{HG}, \neq \text{Khomani})$ , all values shown.*

HungaryGamba_EN	Bar8	Loschbour	$\neq$ Khomani	0.0205	4,189
HungaryGamba_EN	Bar31	Loschbour	$\neq$ Khomani	0.0271	3,786
HungaryGamba_EN	Klei10	Loschbour	$\neq$ Khomani	0.0182	2,551
HungaryGamba_EN	Rev5	Loschbour	$\neq$ Khomani	0.0155	2,397
HungaryGamba_EN	Pal7	Loschbour	$\neq$ Khomani	0.0140	2,205
LBK_EN	Pal7	Loschbour	$\neq$ Khomani	0.0045	0,807
LBK_EN	Bar31	Loschbour	$\neq$ Khomani	0.0201	2,977
LBK_EN	Bar8	Loschbour	$\neq$ Khomani	0.0133	2,738
LBK_EN	Klei10	Loschbour	$\neq$ Khomani	0.0132	1,984
LBK_EN	Rev5	Loschbour	$\neq$ Khomani	0.0061	1,066
LBKT_EN	Pal7	Loschbour	$\neq$ Khomani	0.0089	0.433
LBKT_EN	Rev5	Loschbour	$\neq$ Khomani	0.0348	1,822
LBKT_EN	Bar8	Loschbour	$\neq$ Khomani	0.0216	1,386
LBKT_EN	Bar31	Loschbour	$\neq$ Khomani	0.0215	1,248
LBKT_EN	Klei10	Loschbour	$\neq$ Khomani	-0.0186	-1,077
Spain_EN	Bar8	Loschbour	$\neq$ Khomani	0.0258	5,239
Spain_EN	Bar31	Loschbour	$\neq$ Khomani	0.0307	4,215
Spain_EN	Klei10	Loschbour	$\neq$ Khomani	0.0267	3,868
Spain_EN	Rev5	Loschbour	$\neq$ Khomani	0.0175	2,771
Spain_EN	Pal7	Loschbour	$\neq$ Khomani	0.0168	2,625
Stuttgart	Bar31	Loschbour	$\neq$ Khomani	0.0226	2,809
Stuttgart	Bar8	Loschbour	$\neq$ Khomani	0.0155	2,190
Stuttgart	Klei10	Loschbour	$\neq$ Khomani	0.0172	2,018
Stuttgart	Pal7	Loschbour	$\neq$ Khomani	0.0112	1,305
Stuttgart	Rev5	Loschbour	$\neq$ Khomani	0.0097	1,105

**Table S27:**  $f_4(\text{Aegean}, \text{Aegean}, \text{CHG}, \neq \text{Khomani})$ , all values shown.

Bar8	Bar31	SATP	$\neq$ Khomani	-0.0298	-3.877
Bar8	Bar31	KK1	$\neq$ Khomani	-0.0279	-3.487
Rev5	Bar31	SATP	$\neq$ Khomani	-0.018	-1.659
Rev5	Bar31	KK1	$\neq$ Khomani	-0.0126	-1.291
Klei10	Pal7	SATP	$\neq$ Khomani	-0.0166	-1.249
Klei10	Pal7	KK1	$\neq$ Khomani	-0.0091	-0.934
Rev5	Bar8	KK1	$\neq$ Khomani	0.0078	1.021
Rev5	Bar8	SATP	$\neq$ Khomani	0.0121	1.188
Pal7	Bar31	KK1	$\neq$ Khomani	0.0219	2.056
Klei10	Bar31	SATP	$\neq$ Khomani	0.0214	2.11
Klei10	Bar31	KK1	$\neq$ Khomani	0.0228	2.369
Pal7	Bar31	SATP	$\neq$ Khomani	0.0334	2.69
Klei10	Rev5	SATP	$\neq$ Khomani	0.0312	2.746
Pal7	Rev5	KK1	$\neq$ Khomani	0.0349	3.298
Pal7	Rev5	SATP	$\neq$ Khomani	0.0476	3.318
Klei10	Rev5	KK1	$\neq$ Khomani	0.0321	3.322
Pal7	Bar8	SATP	$\neq$ Khomani	0.0671	5.69
Pal7	Bar8	KK1	$\neq$ Khomani	0.0528	5.752
Klei10	Bar8	KK1	$\neq$ Khomani	0.0491	6.07
Klei10	Bar8	SATP	$\neq$ Khomani	0.0535	6.398



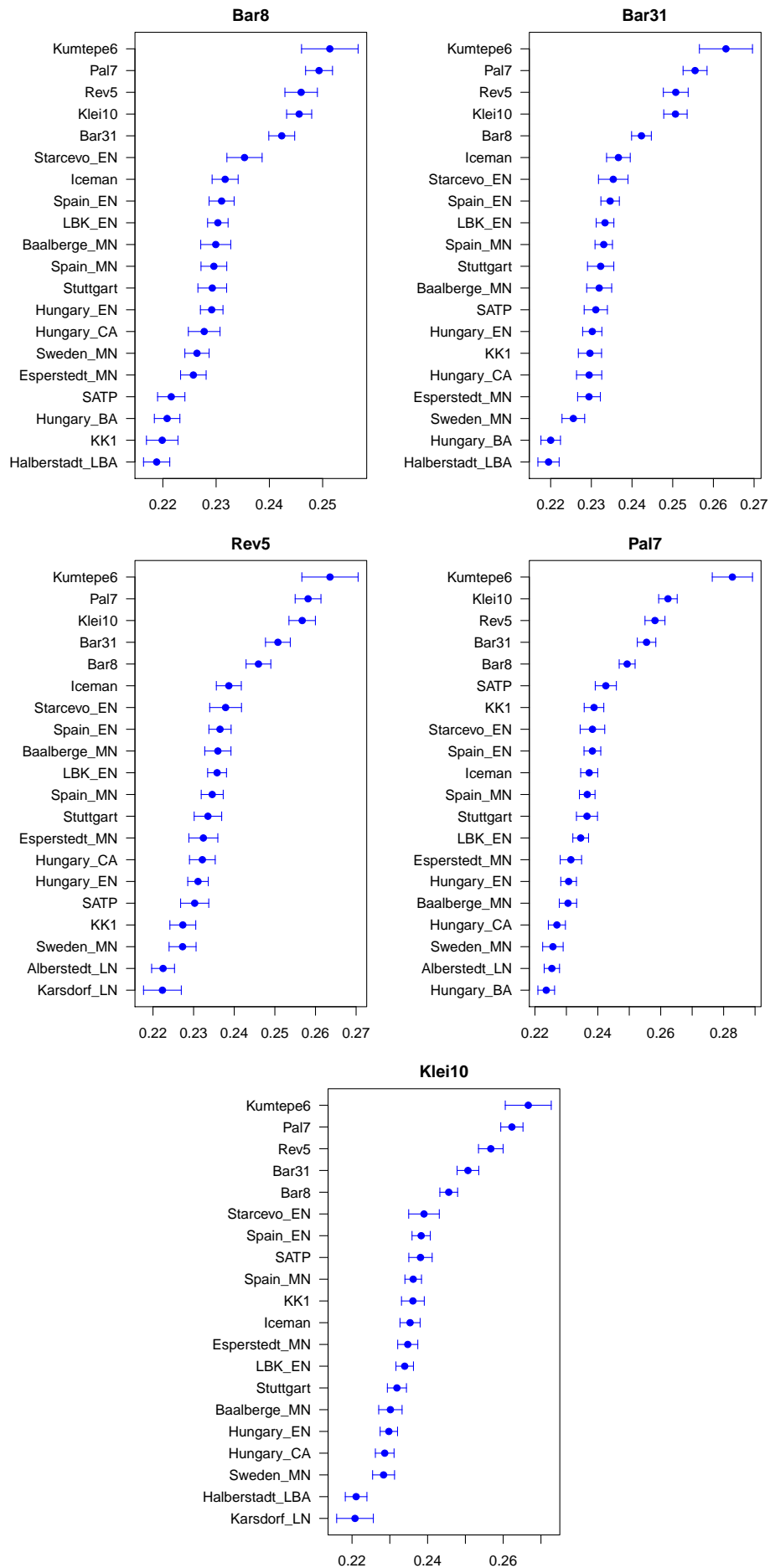


Figure S8:  $f_3(\neq \text{Khomani}; \text{Ancient\_population}, \text{Greek/Anatolian})$ . The highest 20 values shown.

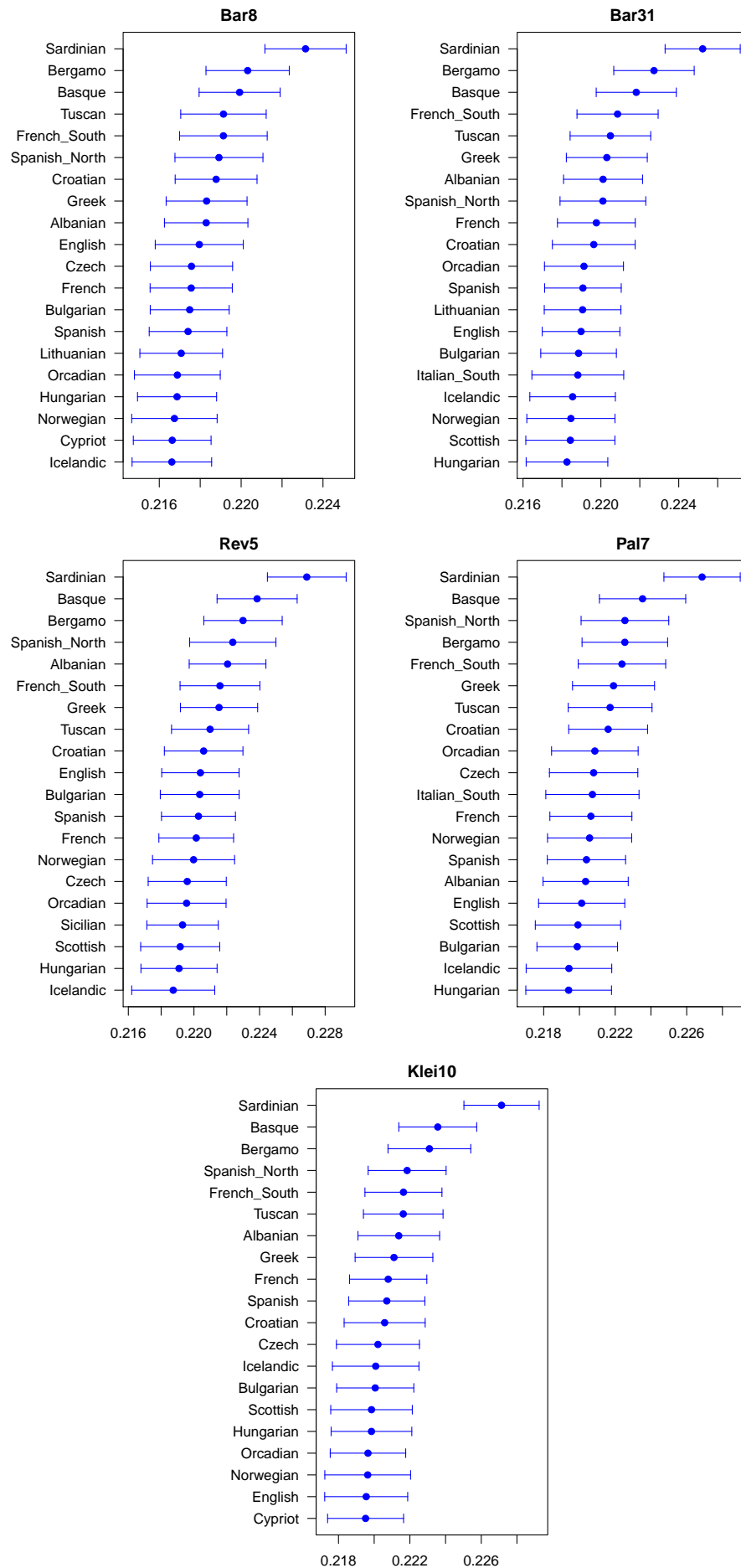
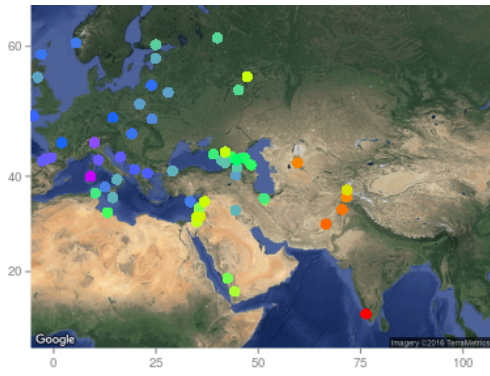
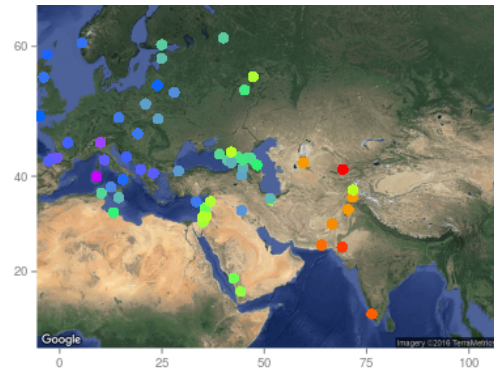


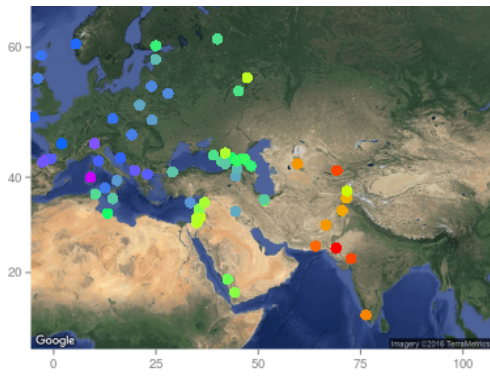
Figure S9:  $f_3(\neq \text{Khomani}; \text{Modern\_population}, \text{Greek/Anatolian})$ . The highest 20 values shown.



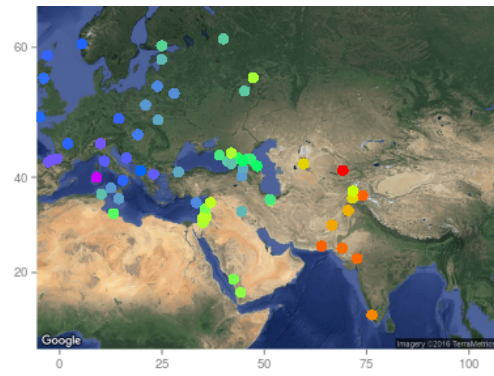
(a) Bar8



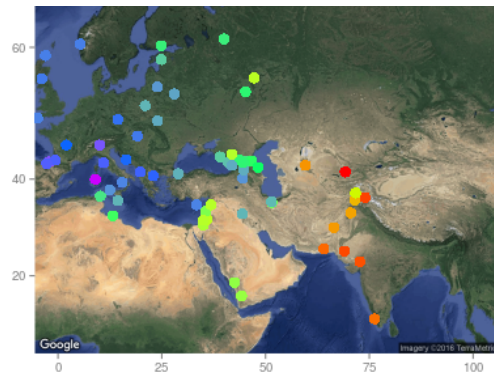
(b) Bar31



(c) Rev5



(d) Pal7



(e) Klei10

**Figure S10:**  $f_3(\neq \text{Khomani}; \text{Modern\_population}, \text{Greek/Anatolian})$ . Values above 0.2 in the relative geographical proximity shown.

# SI8 Proportions of ancestral clusters in Neolithic populations of Europe

Zuzana Hofmanová, Krishna R. Veeramah

## Introduction

We used ADMIXTURE [139] in order to identify ancestral clusters in European samples from the early, middle and late Neolithic. This model-based approach for estimating ancestry in large autosomal datasets does not require *a priori* assignment of individuals to populations for analysis (the so called *unsupervised* approach). This approach can help guide other population genetics methods that require defined populations and allows the identification of migrants, admixed individuals and whole distinct populations. However, the analysis can also be performed by incorporating a subset of individuals of known population ancestry, which can improve the inference of ancestry of other unknown individuals (the so called *supervised* approach). We applied both supervised and unsupervised approaches in our analysis.

## Methods

Allele presence calls were filtered for a Q30 score or greater (see SI5) and merged with the dataset of Haak *et al.* [102] with addition of other ancient DNA samples as in SI7. Known relatives of individuals from the reference dataset (the same as in SI7) were excluded from the analysis (relatives marked in Haak *et al.* [102]). Similarly, SNPs that showed evidence of linkage disequilibrium were removed using PLINK [124]. The maximum  $r^2$  value was set to 0.5 and SNPs were analyzed in sliding windows (window size of 200 SNPs, sliding 50 SNPs per step). As ADMIXTURE is designed primarily for diploid genotype data, the program was run by treating each haploid allele presence call as a homozygote.

ADMIXTURE analysis was initially performed unsupervised and was limited to Neolithic and hunter-gatherer samples. The number of clusters to be estimated varied from  $K=2$  to  $K=8$  and the analysis for each  $K$  consisted of 100 independent runs with differing seeds. The cross-validation error (5-fold) was calculated to determine the optimal  $K=2$ . Results for each  $K$  were matched in CLUMPP [140] and plotted in DISTRUCT [141]. Supervised ADMIXTURE analysis was also performed for  $K=2$  where the allele frequencies were assumed to be known for two populations, with samples from Anatolia (Bar8 and Bar31) and Motula serving as proxies for the ancestral farmer and hunter-gatherer populations respectively. Similarly, the analysis for  $K=2$  to  $K=8$  was performed for the dataset including CHG individuals (KK1, SATP) from Jones *et al.* [129] and supervised runs were performed for  $K=3$  with the CHG set to as an additional known cluster. Additionally, the analysis for  $K=2$  to  $K=8$  and the supervised analysis for  $K=3$  (supervised for CHG, Motula and Anatolian) was repeated while including Yamnaya individuals, a group that are considered to be

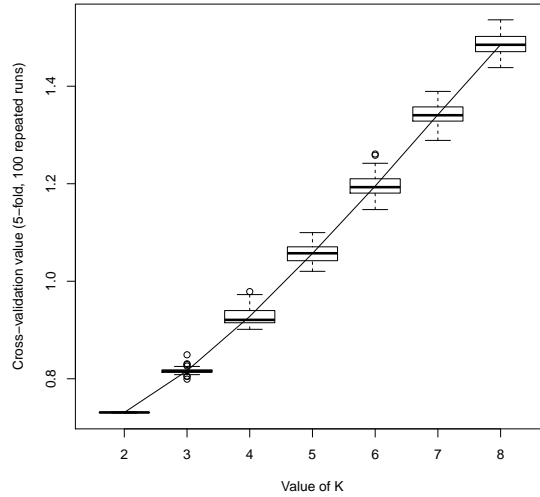
the descendants of a major hunter-gatherer migration from the east during late Neolithic [102]. It should be noted that for some runs the supervised version of ADMIXTURE did not output 100% assignment to a single cluster for some fixed source individuals (though values were very close, 99%), which we believe may be due to a potential bug resulting from a rounding error. Consistent with the assumed model, in cases where this occurred individuals were restored to having 100% ancestry from a single cluster when plotted.

## Anatolian samples form an ancestral cluster to Neolithic farmers in Europe

The results of our ADMIXTURE analysis for Neolithic and Mesolithic samples are shown in Figure S32. The cross-validation error was lowest at  $K=2$  (Figure S12). The sample clustering is highly similar between the supervised and unsupervised run (see Figure S14 and Figure S32). This suggests that the Anatolian samples could be considered as good proxies for the ancestral farmer component, though we note that most other early Neolithic farmers also show the same ancestry component with no evidence of admixture with hunter-gatherers. The only exceptions are NE1, NE3 and NE4 (data from Gamba *et al.* [81]). This result agrees with the  $f_4$ -statistic analysis (see SI7, Table S26), where the HungaryGamba\_EN group containing these samples also demonstrates an apparent signal of admixture with hunter-gatherers. Interestingly an older Neolithic sample from the same region (KO2 from Gamba *et al.* [81]) demonstrates no evidence of hunter-gatherer admixture, while another sample of the same age (KO1 from Gamba *et al.* [81]) is genetically most similar to hunter-gatherers. While hunter-gatherer ancestry is largely absent in Early Neolithic farmers according to ADMIXTURE results, it is increasingly apparent transitioning into the Middle and Late Neolithic. It should be noted that Kumtepe4 is also showing apparent admixture with the non-farmer cluster, however under higher  $K$ , it is obvious that there is no high affinity of Kumtepe4 to Western hunter-gatherers.

## CHG affinities to farmers

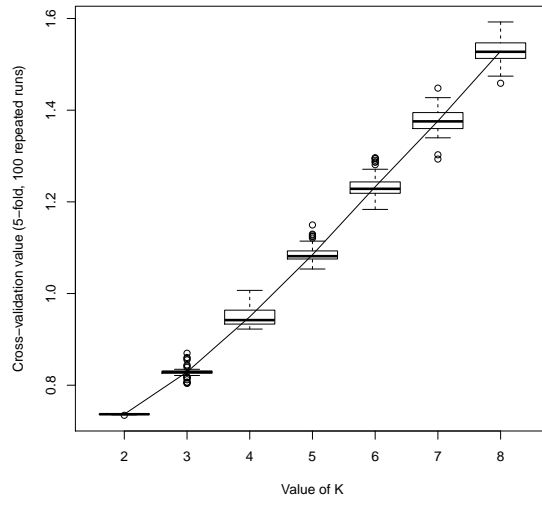
Results for  $K=2$  to  $K=8$  when including CHG samples are shown in Figure S32. The cross-validation error did not change with the addition of the CHG samples (the lowest for  $K=2$ ) (see Figure S11). For the most likely clustering of  $K=2$ , the main conclusion of all Early Neolithic samples clustering with Aegeans was maintained. For  $K=3$ , all Neolithic samples demonstrated mixed ancestry with at least some CHG-defined component in addition to the WHG-defined component described above. Interestingly, the CHG cluster was found at a higher proportion in Aegeans than other Early Neolithic samples, especially for Kumtepe4. The difference between Kumtepe4 and earlier Aegeans in terms of higher CHG influence was also observed using  $f$ -statistics (SI7).



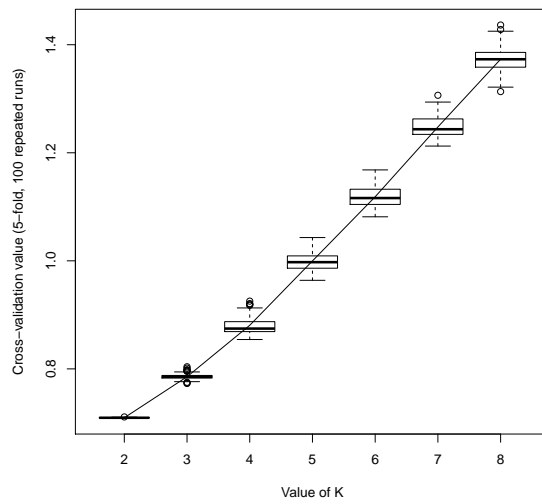
**Figure S11:** Cross-validation error (5-fold) of 100 iterations of unsupervised admixture run with CHG samples.

### Yamnaya signal in Late Neolithic

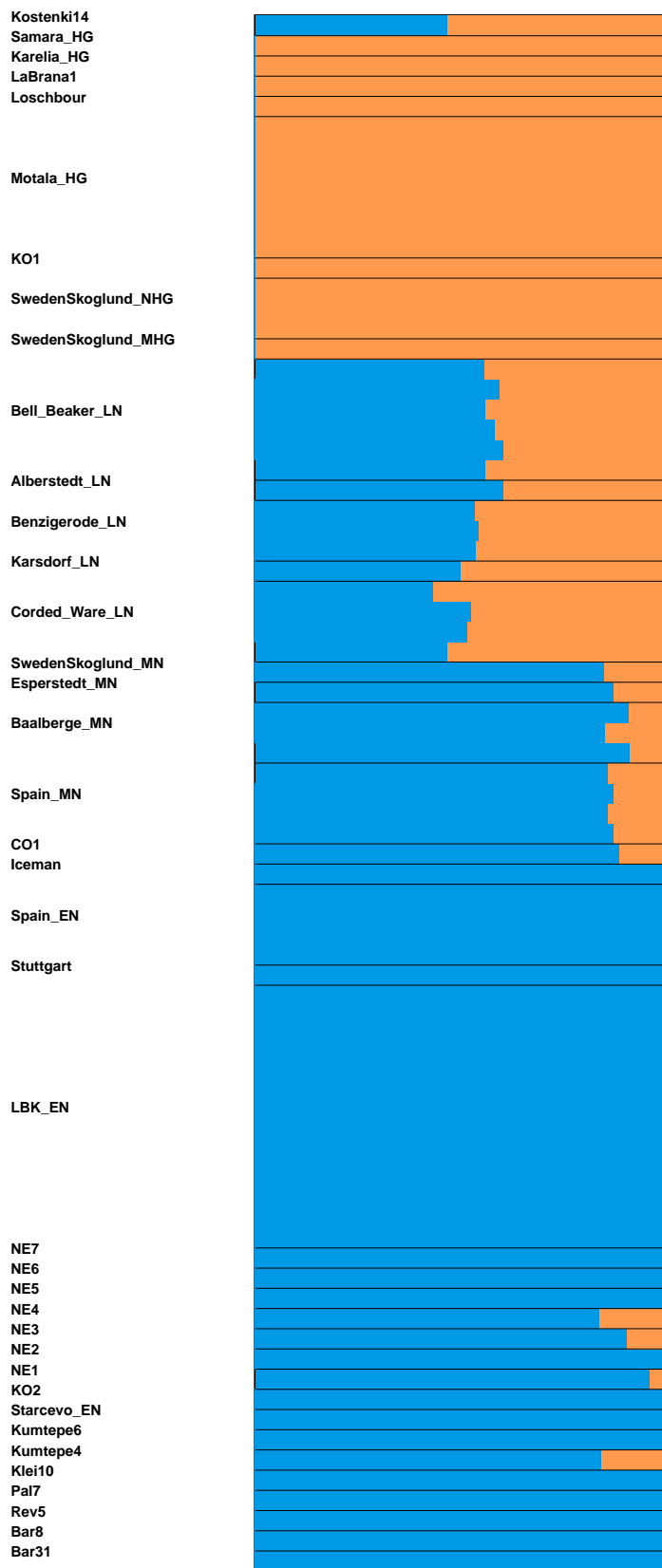
The results of our ADMIXTURE analysis for the dataset including also Yamnaya samples are shown in Figure S32. The cross-validation error was lowest for  $K=2$  (Figure S13). Supervised (Figure S16) and unsupervised analyses are again highly concordant (see Figure S32). Early Neolithic farmers again demonstrate almost no evidence of hunter-gatherer admixture, while it is still observed in the Middle Neolithic farmers. However, much of the Late Neolithic hunter-gatherer ancestry from the previous analysis is replaced by Yamnaya ancestry. These results are consistent with Haak *et al.* [102] who demonstrated a resurgence of hunter-gatherer ancestry followed by the establishment of eastern hunter-gatherer ancestry.



**Figure S12:** Cross-validation error (5-fold) of 100 iterations of unsupervised admixture run with Neolithic samples.

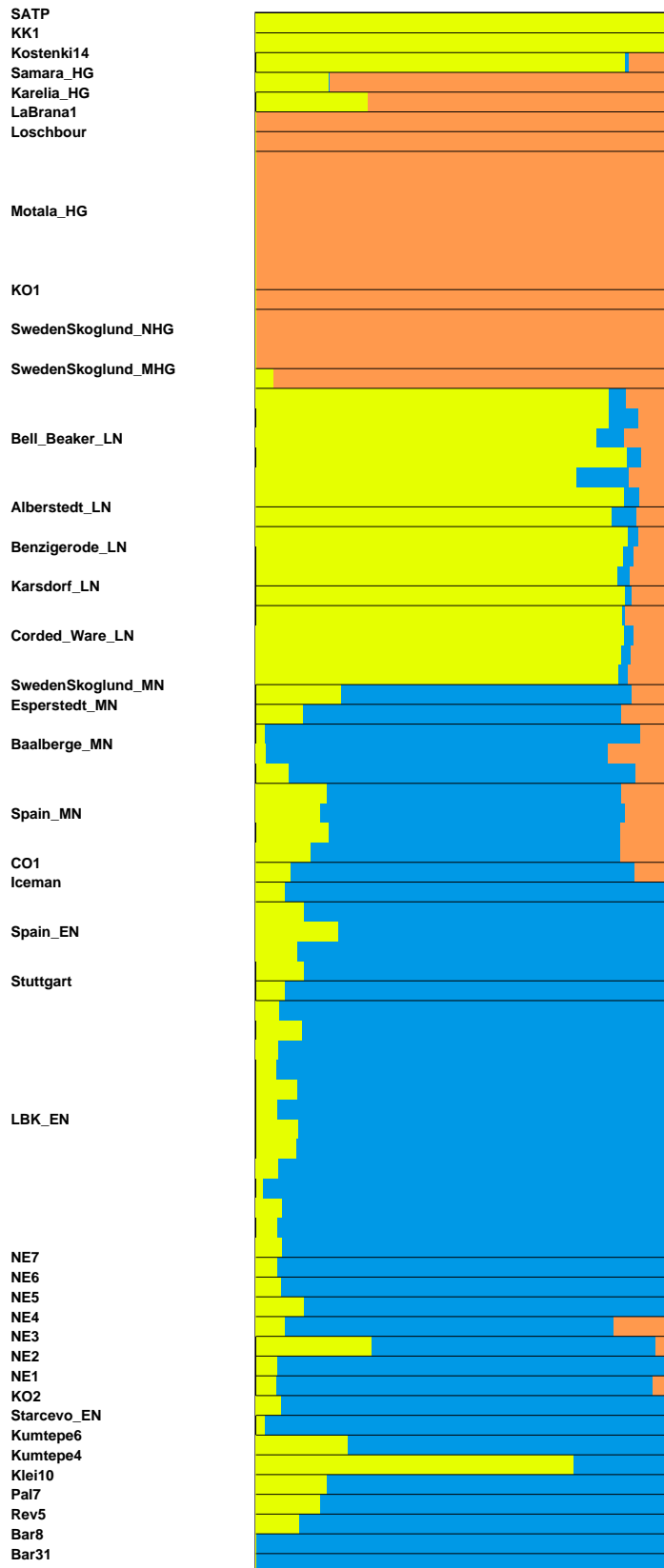


**Figure S13:** Cross-validation error (5-fold) of 100 iterations of unsupervised admixture run with Yamnaya samples.

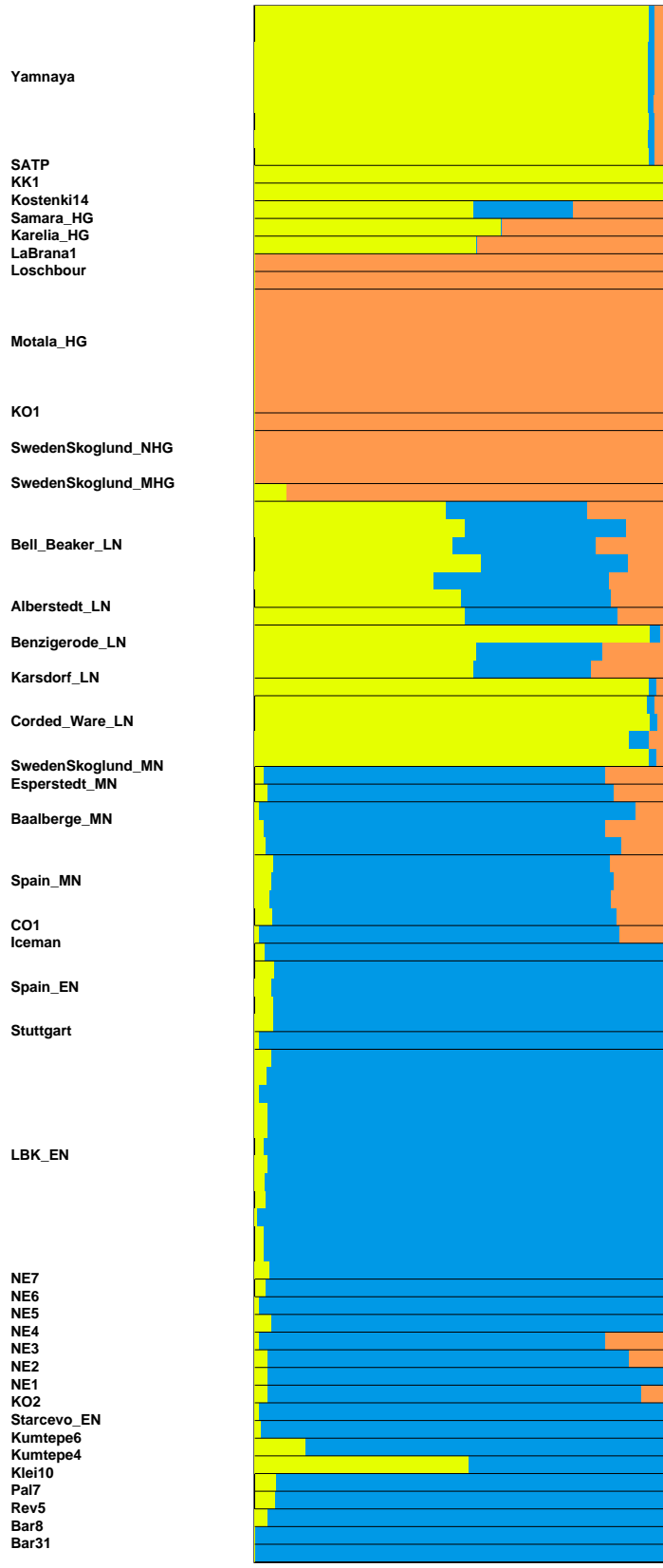


**Figure S14:** Supervised run of ADMIXTURE. The clusters to supervise were chosen to best fit the presumed ancestral populations (for HG Motala and for farmers Bar8 and Bar31).





**Figure S15:** Supervised run of ADMIXTURE. The clusters to be supervised were chosen to best fit the presumed ancestral populations (for WHG Motala, for CHG KK1 and SATP and for farmers Bar8 and Bar31).



**Figure S16:** Supervised run of ADMIXTURE. The clusters to be supervised were chosen to best fit the presumed ancestral populations (for HG Motala and for farmers Bar8 and Bar31 and for later Eastern migration Yamnaya).

## SI9 Population continuity

David Díez-del-Molino, Yoan Diekmann, Mark G. Thomas

### Method overview

A population can be strictly defined as continuous through time if it has not experienced admixture or replacement from other populations. Although this definition is rather strict and unlikely to hold for most population histories over long time periods, it does provide a natural null-hypothesis, and has been tested using single locus genetic systems such as mitochondrial DNA [e.g. 39, 101]. However, challenges remain in formally testing population continuity using ancient genomic data, although relevant efforts are being made [e.g. 142].

Using this rationale, we applied a forward simulation approach to test for population continuity using a single ancient genome and a sample of genomes from a modern population. The test is based on comparison of observed and expected proportions of allele sharing classes between the ancient and modern genomes (see Table S28). We consider only haploid calls for the ancient genome to avoid genotype calling biases due to low coverage. Under the null-hypothesis of genetic drift in a continuous population, we use the shape of the site frequency spectrum (SFS) of a modern population as the reference from which to simulate allele frequency trajectories, and generate distributions of the expected proportions of allele sharing classes. We applied the following steps:

- (1) We consider the overlapping positions between the haploid calls with  $GQ > 30$  of our ancient genomes and the biallelic calls of a sample of modern genomes [41] to estimate the population allele frequencies (Table S28). Alleles and frequencies were polarized to ancestral and derived states by comparing them with the chimpanzee genome (panTro2).
- (2) In order to incorporate the sampling uncertainty in estimates of modern allele frequencies, we sampled 100 frequency vectors using the beta distribution and the Jeffreys prior [143] from the distribution of allele frequencies of that SFS.
- (3) Using these frequency vectors as a starting point, we simulated forward the genetic drift process in order to generate possible allele trajectories. In each generation allele frequencies were updated using a binomial sampling based on the frequency in the previous generation. The two explored parameters are the ancient ( $N_{e,a}$ ) and modern ( $N_{e,m}$ ) effective population sizes. We assumed exponential growth between the ancient and modern population to set the population size at each generation.
- (4) In each simulation we sampled a haploid genome from the initial frequency vector and a diploid genome from the final generation's simulated allele frequency vector. In order to compare the observed with the simulated data we defined six allelic sharing classes formed by all possible combinations of haploid/diploid genotypes for the same position (Table S28). Allele sharing fraction

values are calculated for both the observed and simulated data as the proportion of all analyzed positions that fall into each one of these six classes.

**Table S28:** *The six possible allele sharing classes for the comparison of the haploid calls of the ancient genome ( $t_0$ ) and the biallelic calls of each modern genome ( $t_n$ ) for the same position.*

	Match		Mismatch		Half match	
	a	b	a	b	a	b
$t_0$	A	D	A	D	A	D
$t_n$	AA	DD	DD	AA	AD	AD

(5) One-tailed p-values were obtained by calculating the proportion of simulated allele sharing class fraction values that are greater than the observed values. These p-values were transformed in two-tailed p-values by applying the formula  $1 - 2 * |0.5 - P|$ . Following Fisher’s method, p-values for each allele sharing class were combined to generate a chi-square statistic  $C_{obs}$  as defined by:

$$C = -2 \sum_{i=1}^k \log(P_i)$$

where  $P_i$  is the estimated p-value of the  $i$ -th allele sharing class of  $k$  classes [144]. To generate an expected distribution of C values under the null hypothesis of population continuity ( $C_{sim}$ ), sharing fraction values from each simulation were compared to those from all other simulations as above [144, 145]. Finally, an overall two-tailed p-value for the continuity test was calculated comparing  $C_{obs}$  and  $C_{sim}$  values.

If the observed sharing class fractions between the ancient and modern population samples can be explained under the null-hypothesis of genetic drift in a continuous population (p-value  $> 0.05$ ), population continuity between those populations can not be rejected. However, if sharing class fractions cannot be explained by drift alone (p-value  $< 0.05$ ) then other demographic processes should be invoked. Because we analyze genome-wide variation and test over a relative short evolutionary period, selection and mutation on individual loci are unlikely to greatly affect the results of this test.

## Results

We only considered loci that were covered in Rev5, Pal7, Klei10 and the modern genomes, and tested them individually against modern Greeks (Table S29). In order to explore the plausible parameter space of  $N_e a$  and  $N_e m$  for each ancient Greek genome the continuity test was performed in a 30x30 grid composed of values of these effective sizes ranging from just 10 individuals to 10 million (a 10th of the population in actual Greece is  $\sim 1.1$  million), spaced equally on a log scale. For each combination of parameters we performed 10,000 simulations; a total of 9,000,000 simulations per

ancient genome (Figure S17). We also tested for continuity of Bar31 and Bar8 against a modern Turkish population sample, as described above.

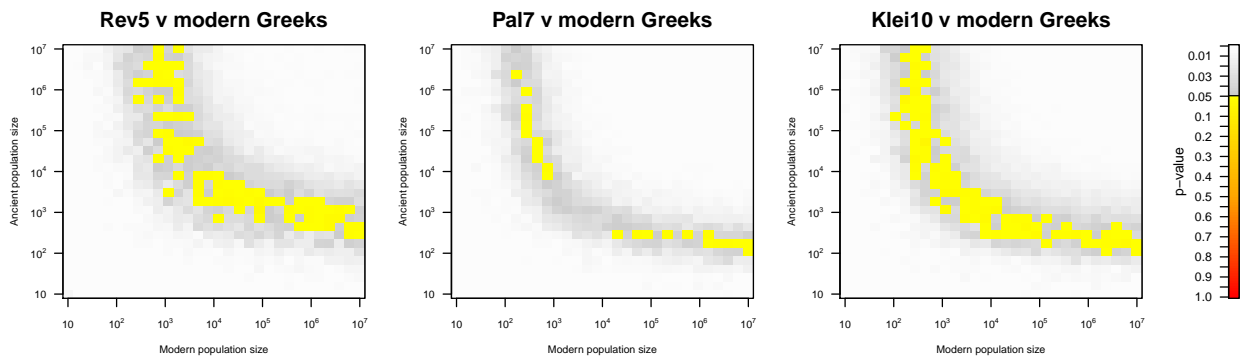
**Table S29:** *Overlapping positions between our ancient genomes and [41] dataset. Generation time was calculated using the mean Cal. age and assuming 25 years per generation.*

Name	Mean Cal. age (BP)	Epoch	Generations	Modern genomes	Overlapping SNPs
Rev5	6,395	Early Neolithic	256	38	200,775
Bar8	6,328	Early Neolithic	253	64	572,950
Ba31	6,221	Early Neolithic	249	64	469,189
Pal7	4,401	Late Neolithic	176	38	246,491
Klei10	4,105	Late Neolithic	164	38	394,860

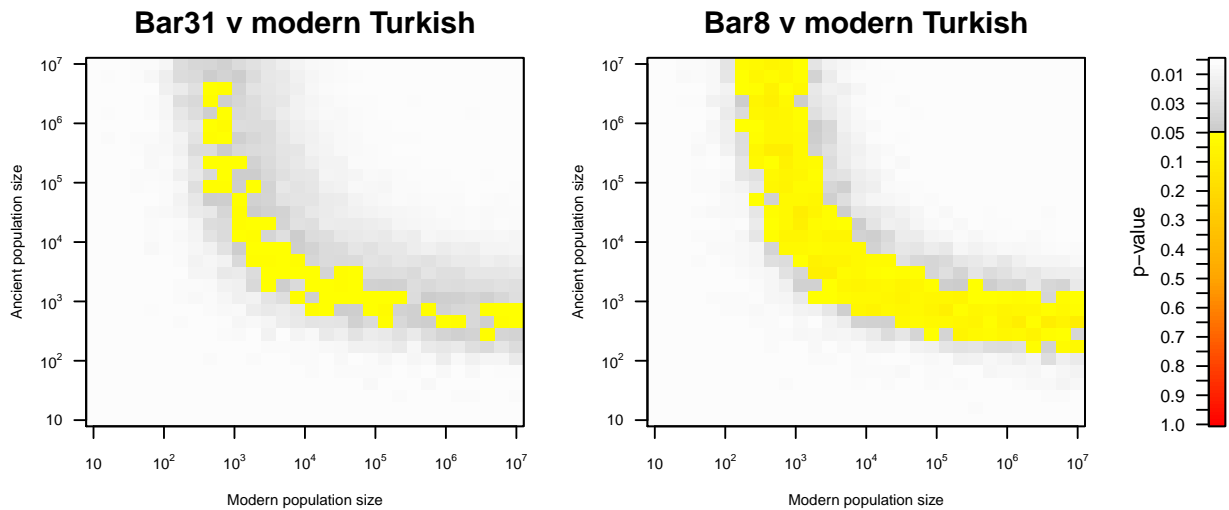
To further examine the range of effective population sizes where continuity could not be rejected, we sliced the grid in 1, 5, 10 and 20% of the modern population size and reported the p-values of the test for each ancient effective size. A modern population size of 11 million and 77 million were assumed for modern Greeks and Turkish populations, respectively.

Grid results do not support a model population continuity between any of the ancient and modern Greek comparisons. However, there are some combinations of the explored parameters for which continuity cannot be rejected (yellow areas, Figure S17). Grid slices indicated that continuity cannot be rejected for Rev5, Pal7 and Klei10 if ancient effective population size was unrealistically small ( $N_e$  ranges 728-1,887, 174-281, and 174-452, respectively; Figure S19). Similarly, population continuity between our Anatolian ancient genomes and modern Turkish populations was rejected for most plausible parameter combinations (Figure S18); population continuity cannot be rejected between Bar31 or Bar8 and the modern Turkish samples for unrealistically small effective ancestral population sizes ( $N_e$  ranges 452-728 and 174-1,172, respectively; Figure S20).

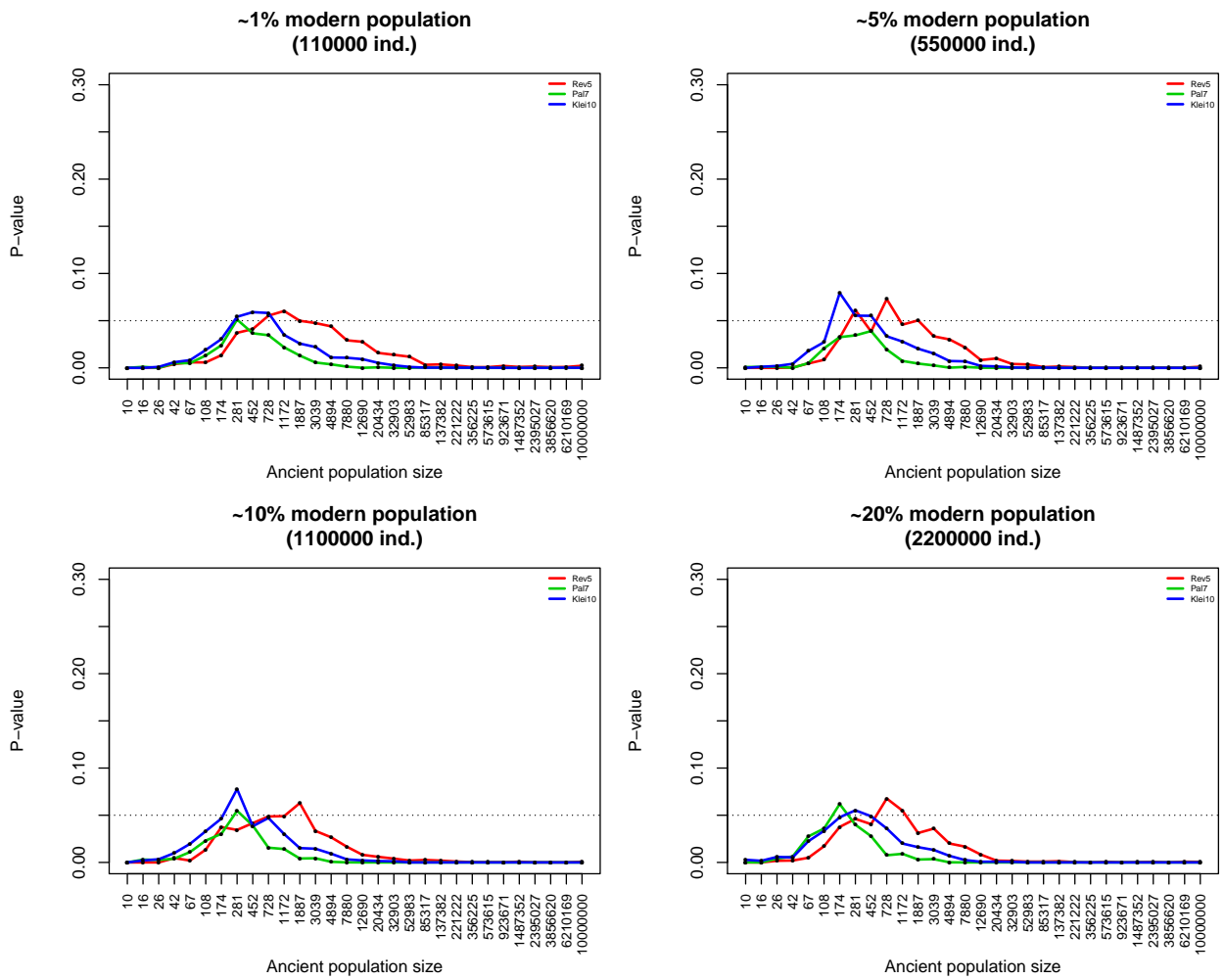
Together, these results indicate that the ancient genomes are not sampled from a population continuous with modern Greeks or Turks. However, genetic drift alone would still be able to explain the differences seen between our ancient genomes and the modern populations samples if the ancient populations were very small.



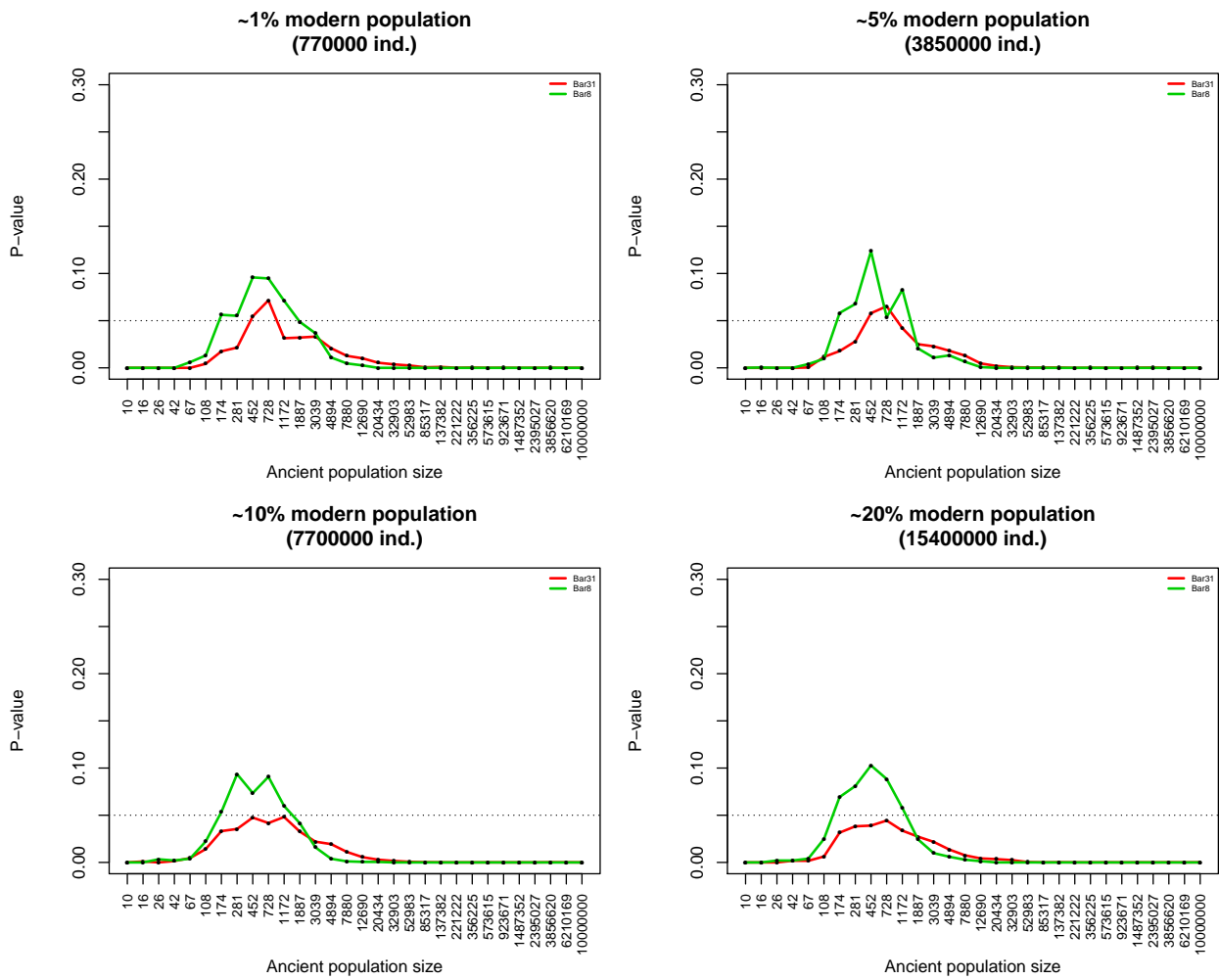
*Figure S17: Continuity grids for Rev5, Pal7 and Klei10, and modern Greeks.*



*Figure S18: Continuity grids for Bar31 and Bar8, and modern Turkish.*



**Figure S19:** Continuity grid slices for Rev5, Pal7 and Klei10 and modern Greeks on 1, 5, 10 and 20% of the population size in modern Greece. Dotted lines represent a p-value threshold of 0.05 above which continuity cannot be rejected.



**Figure S20:** Continuity grid slices for Bar31 and Bar8, and modern Turkish on 1, 5, 10 and 20% of the population size in modern Turkey. Dotted lines represent a p-value threshold of 0.05 above which continuity cannot be rejected.



# SI10 Comparing allele frequency patterns among samples using a mixture model

Lucy van Dorp\*, Saioa Lopez\*, Garrett Hellenthal

\* = contributed equally

## Introduction

The aim of this section is to use statistical models to represent modern and aDNA groups or individuals as mixtures of other sampled groups or individuals. These mixture patterns identify which sampled groups are most related to one another genetically, reflecting shared common ancestry relative to other groups due to e.g. admixture or other historical processes such as shared drift. Because of the low coverage of many aDNA genomes considered in this study, we treat SNPs as independent of each other (i.e. “unlinked”) in this section. A major advantage of our mixture model approach over commonly-used techniques to measure genetic structure, such as  $F_{ST}$  [146] or  $f$ -statistics [134] is that the DNA of each group/individual can be described as a mixture of the DNA of other groups/individuals, in contrast to being compared to only one or a small number of groups at a time.

## Methods

### Description of dataset analysed

We merged our new Neolithic samples from Greece (*Klei10, Pal7, Rev5*) and Anatolia (*Bar8, Bar31*) with the dataset from [41], which contained aDNA from ancient hominid groups Altai Neanderthal (*ALT*) and Denisova (*DEN*), a Mesolithic hunter-gatherer genome from Luxembourg (“Luxembourg\_Mesolithic” or Loschbour, *LOS*) and a Neolithic genome from Germany (“Stuttgart-LBK380”, *LBK*). Throughout we use the labels from [41] to refer to the groups in this dataset. We also merged this with Hungarian Neolithic (*NE1*) and Bronze Age (*BR2*) samples [81] taken from [102].

From this merged dataset, we excluded SNPs with genotype call rates  $<95\%$  using PLINK v1.9 [123]. We also removed 25 individuals (Href, Chimp, Gorilla, Orang, Macaque, Marmoset, Denisova\_light, Vindija\_light, Mez1, Otzi, Saqqaq, MA1, AG2, Skoglunk\_HG, Skoglund\_farmer, Motala\_merge, Motala12, Labrana, Bolivian81 (Bolivian\_Cochabamba), TGBS21 (Chane), Sesk\_47 (Chukchi\_Reindeer), BEL57 (Italian\_South), SD60\_297 (Saami\_WGA), tic\_95\_10 (Ticuna), wayu30 (Wayuu)) from the original dataset presented in [41] that either did not have heterozygous calls, were from groups containing only a single individual, or were not of interest in this study. After accounting for overlap and low coverage in the aDNA samples (see below), in total we analysed 2356 individuals at 35,188 total SNPs.

## Representing modern and aDNA calls

Let  $H_i^1, \dots, H_i^L$  be the observed data for chromosome  $i$  of a sampled individual, where  $H_j^i$  is the probability the given SNP is of allele type  $a_l$  at SNP  $l$ . (Here  $a_l$  is arbitrarily chosen to be one of the two possible allele types at SNP  $l$  according to the dataset from [41].) Here each individual from [41] and [102] has two haploid chromosomes with  $H_l^i \in \{0, 1\}$  for all non-missing data. In contrast, each new aDNA sample (i.e. *Rev5*, *Pal7*, *Klei10*, *Bar8*, *Bar31*) is represented as a single haploid chromosome with  $0 \leq H_l^i \leq 1$  based on the posterior probability of being a particular allele type according to our calling algorithm, subject to the modifications below.

We first removed SNPs where for any new aDNA sample:

- coverage is less than 2
- $Q_i > 0.001$ , where  $Q_i$  is equal to 1 minus the maximum posterior probability at SNP  $i$  across the four possible allele types ( $A, G, C, T$ )
- the allele type with maximum posterior probability does not match the two possible allele types at SNP  $i$  in the dataset from [41]

For SNPs passing these criteria, if  $Q_i = 0$  we set  $H_l^i = 1.0$  or  $H_l^i = 0.0$ , depending on which allele matching that in the “Lazaridis” dataset [41] had the highest posterior probability. For SNPs with  $Q_i > 0$ , we checked whether the allele type with the second highest posterior probability matched the other possible allele type according to the dataset from [41]. If so, we set  $H_l^i = 1.0 - Q_i$  (respectively  $H_l^i = Q_i$ , depending on which allele matching the dataset from [41] had highest posterior probability); otherwise we set  $H_l^i = 1.0$  (respectively  $H_l^i = 0.0$ ). Due to these strict criteria, note that  $H_l^i$  will nearly always be 0 or 1, or extremely close to these values, and can thus be thought of as such throughout the following.

## Inferring “allele matching profiles”

We followed the “unlinked” approach described in [147] (e.g. the “unlinked coancestry matrix”) to compare the alleles of a “recipient” haploid chromosome to that of a set of “donor” haploid chromosomes, while accounting for uncertainty in the calls. Again let  $H_i^1, \dots, H_i^L$  be the observed data for chromosome  $i$  of a sampled individual, where  $H_j^i$  is the probability the given SNP is of allele type  $a_l$  at SNP  $l$ , as described in the previous section. For recipient chromosome  $i$ , let  $X_l^i(d)$  be the score assigned to donor  $d \in [1, \dots, D]$  at SNP  $l$ , with  $D$  the total number of donor chromosomes and:

$$X_l^i(d) = \frac{H_l^i H_l^d + (1.0 - H_l^i)(1.0 - H_l^d)}{\sum_{j=1}^D H_l^i H_l^j + (1.0 - H_l^i)(1.0 - H_l^j)}. \quad (3)$$

(Donors  $d$  with missing values of  $H_l^d$  at SNP  $l$  are not allowed to contribute to recipient  $i$  at that SNP, which can diminish their overall score to recipient  $i$ , an issue we cope with using additional modeling in the next section.) We then calculate the total genome-wide allele matching score for recipient  $i$  and donor  $d$  as  $X^i(d) = \sum_{l=1}^L X_l^i(d)$ . We can then sum  $X^i(d)$  across donor chromosomes  $d$  that belong to the same “group” (e.g. population label, which we use in these analyses). Similarly, for recipient individuals with two haploid chromosomes, we can sum  $X^i(d)$  across these two haploid chromosomes to get a final vector of scores for that recipient individual, or we can sum  $X^i(d)$  across all  $i$  from a given recipient group to get a final vector of scores for that group.

In the end, assuming we partition our  $D$  donors into  $K$  groups (here each of the  $K$  groups refers to a distinct population label), we define the “allele matching profile” for recipient  $r$  as  $f^r \equiv \{f_1^r, \dots, f_K^r\}$ , with:

$$f_k^r = \frac{\sum_{d=1}^D 1_{[d,k]} X^r(d)}{\sum_{j=1}^K [\sum_{d=1}^D 1_{[d,j]} X^r(d)]}, \quad (4)$$

where  $1_{[d,k]} = 1$  if donor  $d$  is assigned to group  $k$  and 0 otherwise. Note that  $\sum_{k=1}^K f_k^r = 1.0$ . As noted above,  $r$  can represent a single haploid chromosome, a single individual, or all haploid chromosomes from a common group.

We used the following two sets of donor and recipient groups, with “ancients” referring to  $\{Klei10, Pal7, Rev5, Bar8, Bar31, LBK, LOS, BR2, NE1\}$  and “moderns” referring to all other groups in the merged dataset (including  $ALT, DEN$ ):

- A. modern groups are used as donors; modern groups and ancients are used as recipients
- B. same as (A), but exclude the following modern groups as donors: Adygei, Armenian, Bulgarian, Cypriot, Georgian, Greek, Hungarian, Palestinian, Syrian, Turkish

Analysis (B) disallows the groups listed above to copy from their close neighbors, which might mask interesting ancestry signals common to all (or a subset of) these groups. In each of (A) and (B), any recipient individual cannot use their own haploid chromosome data as a “donor”. For this reason, when constructing our allele matching profiles we used a “leave-one-out” approach analogous to that described in [121]. In particular if each donor group  $\{1, \dots, K\}$  contains  $\{n_1, \dots, n_K\}$  individuals, respectively, with  $N = \sum_{k=1}^K n_k$  total donor individuals, we fix the set of donors to contain  $n_k - 1$  individuals from each of the  $K$  groups (i.e. giving  $N - K$  donor individuals in total). The reasoning behind this “leave-one-out” approach is to make the final “allele matching profiles”  $f^r$  comparable across all recipient groups (see next section). For example, under analysis (A) each recipient Greek individual can only use  $n_{\text{Greek}} - 1$  other Greek individuals as donors, and therefore we fix every other recipient individual to use only  $n_{\text{Greek}} - 1$  Greek individuals as donors. Exceptions to this “leave-one-out” rule are the two ancient hominid genomes  $ALT$  and  $DEN$ , for which we have only a single

sample of each; for this reason these genomes are used as donors for every group and otherwise are *not* used as contributors to the mixture model described in the next section.

### Inferring final “proportions of ancestry” based on the “allele matching profiles”

Our inferred “allele matching profiles” suffer some limitations. For example *a priori* a donor group  $d$  with a disproportionately large number of sampled individuals (or lower amount of missingness; see above) may have relatively higher values of  $f_d^r$  across all recipient groups  $r$ , potentially leading to a biased interpretation of results. To cope with this, we use additional mixture modeling described in this section to “clean” the raw  $f^r$  inference, as in [121, 148, 149]. We can also use this technique to compare a group’s “allele matching profiles” to that of any other groups we include in the mixture.

As before (following similar notation to that in [121, 149]), let  $f^r \equiv \{f_1^r, \dots, f_K^r\}$  be the observed “allele matching profile” inferred using equation (4) for recipient group  $r$ . Note that for analyses (A) and (B) described in the previous section, we have analogous allele matching profiles for all other recipient groups  $j \neq r \in [1, \dots, R]$ . To measure the relative amount of drift (or “self-copying”) in group  $r$ , we introduce a  $K$ -vector  $f^{r*}$  with  $f_r^{r*} = f_r^r$  and all other entries 0. We “clean” the painting of group  $r$  using the following linear model:

$$f^r = \left[ \sum_{s=1}^S \beta_s^r f^s \right] + \beta_{\text{SELF}}^r f^{r*} + \epsilon, \quad (5)$$

where  $s = 1, \dots, S$  represents a set of “surrogate” groups used to describe the ancestry of group  $r$ . Specifically, the set of surrogates can contain all other  $R - 1$  recipient groups, or it may contain any subset of these  $R - 1$  total groups. We explore several combinations of surrogates below. Here  $\epsilon$  is a vector of errors, and we seek the estimates  $(\hat{\beta}_1^r, \dots, \hat{\beta}_S^r, \hat{\beta}_{\text{SELF}}^r)$  to replace  $(\beta_1^r, \dots, \beta_S^r, \beta_{\text{SELF}}^r)$ , respectively, in equation (5) that minimize  $\epsilon$  using least-squares. (Note that  $f_r^{r*} = 0$  if recipient group  $r$  was not included among the  $K$  donor groups, so that  $\beta_{\text{SELF}}^r = 0$  in these cases.) We use the non-negative-least-squares “nls” package in R to estimate the  $\beta$ s under the constraints that each  $\hat{\beta}_s^r \geq 0$ ,  $\hat{\beta}_{\text{SELF}}^r \geq 0$ , and  $(\hat{\beta}_{\text{SELF}}^r + \sum_{s=1}^S \hat{\beta}_s^r) = 1.0$ . We refer to the set of  $\{\hat{\beta}_1^r, \dots, \hat{\beta}_S^r, \hat{\beta}_{\text{SELF}}^r\}$  values as our inferred “proportions of ancestry” for group  $r$  conditional on this set of surrogates.

To measure uncertainty in the  $\hat{\beta}$ s, as in [149] we take an approach analogous to [134] and calculate standard errors using a weighted Block Jackknife [150] approach that removes each chromosome one-at-a-time (from all donor and recipient groups) and re-calculates the  $\hat{\beta}$ s, weighting each jackknife sample by that chromosome’s number of SNPs. An alternative approach to measure uncertainty is used in [148], who instead used bootstrap re-samples of individuals’ chromosomes. In contrast to that approach, the jackknife technique we use here is applicable when evaluating uncertainty in a single genome.

For each allele-matching analysis (A) and (B), we performed the following four mixture model analyses (though here “modern” groups exclude *ALT, DEN*, who are not used as surrogates for reasons described above):

- (I) “all moderns” – form each ancient and modern genome using all modern groups as surrogates
- (II) “all moderns + ancients” – form each ancient and modern genome using all modern+ancient groups as surrogates
- (III) “ancients + Yoruba” – form each ancient and modern genome using all other ancient genomes, plus the modern Yoruba, as surrogates
- (IV) “ancients (excluding *BR*) + Yoruba” – form each ancient and modern group using the modern Yoruba and all other ancient genomes except *BR2* as surrogates

In each case, a group cannot use itself as a surrogate or else it would match itself exactly. Under allele-matching analysis (B), the same groups we disallow as donors are also disallowed as surrogates for mixture model analyses (I) and (II). For analyses (III) and (IV), we were interested in how modern and ancient groups relate ancestrally to different sets of ancient genomes. We also included the Yoruba as a surrogate in (III) and (IV), since our ancient samples contain no proxies for sub-Saharan Africa and e.g. several West Eurasian groups we use here have been shown to have recent African admixture [121].

For analyses (I) and (II), if the final inference included more than ten surrogate groups with  $\hat{\beta}_s^r > 0$ , we did an altered procedure to mitigate effects of over-fitting. In particular we sequentially included surrogates that improved the total variation distance (TVD) measure (e.g. used in [148]) between  $\hat{f}^r$ , the inferred allele matching profile of recipient group  $r$  based on the inferred best fit to equation (5), and  $f^r$ , the actual allele matching profile of recipient group  $r$ . To do so, we measure TVD comparing two profiles  $x, y$  using:

$$TVD(f^x, f^y) = 0.5 \sum_{k=1}^K |f_k^x - f_k^y|, \quad (6)$$

and we performed the following procedure:

1. Calculate TVD between each surrogate group and the recipient group  $r$ .
2. Take the ten surrogates with the lowest TVD scores from step 1.
3. From the group of surrogates not among those selected in step 2, sequentially add each surrogate, one-at-a-time. Find the added surrogate among these for which  $TVD(f^r, \hat{f}^r(s))$  is

smallest, where  $\hat{f}^r(s)$  is the inferred best fit to equation (5) when including newly added surrogate  $s$  and the surrogates already included from step 2. Add the surrogate with lowest such  $TVD(f^r, \hat{f}^r(s))$  to the list of surrogates created in step 2.

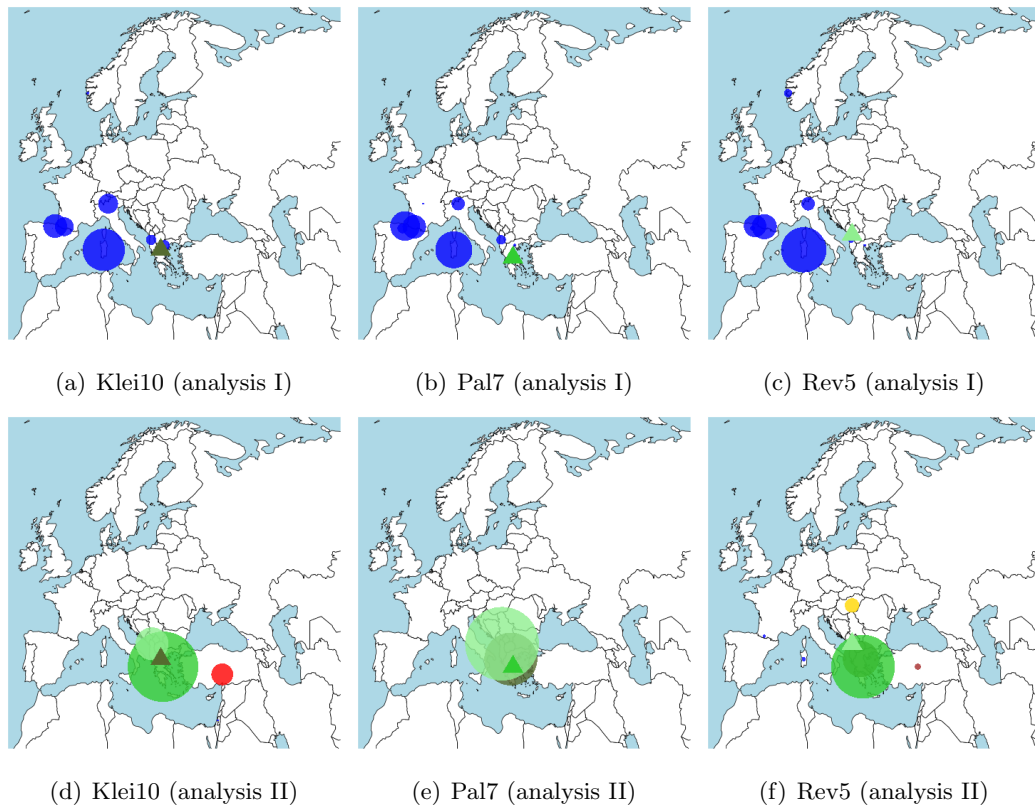
4. Repeat step 3 until the proportion change in  $\min_s[TVD(f^r, \hat{f}^r(s))]$  relative to the previous iteration is  $< 0.2$ . I.e. repeat step 3 until, when adding each of the not-yet-included surrogates (call this set of surrogates  $\Omega$ ) one at a time, there is always a less than 20% reduction in fitted TVD relative to the previous iteration that did not include any surrogates from  $\Omega$ . We fixed these final surrogates when performing jack-knifing to get standard errors around the inferred  $\beta$  values.

For analyses (III) and (IV), we also used a slightly alternative version of equation (5) that matches that used in [121]. In particular when inferring coefficients for each group  $r$ , we set  $f_r^r = 0$  (and rescaled so that  $\sum_{k=1}^K f_k^r = 1.0$ ),  $f_r^s = 0$  for all  $s \in [1, \dots, S]$  (similarly rescaling each  $f^s$  vector to sum to 1.0), and  $\beta_{\text{SELF}}^r = 0$ , i.e. we disregarded any “self-copying” in group  $r$ , which gave more consistent results with such a limited set of surrogate groups.

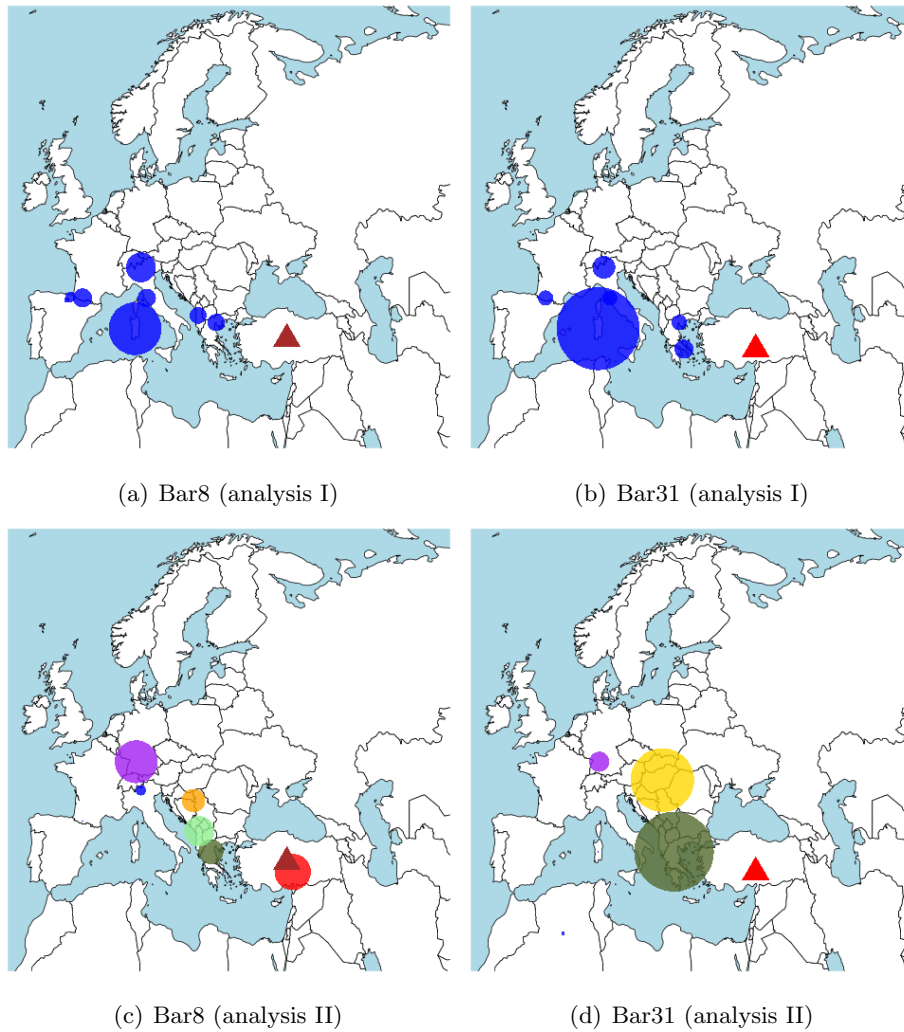
## Results

### Inferred proportions of ancestry for aDNA samples

Our inferred proportions of ancestry under allele-matching analysis (A) are provided for all ancient samples (i.e.  $\{Klei10, Pal7, Rev5, Bar8, Bar31, LBK, LOS, BR2, NE1\}$ ) for each of analyses (I) and (II) in Figures S21-S24, with summarized analysis (II) results in Figure S30, and for analyses (III) and (IV) in Figure S29. Table S30 highlights the proportion of contributions from populations of Southern Europe, Northern Europe, the Levant and Caucasus for analysis (I). Estimates for all analyses (I)-(IV) and allele-matching analyses (A)-(B), along with jackknife-based standard errors, are provided in Dataset S3. For the figures here, we provide results only under allele-matching analysis (A), though note that results were very similar (as expected) under allele-matching analysis (B) (see Dataset S3). In all figures the positioning of the ancient samples is indicative so that the individual samples can be more easily observed in each scenario.

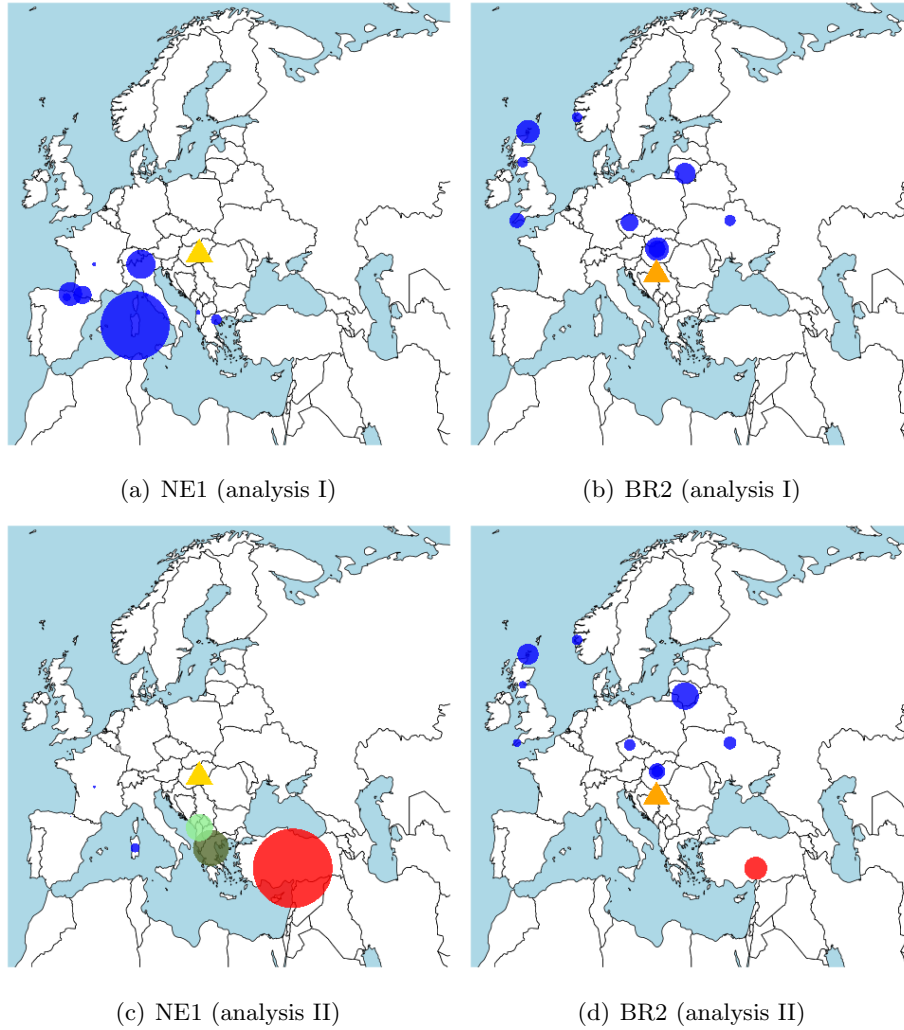


**Figure S21:** *Inferred proportions of ancestry for each Neolithic Greek sample when using all moderns (top row; analysis I) versus all modern+ancient groups (bottom row; analysis II) as surrogates under allele-matching analysis (A). Circles are proportional to the inferred proportions from modern samples (blue) and aDNA samples from Greece (green), Anatolia (red), Hungary (BR2, orange; NE1, yellow), Stuttgart (LBK; purple) and Luxembourg (LOS; black). For regions with multiple aDNA samples (i.e. Greece, Anatolia, Hungary), colors for the samples are darker the younger the sample. Triangles represent the sampling location of the depicted target sample (and also provide the key for that sample's color, as also provided in Figure S29).*

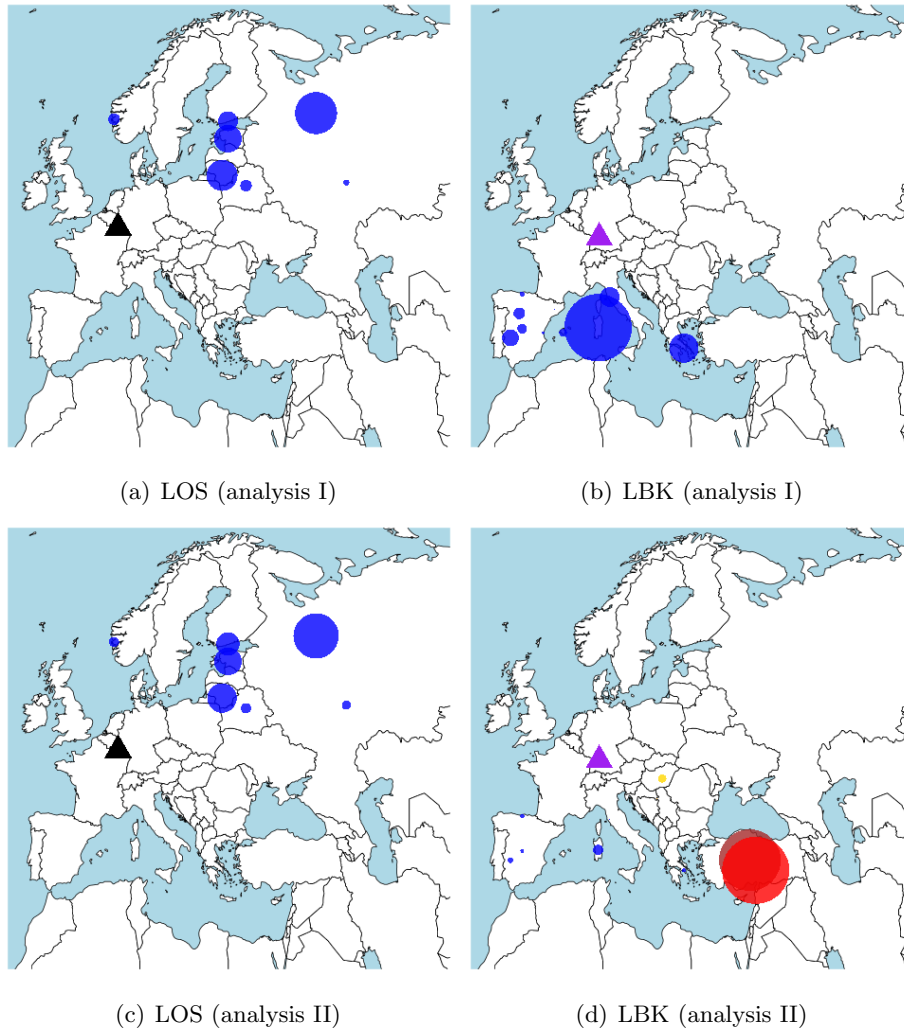


**Figure S22:** *Inferred proportions of ancestry for each Neolithic Anatolia sample when using all moderns (top row; analysis I) versus all modern+ancient groups (bottom row; analysis II) as surrogates under allele-matching analysis (A). See caption to Figure S21 for legend.*





**Figure S23:** *Inferred proportions of ancestry for Hungarian Neolithic (NE1) and Bronze Age (BR2) samples from [81] when using all moderns (top row; analysis I) versus all modern+ancient groups (bottom row; analysis II) as surrogates under allele-matching analysis (A). See caption to Figure S21 for legend.*



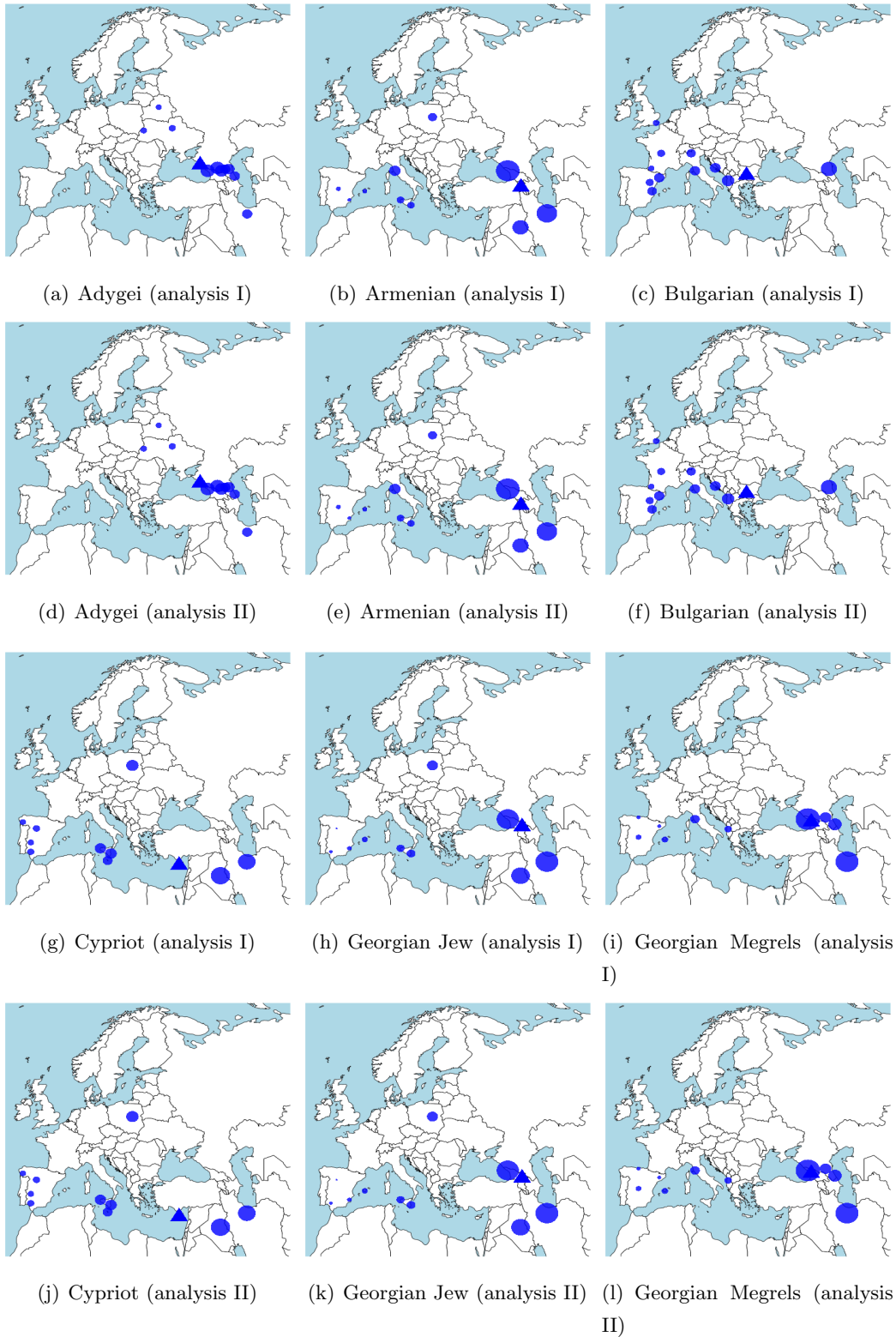
**Figure S24:** *Inferred proportions of ancestry for Luxembourg hunter-gatherer (LOS) and Germany Neolithic (LBK) aDNA samples from [41] when using all moderns (top row; analysis I) versus all modern+ancient groups (bottom row; analysis II) as surrogates under allele-matching analysis (A). See caption to Figure S21 for legend.*

**Table S30:** *Inferred proportions of ancestry under allele-matching analysis (A) for all the ancient samples when using all the modern populations as surrogates (analysis I). <sup>1</sup>=Italy, Sardinia, Spain, South France; <sup>2</sup>=England, Scotland, Finland, Norway, Orkney Islands, Lithuania, Estonia, Czechoslovakia, Iceland, Poland; <sup>3</sup>=Israel, Cyprus, Jordan, Lebanon, Syria, Turkey; <sup>4</sup>= Armenia, East Ukrania, Georgia.*

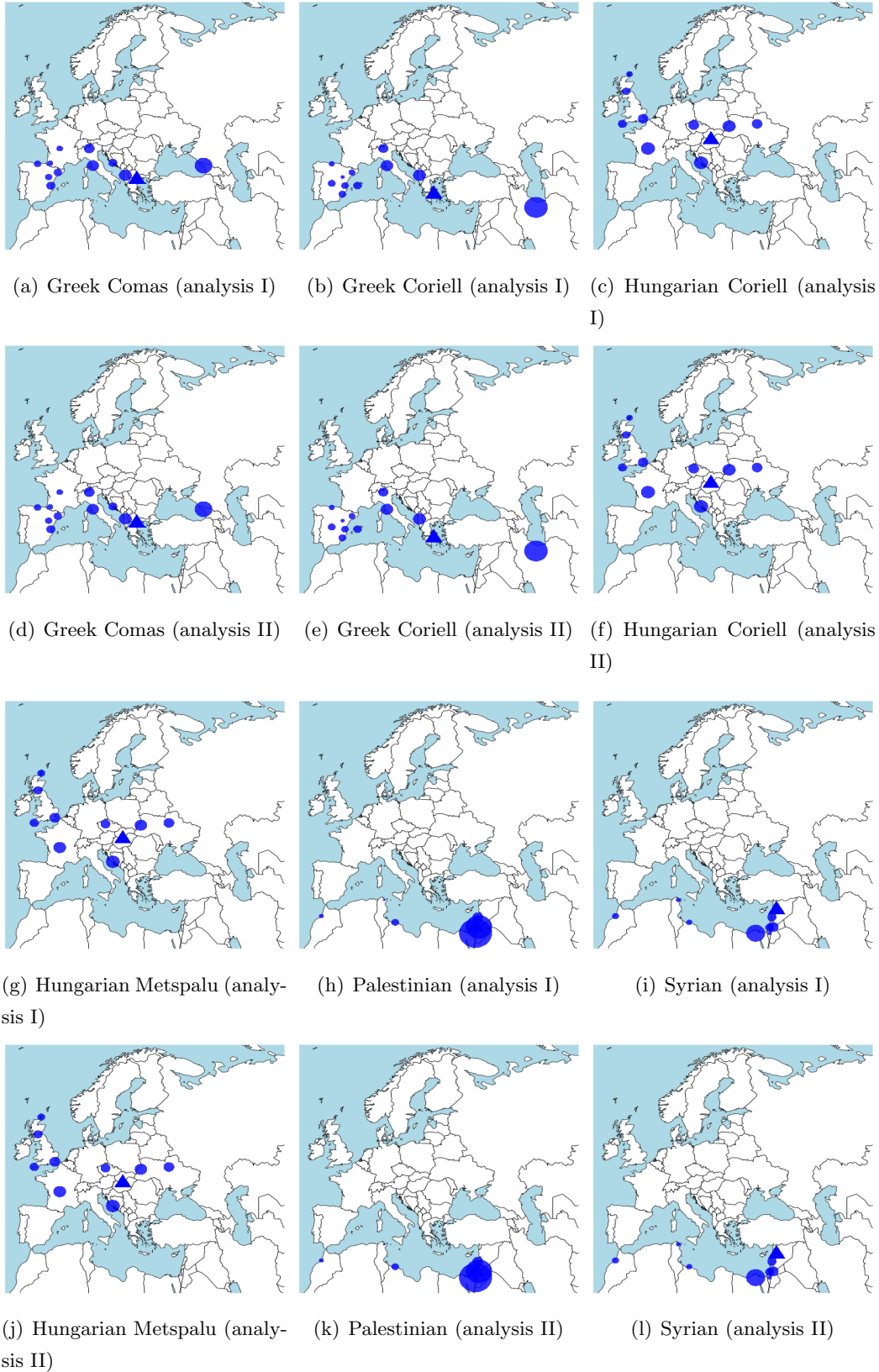
	Loschbour	LBK	Bar31	Bar8	Klei10	Pal7	Rev5	BR2	NE1
Southern Europe <sup>1</sup>	0.00	71.14	71.05	68.98	81.92	88.67	86.71	0.00	86.89
Sardinia	0.00	39.54	48.70	31.00	32.97	28.84	35.36	0.00	40.50
Northern/Central Europe <sup>2</sup>	63.04	0.00	0.00	0.00	2.20	0.00	6.41	69.48	0.00
Levant <sup>3</sup>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Caucasus <sup>4</sup>	0.00	0.00	0.00	0.00	0.00	0.00	0.00	6.66	0.00
Others	36.96	28.86	28.95	31.02	15.88	11.33	6.88	23.86	13.11
Greece	0.00	17.19	20.15	10.23	7.78	2.37	2.09	0.00	6.23

### Inferred proportions of ancestry for modern samples

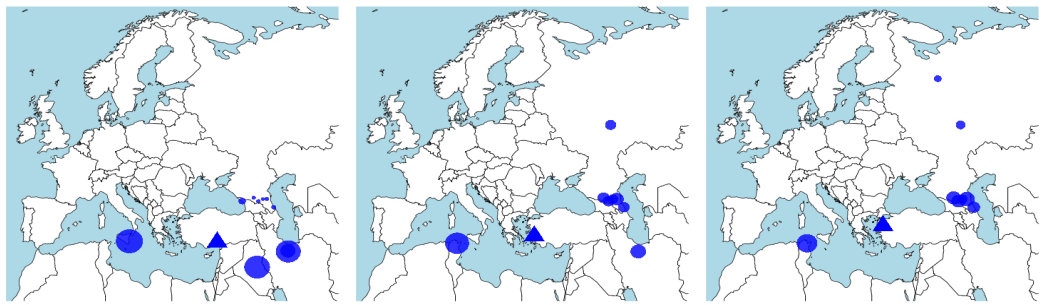
Our inferred proportions of ancestry under allele-matching analysis (B) are provided for the modern groups {Adygei, Armenian, Bulgarian, Cypriot, Georgian, Greek, Hungarian, Palestinian, Syrian, Turkish} for each of analyses (I) and (II) in Figures S25-S28, with summarized analysis (II) results for these and additional modern groups in Figure S30, and under allele-matching analysis (A) for analyses (III) and (IV) for these and additional modern groups in Figure S29. Estimates for all modern groups for all analyses (I)-(IV) and allele-matching analyses (A)-(B), along with jackknife-based standard errors, are provided in Dataset S3.



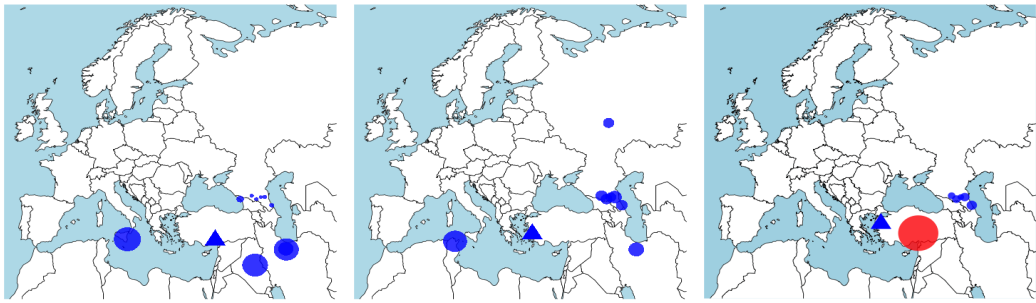
**Figure S25:** *Inferred proportions of ancestry for modern groups when using modern (first and third rows; analysis I) versus modern+ancient groups (second and fourth rows; analysis II) as surrogates under allele-matching analysis (B). See caption to Figure S21 for legend.*



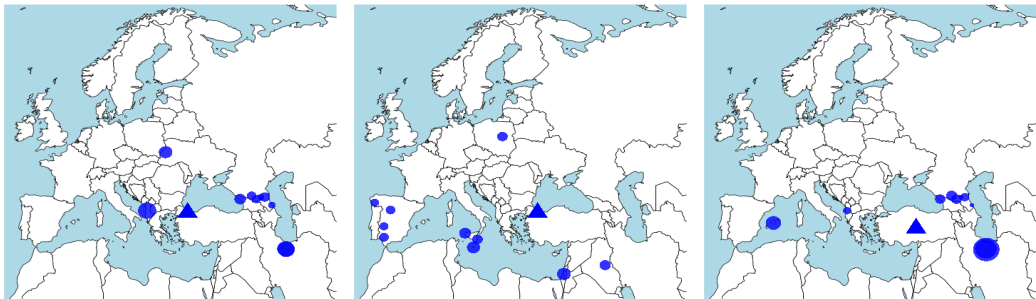
**Figure S26:** *Inferred proportions of ancestry for modern groups when using modern (first rows; analysis I) versus modern+ancient groups (second and fourth rows; analysis II) as surrogates under allele-matching analysis (B). See caption to Figure S21 for legend.*



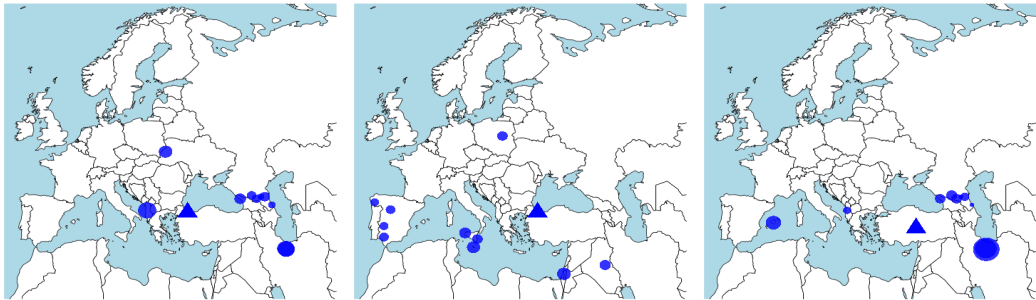
(a) Turkish Adana (analysis I) (b) Turkish Aydın (analysis I) (c) Turkish Balıkesir (analysis I)



(d) Turkish Adana (analysis II) (e) Turkish Aydın (analysis II) (f) Turkish Balıkesir (analysis II)

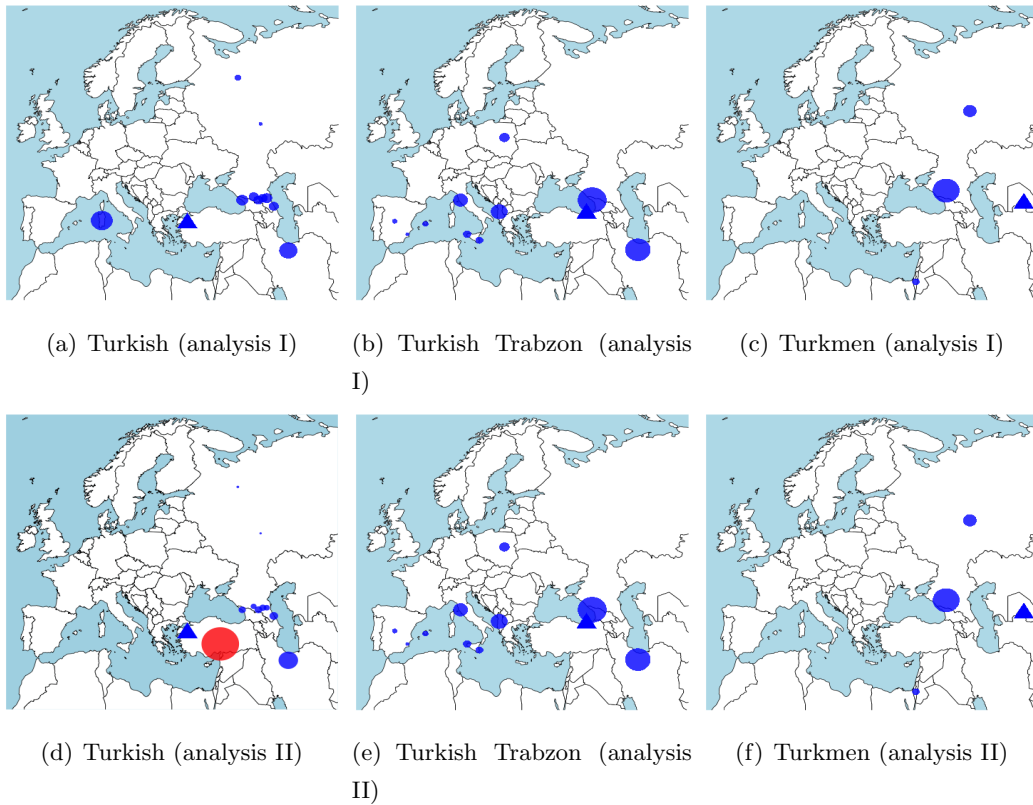


(g) Turkish İstanbul (analysis I) (h) Turkish Jew (analysis I) (i) Turkish Kayseri (analysis I)

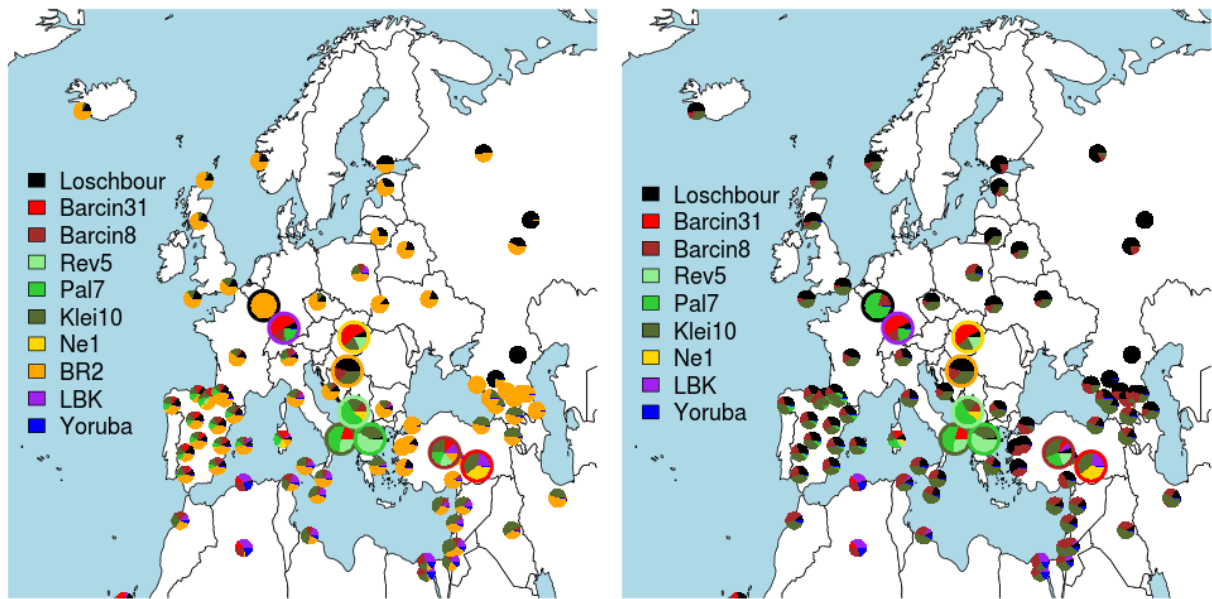


(j) Turkish İstanbul (analysis II) (k) Turkish Jew (analysis II) (l) Turkish Kayseri (analysis II)

**Figure S27:** *Inferred proportions of ancestry for modern groups when using modern (first and third rows; analysis I) versus modern+ancient groups (second and fourth rows; analysis II) as surrogates under allele-matching analysis (B). See caption to Figure S21 for legend.*



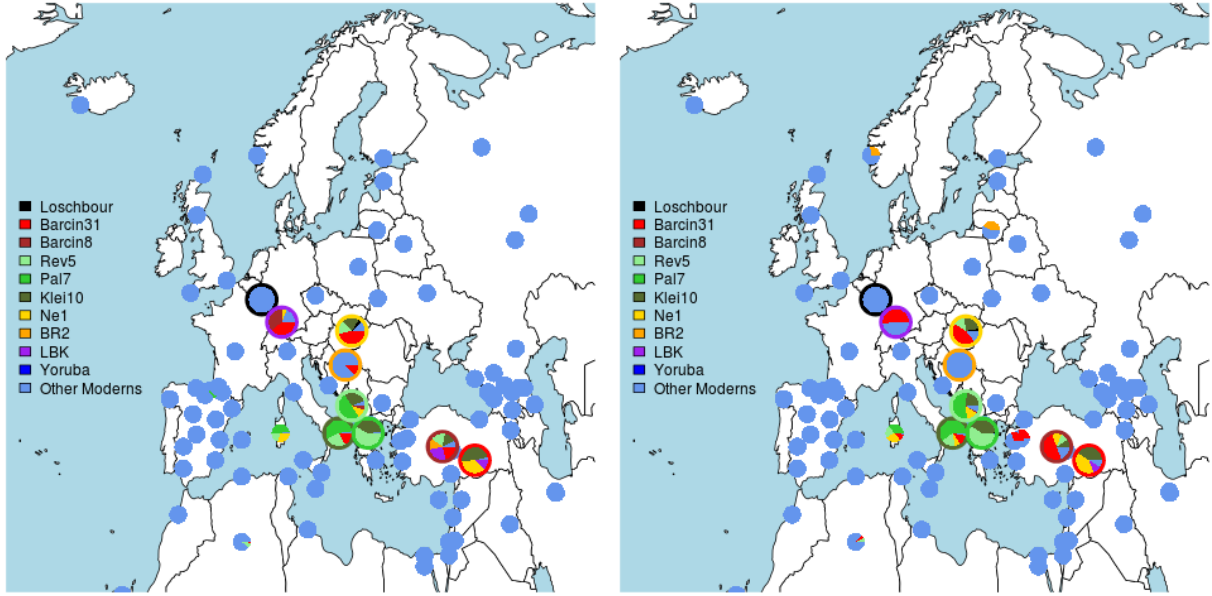
**Figure S28:** Inferred proportions of ancestry for modern groups when using modern (first row; analysis I) versus modern+ancient groups (second row; analysis II) as surrogates under allele-matching analysis (B). See caption to Figure S21 for legend.



(a) all ancients and Yoruba as surrogates (analysis III) (b) all ancients except *BR2* and Yoruba as surrogates (analysis IV)

**Figure S29:** Inferred proportions of ancestry for modern and ancient groups when using all ancients and the modern Yoruba samples (analyses III-IV) as surrogates under allele-matching analysis (A). For (b), *BR2* is excluded as a surrogate. The larger pie charts represent the aDNA samples, with borders corresponding to the legend at left. In both (a) and (b) samples are ordered in the legend according to a mix of age and geography.





(a) ancients and moderns as surrogates (analysis AII)    (b) ancients and moderns as surrogates (analysis BII)

**Figure S30:** Inferred proportions of ancestry for modern and ancient groups when using modern groups and all ancients as surrogates for analysis AII (A) and analysis BII (B). Contributions from Yoruba are shown in dark blue; contributions from all other modern surrogates are shown in light blue. The larger pie charts represent the aDNA samples, with borders corresponding to the legend at left. In both (a) and (b) samples are ordered in the legend according to a mix of age and geography.

## Discussion

We note the following observations about our analysis of the ancient samples

$\{Klei10, Pal7, Rev5, Bar31, Bar8, LBK, LOS, NE1, BR2\}$ :

1. When inferring proportions of ancestry for any given ancient sample  $r$ , only modern groups geographically near to where  $r$  was sampled typically give  $\hat{\beta}_s^r > 0$ , with the vast majority of the  $>200$  modern groups  $s$  not contributing at all (Figures S21-S24, top row). This encouragingly suggests that our aDNA calls appear relatively stable. As seen previously in early farmer samples [41], modern-day Sardinians give the highest inferred proportion for all ancient samples, except the hunter-gatherer  $LOS$  and the Bronze Age sample  $BR2$ .
2. With the exceptions of  $LOS$  and to a lesser extent  $BR2$ , when inferring each ancient sample's proportions of ancestry, incorporating the other ancient samples as surrogates drastically reduces the contributions from the modern groups (Figures S21-S24, compare the bottom rows to the top rows).
3. The Greek samples appear to be genetically similar to each other, each inferring their largest contributions from the other Greek samples with the sample closest in age favoured (Figures S21a-f and Figure S29). The youngest aDNA Greek genome  $Klei10$  receives contributions from the Anatolian genome  $Bar31$  (Figure S21d). The oldest Greek genome  $Rev5$  also receives a

contribution from Anatolia, this time from the younger Anatolian genome *Bar8* but also from the Neolithic Hungarian *NE1* (Figure S21f).

4. The older Anatolia genome *Bar31* looks genetically similar to the Greek genome *Klei10*, the Neolithic Hungarian *NE1* and – to a lesser extent – the Neolithic German *LBK* (Figure S22d, Figures S29-S30). In contrast, the younger Anatolia genome *Bar8* (roughly) appears similar to the Aegean Neolithic genomes *Bar31*, *Klei10*, *Rev5*, the Germany Neolithic *LBK* and – to a lesser extent – the Hungarian Bronze Age sample *BR2* (Figure S22c, Figures S29-S30).
5. Interestingly, the Neolithic Hungarian genome *NE1* from [81] looks most genetically similar to the older Anatolia genome *Bar31*, with notable contributions from the Neolithic Greek genomes *Rev5* and *Klei10* and little contribution from any other Neolithic surrogates (Figure S23c, Figures S29-S30). In contrast, under analysis (II) the Bronze Age Hungarian genome *BR2* from [81] looks genetically like a mixture of neighbouring modern European groups, with relatively little contributions from any Neolithic surrogates except for the oldest Anatolian genome *Bar31*, Figure S23d.
6. As reported previously [41], the hunter-gatherer genome *LOS* from [41] looks genetically most similar to modern groups from eastern Europe, with no substantial contributions from any aDNA samples under analysis (II) (Figure S24a and c, Figure S30). *LOS* contributes more to the Germany *LBK* and Hungary *NE1* genomes out of all Neolithic samples under analyses (III) and (IV) (Figure S29). As a hunter-gatherer, *LOS* likely looks most like the Bronze Age sample *BR2* and Neolithic Greek genome *Pal7* under these analyses because there are no other Neolithic hunter-gatherers in the dataset (Figure S29).
7. Interestingly, the Germany Neolithic genome *LBK* from [41] looks genetically very similar to both Anatolian aDNA genomes, and particularly the older genome *Bar31*, with only a small or no contribution from the geographically closer Neolithic genome *NE1* from Hungary (Figure S24d, Figures S29-S30).
8. In nearly every analysis the Germany Neolithic genome *LBK* contributes relatively less to the Neolithic Aegean genomes than the Neolithic Aegeans contribute to *LBK*, particularly comparing to the Anatolians (e.g. Figures S21, S22, S24, S29, S30). This observation is consistent with founder effects in the Germany Neolithic sample, after deriving from a source genetically similar to the Anatolia or Greek Neolithic samples. However, we caution against using this as strong evidence, as which populations are included in the mixture will influence results. For example here we have few ancient Neolithic samples from locations proximal to *LBK*, while including many geographically neighbouring aDNA samples from Anatolia and Greece.
9. Similarly, the Hungary Neolithic genome *NE1* contributes relatively less to the Neolithic Aegean (Greek and Anatolia) genomes than the Neolithic Aegeans contribute to *NE1* (e.g. Figures S21, S22, S23, S29, S30), with the exception of *Bar8* in analysis (BII). However, this effect is much less pronounced than the similar comparison between *LBK* and the Anatolians.

Interestingly, *LBK* and *NE1* do not contribute to each other under any analysis except a  $< 5\%$  contribution from *NE1* to *LBK* under analysis (AII), plausibly suggesting any founder effect influencing *LBK* did not as strongly affect *NE1*. At the very least, this suggests a closer relationship between the Aegean genomes and *NE1* relative to *LBK*.

For our comparisons of the modern groups to the ancient samples, we note the following observations:

1. As expected, modern groups (e.g. Adygei, Armenian, Bulgarian, Cypriot, Georgian, Greek, Hungarian, Palestinian, Syrian, Turkish) typically receive small contributions from aDNA samples, relative to the total contribution from other modern groups, when including both moderns and ancients as surrogates (Figure S25, S26, S27, S28, S30).
2. An exception to this are two Turkish groups (Turkish, Turkish\_Balikesir), which each receive a substantial contribution ( $\approx 30\%$ ) from the Neolithic *Bar31* sampled in modern-day Turkey. Notably, out of all modern groups Sardinians receive by far the highest genetic contributions ( $> 97\%$ ) from Neolithic genomes under analysis (II), suggesting a strong genetic affinity between modern-day Sardinians and these ancient samples even relative to any other modern-day groups (Figure S30).
3. When analysing the modern groups using only the aDNA genomes (plus the Yoruba) as surrogates while excluding the Bronze Age sample *BR2* (analyses IV, Figure S29b), most moderns are inferred to be mixtures of the hunter-gatherer genome *LOS* and the youngest (generally higher coverage) Neolithic genomes, consistent with previous findings [41].
4. Under analysis (IV), the youngest Anatolian Neolithic sample *Bar8* and youngest Greek Neolithic sample *Klei10* are the highest contributing surrogates to most modern groups in Europe and surrounding the Mediterranean, with the interesting exception of some North African modern groups that infer relatively higher contributions from the Germany Neolithic *LBK* (Figure S29b). These observations are consistent if we include only *Bar8* or only *Klei10* out of all Aegean surrogates in analysis (IV). This preferential matching of modern groups to the Aegean Neolithic samples over both *LBK* and *NE1* should not reflect any bias in genotype calling protocol, given that the Aegean genotypes were called separately from publicly available data that included *LBK*, *NE1* and all moderns [41, 102]. Instead this provides additional evidence that the Aegean aDNA samples better represent the ancestors of modern groups in the region, as *LBK* (and perhaps *NE1*) were subjected to founder (or other) effects that make them less representative.

## SI11 Runs of Homozygosity

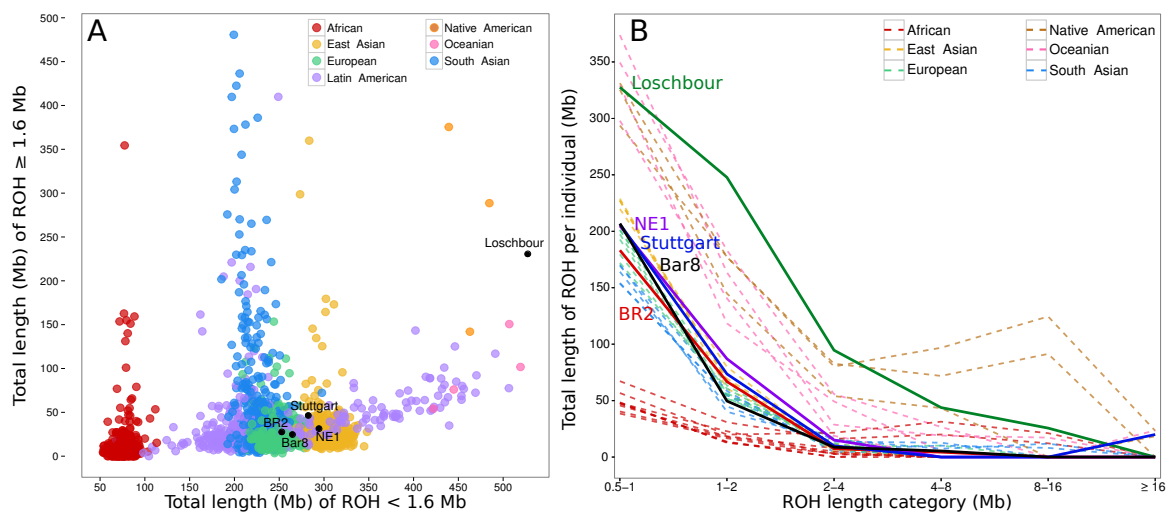
Lara M Cassidy & Daniel G Bradley

One sample from this study, Bar8, had sufficient genomic depth of coverage to analyse the proportion of the genome under runs of homozygosity (ROH). The distribution of ROH in the genome is informative of past population demography [151], [152]. Long stretches of homozygosity indicate recent endogamy in an individual's ancestry while an excess of shorter runs across the genome are a result of more ancient bottlenecks in the population's past. ROH analysis on Bar8 was carried out alongside four other high coverage ancient European samples, the Mesolithic individual Loschbour [41], two Early Neolithics, NE1 from Hungary [81] and Stuttgart, a LBK farmer from Germany [41], and BR2, a Bronze Age Individual from Hungary [81]. Also included were the 2504 modern individuals from the Phase 3 v5 1000 Genomes release [153] and 23 individuals from the Simons Genome Diversity Project [154].

Diploid genotype calls (SI5) from the five ancient genomes were filtered for a depth of coverage of 10X or above and a genotype quality of 30 or above. Only biallelic autosomal transversion SNPs were considered. These SNPs also required a minor allele frequency above 1% in at least one of the five 1000 Genomes super population groupings. The filtered calls from the ancient samples were subsequently merged with the same set of SNPs from 2504 individuals of the Phase 3 v5 1000 Genomes release [153] and 23 individuals from the Simons Genome Diversity Project [154] using PLINK v1.90 [124]. Only SNPs called securely across all individuals were retained, leaving 1,447,024 transversion SNPs for ROH analysis.

ROH were identified using PLINK v1.90 [124] with the specifications used in [81] (-homozyg, -homozyg-density 50, -homozyg-gap 100, -homozyg-kb 500, -homozyg-snp 50, -homozyg-window-het 1, -homozyg-window-snp 50 and -homozyg-window-threshold 0.05). ROH in each individual were divided into two classes, long ROH (>1.6 Mb), indicative of recent endogamy, and short to intermediate (<1.6 Mb), a result of more ancient population bottlenecks [152]. The summed total length of ROH for both classes was then calculated for each individual. These values are plotted against each other in Figure S31.A.

To investigate the length distribution of ROH in more detail we selected single individuals from nineteen 1000 Genomes populations, not known to have experienced recent historical admixture. These individuals displayed the approximate median value of their population for total fraction of the genome under ROH. Thirteen individuals from the Simons Genome Diversity Project were also selected for analysis. For each individual ROH were placed into various size bins and the total length of ROH in each bin was calculated. Results are displayed in Figure S31.B.



**Figure S31:** Estimated ROH distributions for high coverage ancient and modern individuals. ROH distribution was determined using 1,447,024 transversion SNPs called securely in all samples. Modern individuals are colored according to geographical region. (A) displays the total combined lengths of short to intermediate ROH ( $< 1.6$  Mb) plotted against those of long ROH ( $> 1.6$  Mb) for 2,527 modern and 5 ancient individuals. (B) shows the total length of ROH in 32 modern and 5 ancient individuals for a series of length categories

## SI12 Functional Markers

Karola Kirsanow

We assessed the five individuals for whom nuclear genomic data was available (Klei10, Pal7, Rev5, Bar8, and Bar31) for diploid genotype at a panel of single nucleotide polymorphisms (SNPs) having known functional associations in modern populations, many of which have been identified as targets of selection in modern and ancient human populations. Genotypes were determined using the diploid genotyping method described in SI5, and further verified through direct observation of BAM files using samtools tview [95]. We included sites having  $\geq 2X$  coverage in the analysis. We remind the reader that phenotypic inference based on single-SNP genotypes in low-coverage data can be precarious; contextualizing low-coverage genotypes into haplotypes or network models can ameliorate but not eliminate this problem.

### Pigmentation

Rev5, Klei10, Bar8, and Bar 31 were all observed to carry homozygous derived alleles at *SLC24A5* rs1426654 (A111T); this SNP, which is associated with skin depigmentation, is nearly fixed in modern Europeans (1000 Genomes EUR frequency  $> 0.99$  [155]), including modern Greek ( $> 0.99$ ) and Turkish ( $> 0.98$ ) populations [156], and known to be under strong natural selection [130, 157]. The Klei10, Bar8, and Bar31 individuals carried the 16-SNP founder haplotype on which the *SLC24A5* A111T is believed to have arisen [158]; Pal7, in which the core *SLC24A5* A111T mutation could not be genotyped, may have also carried this haplotype (see Table S31). The haplotype of the Rev5 individual could not be identified due to insufficient coverage.

Derived alleles of the *SLC45A2* rs16891982 L374F mutation, also associated with skin depigmentation and identified as under selection in modern and ancient Europeans [157, 159], were observed in the Klei10, Pal7, Bar31, and Bar8 individuals. This allele is nearing fixation in modern Europeans (1000 Genomes EUR frequency 0.94 [160]), and is at high frequency in modern Greek (0.86) and Turkish (0.68) populations [156]. The Klei10, Bar8, and Bar31 genotypes at a panel of *SLC45A2* markers are compatible with the inference that these individuals carried at least one copy of the putative *SLC45A2*, *L374F* founder haplotype [161], however, the haplotype determination is inconclusive due to incomplete and unphased data.

The Klei10 individual appears to have been heterozygous for the derived allele at *HERC2* rs12913832, a mutation associated with iris depigmentation and having evidence for differential selection at different points in prehistory [130, 159]. Qualified support for the presence of a derived allele in the Klei10 individual comes from the observation of a derived allele at rs1129038 which is almost completely linked with the causal rs12913832 SNP in modern populations [162]. However, conclusive haplotype determinations for the 13-SNP *OCA2-HERC2* profile found in almost all modern blue-eyed individuals could not be made for the ancient Aegean individuals [163] (see Table S31).

In order to compare variation in pigmentation-associated loci in the ancient Aegean individuals with that observed in modern Europeans, we typed each individual at a panel of sites developed for forensic eye and hair pigmentation phenotype prediction (Hirisplex [164]). Coverage was sufficient to enable eye and hair reconstructions for the Klei10, Pal7, Bar8, and Bar31 individuals. All individuals were reconstructed as having a higher probability of non-blue eyes, however the Bar8 and Bar31 Anatolians were more likely to have had dark hair of an indeterminate color, while the Klei10 individual may have had lighter hair (Pal7 hair shade was inconclusive) (see Table S32). We caution against over-interpreting the Hirisplex results, which are based on modern patterns of variation, including the fact that rs1426654 is nearly fixed in modern Europe and therefore uninformative, which is not true of our ancient Aegean samples. The power of the reconstruction is further compromised by missing genotypes at several sites (see Table S31).

## Lactase Persistence

Klei10, Bar8, and Bar31 have the genotypes associated lactase non-persistence at the two lactase-persistence-associated loci under selection in modern Europeans (rs182549 and rs4988235) [165]; Pal7 could not be genotyped at rs182549 but has the ancestral genotype at rs4988235; Rev5 could not be genotyped at rs4988235 but has the ancestral genotype at rs182549. These Neolithic Aegeans predate the period during which the lactase persistence mutation is thought to have reached appreciable frequency in Europe, around 4kya [81, 130, 166, 167, 168, 169, 170]. The *LCT* T-13910 and A-22018 mutations remain at relatively low frequencies in modern Greek, Turkish, and Sardinian populations relative to modern central and northern Europeans (see Tables S33 & S35).

## Markers selected in Ancient Eurasians

We additionally typed the Neolithic Aegeans at a panel of twelve SNPs identified by Mathieson *et al.* [130] as being under selection in ancient Eurasians (see Table S33). In addition to the *HERC2* and *SLC45A2* pigmentation loci, derived alleles can be observed in multiple ancient Aegeans at loci associated with skin pigmentation (rs7119749 in *GRM5*), vitamin D status (rs7944926 near *DHCR7/NADSYN1*) and susceptibility to celiac disease (rs272872 near *SLC22A4*).

Following the observation that loci associated with inflammatory and metabolic disease show evidence of selection in both modern and ancient genomic datasets [130, 171, 172], we genotyped our samples at a panel of these markers for which selective sweep ages have been estimated [171, 172], as well as two loci in the *HBB* gene associated with malaria resistance/susceptibility to anemia (see Table S34).

Examination of the *TCF7L2* SNPs indicates that the two Neolithic Anatolian samples, Bar8 and Bar31, are likely to have carried at least one copy of haplotype (defined by rs7903146 C and rs10885406 A and tagged by rs7924080 T) conferring some protection against type-2 diabetes. The

presence of the derived state of the tag SNP in the Klei10 and Rev5 individuals suggests they may have also carried a copy of the haplotype. This haplotype is believed to have been a target of natural selection in Europeans (as well as in East Asians and West Africans), with evidence for a selective sweep occurring around 11,900 years ago in Europe [172].

The Neolithic Aegeans display derived selected alleles at a number of positions associated with susceptibility to inflammatory disease (e.g. Crohn’s disease, rheumatoid arthritis, celiac disease, multiple sclerosis, ulcerative colitis), including at a network of markers estimated to have undergone selection between 2.6-1.2 kya in Europeans [171], suggesting that selection acted on standing variation present at appreciable frequency. We additionally typed our samples at three loci in the *SLC22A4/SLC22A5* IBD5 haplotype associated with Crohn’s disease, which is estimated to have swept to higher frequency in Eurasians approximately 12,500 years ago, possibly in relation to the transition to an agriculturalist diet [173, 174]. The Aegeans were not observed to carry any of the hitchhiking deleterious alleles.

## Additional Functional Markers

No derived alleles were observed for any of the *ALDH2* or *ADH1B* alcohol-metabolism-associated mutations under selection in modern Asians [175, 176], and all 5 ancient Aegeans carried the common European ancestral genotypes at the *EDAR* locus under selection in Asians [157] (see Table S32).

We also genotyped a panel of SNPs in the *NAT2* in order to infer the acetylation status (slow > intermediate > rapid) of the ancient Aegean individuals using the NAT2pred online tool (NAT2pred.rit.allembany.edu) [177]. *NAT2* is involved in xenobiotic metabolism, and may have experienced selective pressure related to the adoption of different dietary lifeways [178, 179, 180]. The Anatolian Bar31 and Bar8 individuals had sufficient coverage for the inference of acetylation status using this method; both individuals are reconstructed as intermediate acetylators. Using the tag SNP rs1495741, a less sensitive method, predicts that Pal7 is a slow acetylator and that Rev5 is a rapid acetylator [181]. A second two-SNP genotyping method supports the prediction of rapid acetylation status for Rev5 and intermediate acetylation status for the two Anatolian samples [182]. The acetylation status of Klei10 could not be reconstructed.

We additionally typed the Aegeans for several mutations predisposing carriers to  $\beta$ -thalassemia, an autosomal recessive disorder characterized by anemia (S33). The mediterranean zone from which our samples derive is part of the modern and historical range of  $\beta$ -thalassemia phenotypes. The relatively high modern incidence of  $\beta$ -thalassemia carriers in the mediterranean (between 2 and 17% in different micro-regions) is thought to be related to heterozygote advantage versus *Plasmodium falciparum* malaria, as mosquitoes carrying the disease are historically endemic to the  $\beta$ -thalassemia zone [183]. We selected the five most prevalent disease-associated SNPs in modern Turkey, which together account for roughly 65% of observed cases (there is a long tail of minor mutations accounting



for the remainder of the case spectrum) [184]. The Neolithic Aegeans were not observed to carry the  $\beta$ -thalassemia mutations at any of these sites.

**Table S31:** Genotypes at SNPs associated with pigmentation. A: Ancestral, D: Derived Alleles.

<i>Hirisplex SNPs</i>							
SNP	Gene	A/D	Klei10	Pal7	Rev5	Bar31	Bar8
n29insa	<i>MC1R</i>	C/insA	C/C (3X)	C/C (2X)	C/C (2X)	-	C/C (2X)
rs11547464	<i>MC1R</i>	G/A	G/G (5X)	-	-	G/G (2X)	G/G (10X)
rs885479	<i>MC1R</i>	G/A	-	-	-	G/G (3X)	G/G (6X)
rs1805008	<i>MC1R</i>	C/T	-	-	-	-	-
rs1805005	<i>MC1R</i>	G/T	G/G (5X)	-	-	G/G (3X)	G/G (5X)
rs1805006	<i>MC1R</i>	C/A	C/C (4X)	-	-	C/C (3X)	C/C (8X)
rs1805007	<i>MC1R</i>	C/T	C/C (2X)	-	-	-	-
rs1805009	<i>MC1R</i>	G/C	G/G (2X)	-	-	G/G (3X)	G/G (8X)
y152och	<i>MC1R</i>	C/A	C/C (2X)	-	-	-	C/C (7X)
rs2228479	<i>MC1R</i>	G/A	G/G (3X)	-	-	G/G (6X)	G/G (5X)
rs1110400	<i>MC1R</i>	T/C	-	-	-	-	T/T (9X)
rs28777	<i>SLC45A2</i>	C/A	A/A(4X)	-	-	A/A (4X)	C/A (5X)
rs16891982	<i>SLC45A2</i>	C/G	G/G (2X)	G/G (3X)	C/C (2X)	C/G (5X)	C/G (9X)
rs12821256	<i>KITLG</i>	T/C	T/T (3X)	T/T (5X)	-	-	T/T (9X)
rs4959270	<i>EXOC2</i>	C/A	-	-	-	-	C/C (4X)
rs12203592	<i>IRF4</i>	C/T	C/C (3X)	-	-	C/C (5X)	C/C (8X)
rs1042602	<i>TYR</i>	C/A	C/A (2X)	C/C (3X)	C/C (3X)	C/C (7X)	C/C (10X)
rs1800407	<i>OCA2</i>	C/T	-	-	-	-	C/C (6X)
rs2402130	<i>SLC24A4</i>	G/A	-	-	-	A/A (10X)	G/A (6X)
rs12913832	<i>HERC2</i>	A/G	A/G (4X)	A/G (3X)	A/A (2X)	A/A (4X)	A/A (12X)
rs2378249	<i>PIGU/ASIP</i>	A/G	-	A/A (3X)	-	-	G/A (7X)
rs12896399	<i>SLC24A4</i>	G/T	G/T (3X)	-	-	G/T (5X)	G/T (7X)
rs1393350	<i>TYR</i>	G/A	G/G (3X)	G/G (3X)	G/G (3X)	G/G (3X)	G/G (15X)
rs683	<i>TYRP1</i>	A/C	C/C (3X)	A/A (2X)	-	A/A (2X)	C/A (10X)
<i>SLC24A5 16-marker haplotype</i>							
rs1834640	<i>SLC24A5</i>	G/A	-	A/A (4X)	A/A (2X)	A/A (8X)	A/A (6X)
rs2675345	<i>SLC24A5</i>	G/A	-	-	-	A/A (7X)	A/A (10X)
rs2469592	<i>SLC24A5</i>	G/A	A/A (4X)	A/A (3X)	-	-	A/A (13X)
rs2470101	<i>SLC24A5</i>	C/T	-	-	-	T/T (2X)	T/T (9X)
rs938505	<i>SLC24A5</i>	C <sup>1</sup> /T	C/C (3X)	C/C (2X)	-	C/C (3X)	C/C (3X)
rs2433354	<i>SLC24A5</i>	T/C	C/C (2X)	C/T (2X)	-	-	-
rs2459391	<i>SLC24A5</i>	G/A	A/A (4X)	A/A (3X)	A/A (2X)	A/A (6X)	A/A (10X)
rs2433356	<i>SLC24A5</i>	A/G	-	-	G/G (5X)	G/G (7X)	G/G (23X)
rs2675347	<i>SLC24A5</i>	G/A	A/A (2X)	-	-	A/A (3X)	A/A (5X)
rs2675348	<i>SLC24A5</i>	G/A	A/A (4X)	-	A/A (3X)	A/A (14X)	A/A (17X)
rs1426654	<i>SLC24A5</i>	G/A	A/A (4X)	-	A/A (5X)	A/A (11X)	A/A (7X)
rs2470102	<i>SLC24A5</i>	G/A	A/A (2X)	A/A (3X)	-	A/A (5X)	A/A (8X)
rs16960631	<i>SLC24A5</i>	A <sup>1</sup> /G	A/A (3X)	A/A (4X)	A/A (3X)	A/A (3X)	A/A (8X)
rs2675349	<i>SLC24A5</i>	G/A	A/A (2X)	-	A/A (3X)	A/A (9X)	A/A (7X)
rs3817315	<i>SLC24A5</i>	T/C	C/C (3X)	-	-	C/C (5X)	C/C (2X)
rs7163587	<i>SLC24A5</i>	T/C	C/C (2X)	-	C/C (2X)	C/C (4X)	C/C (13X)
<i>SLC45A2 12-marker haplotype</i>							
rs732740	<i>SLC45A2</i>	A/G	-	A/A (2X)	A/A (3X)	-	A/A (4x)
rs250413	<i>SLC45A2</i>	G/A	G/G (2X)	-	-	G/A (4X)	-
rs181832	<i>SLC45A2</i>	G/A	A/A (2X)	-	-	A/A (4X)	A/G (6X)
rs3776549	<i>SLC45A2</i>	C/T	-	-	C/C (2X)	C/C (8X)	C/C (12X)
rs3756462	<i>SLC45A2</i>	A/G	-	-	G/G (3X)	A/A (9X)	-
rs26722	<i>SLC45A2</i>	C/T	C/C (4X)	-	-	C/C (5X)	C/C (15X)
rs2287949	<i>SLC45A2</i>	C/T	C/C (2X)	-	C/T (4X)	C/C(2X)	C/C (3X)
rs250417	<i>SLC45A2</i>	G/C	-	-	G/G (2X)	C/C (2X)	C/C (3X)
rs16891982	<i>SLC45A2</i>	C/G	G/G (2X)	G/G (3X)	C/C (2X)	C/G (5X)	C/G (9X)
rs40132	<i>SLC45A2</i>	A/G	A/A (4X)	-	-	A/A (2X)	A/A (9X)
rs35394	<i>SLC45A2</i>	T/C	-	-	-	-	T/T (2X)

<sup>1</sup>ancestral allele part of the C11 haplotype

SNP	Gene	A/D	Klei10	Pal7	Rev5	Bar31	Bar8
rs3733808	<i>SLC45A2</i>	C/G	C/C(2X)	-	-	C/C (2X)	-

**HERC2 13-marker haplotype**

rs4778241	<i>HERC2</i>	A/C	A/A (2X)	-	-	A/A (7X)	A/C (16X)
rs1129038	<i>HERC2</i>	C/T	C/T (4X)	-	-	C/C (4X)	C/C (3X)
rs12593929	<i>HERC2</i>	G/A	-	-	-	-	G/G (5X)
rs12913832	<i>HERC2</i>	A/G	A/G(4X)	A/A(3X)	A/A (2X)	A/A (4X)	A/A (12X)
rs7183877	<i>HERC2</i>	C <sup>2</sup> /A	C/C (2X)	-	C/C (4X)	C/C (2X)	C/C (10X)
rs3935591	<i>HERC2</i>	T/C	C/C (3X)	-	-	T/T (10X)	C/C (6X)
rs7170852	<i>HERC2</i>	T/A	A/A (3X)	T/T (2X)	-	T/T (4X)	A/A (6X)
rs2238289	<i>HERC2</i>	G/A	-	-	-	G/G (2X)	A/A (7X)
rs3940272	<i>HERC2</i>	T/G	-	-	-	-	-
rs8028689	<i>HERC2</i>	C/T	T/T (3X)	C/T (3X)	-	-	T/T (9X)
rs2240203	<i>HERC2</i>	C/T	-	-	-	T/T (3X)	T/T (13X)
rs11631797	<i>HERC2</i>	A/G	-	A/G (2X)	-	A/A (2X)	-
rs916977	<i>HERC2</i>	T/C	C/C (2X)	C/T (2X)	C/C (2X)	T/T (8X)	C/C (7X)

<sup>2</sup>ancestral allele part of the blue-eye-associated haplotype

*Table S32: Hirisplex model results.*

	<b>Klei 10</b>		<b>Pal7</b>		<b>Bar31</b>		<b>Bar8</b>	
	<i>p-value</i>	<i>AUC Loss</i>	<i>p-value</i>	<i>AUC Loss</i>	<i>p-value</i>	<i>AUC Loss</i>	<i>p-value</i>	<i>AUC Loss</i>
blue eye	0.13	0.002	0.001	0.014	0	0.002	0	0
intermediate eye	0.154	0.009	0.026	0.038	0.011	0.009	0.007	0
brown eye	0.716	0.005	0.973	0.011	0.989	0.005	0.993	0
blond hair	0.409	0.021	0.000	0.04	0	0.037	0	0.028
brown hair	0.45	0.02	0.000	0.042	0	0.031	0	0.024
red hair	0.025	0.077	0.000	0.288	0	0.213	0	0.178
black hair	0.116	0.008	0.000	0.033	0	0.027	0	0.018
light hair	0.756	0.007	0.471	0.016	0.217	0.015	0.035	0.008
dark hair	0.244	0.007	0.529	0.016	0.783	0.015	0.965	0.008

**Table S33:** Genotypes at other functional markers under selection in modern and ancient Eurasian populations. A: Ancestral, D: Derived Alleles.

SNP	Gene	A/D	Klei10	Pal7	Rev5	Bar31	Bar8
rs3827760	<i>EDAR</i>	A/G	A/A (2X)	-	A/A (2X)	A/A (4X)	A/A (7X)
rs4988235	<i>LCTa</i>	G/A	G/G (4X)	G/G (4X)	-	G/G (4X)	G/G (9X)
rs182549	<i>LCTb</i>	C/T	C/C (4X)	-	C/C (3X)	C/C (3X)	C/C (7X)
rs3811801	<i>ADH1Ba</i>	G/A	G/G (3X)	G/G (4X)	G/G (2X)	-	G/G (6X)
rs1229984	<i>ADH1Bb</i>	C/T	C/C (3X)	-	-	C/C (3X)	-
rs671	<i>ALDH2</i>	G/A	G/G (2X)	-	-	G/G (2X)	G/G (2X)
rs1801279	<i>NAT2</i>	G/A	-	G/G (2X)	-	G/G (5X)	G/G (3X)
rs1041983	<i>NAT2</i>	C/T	-	-	C/T (2X)	C/C (3X)	C/T (5X)
rs1801280	<i>NAT2</i>	T/C	-	-	T/T (2X)	C/T (5X)	T/T (4X)
rs1799929	<i>NAT2</i>	C/T	C/T (3X)	-	-	C/C (5X)	C/C (10X)
rs1799930	<i>NAT2</i>	G/A	-	-	A/G (3X)	G/G (5X)	G/A (9X)
rs1208	<i>NAT2</i>	A/G	-	A/A (2X)	A/A (3X)	G/A (7X)	A/A (6x)
rs1799931	<i>NAT2</i>	G/A	-	-	-	G/G (4X)	G/G (5X)
rs1495741	<i>NAT2</i>	G/A	-	A/A (3X)	G/G (2X)	G/A (8X)	G/A (10X)
rs2269424 <sup>3</sup>	<i>MHC region</i>	G/A	-	-	-	A/A (2X)	G/G (3X)
rs174546 <sup>3</sup>	<i>FAD1/</i> <i>FADS2</i>	T/C	T/T (2X)	T/T (2X)	-	T/T (4X)	T/T (5X)
rs4833103 <sup>3</sup>	<i>TLR1/</i> <i>TLR6/</i> <i>TLR10</i>	C/A	C/C (2X)	-	C/C (2X)	C/C (3X)	C/C (7X)
rs653178 <sup>3</sup>	<i>ATXN2/</i> <i>SH2B3</i>	T/C	T/T (3X)	-	T/T (2X)	T/T (4X)	C/T (9X)
rs7944926 <sup>3</sup>	<i>DHCR7/</i> <i>NADSYN1</i>	G/A	A/G (3X)	-	A/A (1X)	A/A (4X)	A/G (7X)
rs7119749 <sup>3</sup>	<i>GRM5</i>	A/G	-	A/A (2X)	A/G (3X)	A/A (3X)	A/G (10X)
rs272872 <sup>3</sup>	<i>SLC22A4</i>	G/A	A/A (3X)	-	G/G (5X)	A/G (6X)	A/G (8X)
rs6903823 <sup>3</sup>	<i>ZKSCAN3/ ZSCAN31</i>	A/G	-	-	-	A/A (3X)	A/A (6X)
rs1979866 <sup>3</sup>	-	C/A	-	-	-	-	A/A (4X)
rs35004220 (IVS-I-110)	<i>HBB</i>	C/T	C/C (2X)	C/C (2X)	C/C (2X)	-	C/C (9X)
rs35724775 (IVS-I-6)	<i>HBB</i>	A/G	A/A (2X)	-	-	-	A/A (6X)
rs35497102 (FSC-8)	<i>HBB</i>	TT/-	TT (4X)	-	-	-	TT (4X)
rs33971440 (IVS-I-1)	<i>HBB</i>	C/T	C/C (4X)	-	-	C/C (2X)	C/C (5X)
rs34690599 (IVS-II-745)	<i>HBB</i>	G/C	G/G (2X)	G/G (3X)	G/G (2X)	G/G (5X)	G/G (10X)

<sup>3</sup>Additional markers identified by Mathieson *et al.* [130] as under selection in prehistoric Eurasians

**Table S34:** Markers related to pathogen resistance, metabolic and inflammatory disease. A: Ancestral, D: Derived Alleles.

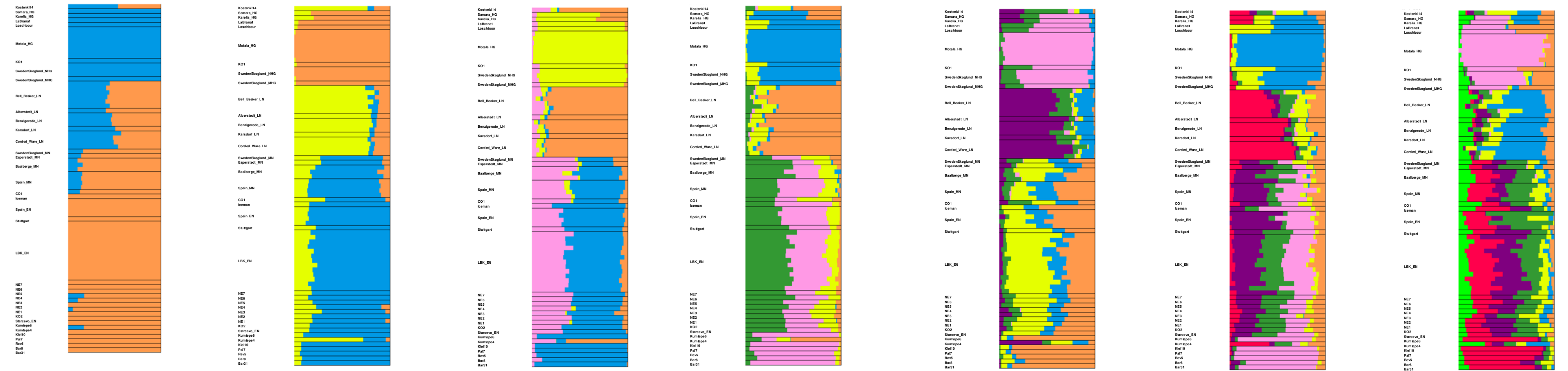
SNP	Gene	A/D	Klei10	Pal7	Rev5	Bar31	Bar8	Estimated sweep age (ya) [171, 172, 173]
rs10786436	<i>HPSE2</i>	C/T	-	-	-	-	C/C (4X)	~ 6,560
rs1132200 <sup>4</sup>	<i>ARHGAP31/STAT1</i>	C/T	C/T (5X)	-	-	C/T (5X)	C/C (7X)	~ 2,200
rs12638253	<i>LEKR1</i>	C/T	T/T (4X)	-	C/T (2X)	C/T (4X)	C/T (6X)	~ 5,100
rs12722489 <sup>4</sup>	<i>IL2RA</i>	C/T	C/C (6X)	C/T (3X)	-	-	C/C (4X)	~ 1,200
rs17696736 <sup>4</sup>	<i>TMEM116</i>	A/G	A/A (4X)	-	A/A (2X)	A/A (5X)	A/A (10X)	~ 1,400
rs17810546 <sup>4</sup>	<i>SCHIP1/IL12A</i>	A/G	-	-	A/A (2X)	A/A (6X)	A/A (9X)	~ 2,310
rs2058660	<i>IL18RAP</i>	A/G	A/A (3X)	-	-	A/A (6X)	A/A (6X)	~ 7,500
rs2188962 <sup>4</sup>	<i>SLC22A5/IRF1</i>	C/T	C/C (6X)	-	-	-	C/T (10X)	~ 1,380
rs2248359	<i>CYP24A1</i>	T/C	C/C (3X)	C/C (2X)	C/C (3X)	C/C (2X)	T/T (9X)	~ 8,500
rs2285795	<i>TRIM10</i>	T/C	C/C (3X)	T/T (2X)	-	C/C (2X)	C/C (7X)	~ 2,280
rs307896	<i>SAE1</i>	G/A	A/A (2X)	-	-	G/A (7X)	G/A (10X)	~ 2,700
rs3129934	<i>BTNL2</i>	C/T	C/C (2X)	-	-	C/C (3X)	C/C (7X)	~ 3,300
rs3131379	<i>VARS/LSM2</i>	G/A	G/G (3X)	-	-	G/G (3X)	G/G (3X)	~ 1,980
rs3184504	<i>SH2B3/ATXN2</i>	C/T	C/C (3X)	T/T (2X)	C/C (2X)	-	C/C (3X)	~ 1,500
rs6822844 <sup>4</sup>	<i>IL2/IL21</i>	G/T	G/G (2X)	G/G (2X)	-	G/G (4X)	G/T (8X)	~ 2,150
rs6897932 <sup>4</sup>	<i>IL7R</i>	C/T	T/T (2X)	C/C (2X)	C/C (2X)	C/T (6X)	C/T (11X)	~ 1,800
rs744166 <sup>4</sup>	<i>STAT3</i>	G/A	-	-	A/G (2X)	G/G (6X)	G/G (9X)	~ 2,600
rs1050152	<i>SLC22A4</i>	C/T	C/C (5X)	-	-	C/C (6X)	C/C (12X)	~ 12,500
rs2631367	<i>SLC22A5</i>	G/C	G/G (2X)	-	-	-	-	~ 12,500
rs11739623	<i>IRF1/IL5</i>	C/T	C/C (4X)	C/C (5X)	C/C (3X)	C/C (4X)	C/C (7X)	~ 12,500
rs7903146	<i>TCF7L2</i>	T/C	-	-	-	C/T (3X)	C/C (3X)	~ 11,900
rs10885406	<i>TCF7L2</i>	G/A	A/G (2X)	A/A (2X)	-	A/G (5X)	A/A (10X)	~ 11,900
rs12255372	<i>TCF7L2</i>	G/T	G/G (3X)	T/T (3X)	G/T (3X)	G/T (12X)	G/G (13X)	~ 11,900
rs7924080	<i>TCF7L2</i>	T/C	T/C (9X)	-	T/C (3X)	T/C (7X)	-	~ 11,900

**Table S35:** Modern distributions of lactase and pigmentation derived allele frequencies. Number of chromosomes in parentheses.

SNP/Gene	Sardinian [156]	Turk [156]	Greek [156]	CEU [155]	FIN [155]	GBR [155]	IBS [155]	TSI [155]
rs4988235 <i>LCTa</i> C > T	0.03 (70)	.03 [185] (98)	.08 (40)	0.74 (198)	0.59 (198)	0.72 (182)	0.46 (214)	0.09 (214)
rs182549 <i>LCTb</i> G > A	0.07 (56)	0.08 (112)	0.146 (84)	0.74 (198)	0.59 (198)	0.72 (182)	0.46 (214)	0.09 (214)
rs1426654 <i>SLC24A5</i> G > A	0.98 (118)	0.98 (260)	1 (184)	1 (198)	0.99 (198)	1 (182)	1 (214)	>0.99 (214)
rs16891982 <i>SLC45A2</i> C > G	0.71 (224)	0.68 (272)	0.86 (1198)	0.98 (198)	0.96 (198)	0.98 (182)	0.82 (214)	0.97 (214)
rs12913832 <i>HERC2</i> A > G	0.22 (300)	0.33 (390)	0.36 (1382)	0.77 (198)	0.91 (198)	0.82 (182)	0.32 (214)	0.42 (214)

<sup>4</sup>SNPs in core selected network identified by Raj *et al.* 2013 [171]

a



b



c



Figure S22: Graphs showing ancestry estimated among various hunter-gatherer and farmer populations using ADMIXTURE for K=12 to K=18. Analyses were performed without CHG and Yamnaya (a), with CHG (b) and with both Yamnaya and CHG (c) genomes from the Yamnaya culture. Early Neolithic genomes cluster with the Argos Neolithic genomes, while genomes from Middle Neolithic show increased amounts of hunter-gatherer ancestry. Genomes from the Late Neolithic additionally demonstrate a substantial amount of ancestry from a group related to the people of the Yamnaya culture.

## Data available online

**Dataset S1** List with all genomes discussed in the text giving relevant information in chronological order.

**Dataset S2** This file contains detailed results of the f3- and f4- statistics (see SI7).

**Dataset S3** This file contains a detailed table of the inferred mixture coefficients when forming the DNA of each target group (rows) as mixtures of that from other surrogate groups (columns) (see SI10).

**3D-figure S4** Interactive version of Figure 2 including the third principal component. See legend of Figure 2 for colour code. The interactive file can be accessed at [https://figshare.com/articles/Hofmanova\\_et\\_al\\_3D\\_figure\\_S4/3188767](https://figshare.com/articles/Hofmanova_et_al_3D_figure_S4/3188767).



## References

1. Reingruber A (2008) *Die Argissa-Magula. Das frühe und das beginnende mittlere Neolithikum im Lichte transägäischer Beziehungen. Die deutschen Ausgrabungen auf der Argissa-Magula in Thessalien II. Beiträge zur ur- und frühgeschichtlichen Archäologie des Mittelmeer-Kulturraumes.* (Dr. Rudolf Habelt GmbH, Bonn) Vol. 35.
2. Perlès C, Takaoğlu T, Gratuze B (2011) Melian obsidian in NW Turkey: Evidence for early Neolithic trade. *Journal of Field Archaeology* 36(1):42–49.
3. Özdoğan M, Başgelen N, Kuniholm P (2012) *The Neolithic in Turkey: New Excavations and New Research. Western Turkey.* (Archaeology and Art Publications, Istanbul) Vol. 4.
4. Çilingiroğlu Ç, Çakırlar C (2013) Towards configuring the Neolithisation of Aegean Turkey. *Documenta Praehistorica* 40:21–29.
5. Weninger B et al. (2014) Neolithisation of the Aegean and Southeast Europe during the 6600–6000 calBC period of Rapid Climate Change. *Documenta Praehistorica* 41:1–31.
6. Özdoğan E (2015) Current research and new evidence for the Neolithization process in western Turkey. *European Journal of Archaeology* 18(1):33–59.
7. Zilhão J (2001) Radiocarbon evidence for maritime pioneer colonization at the origins of farming in west Mediterranean Europe. *Proceedings of the National Academy of Sciences* 98(24):14180–14185.
8. Kotsakis K (2001) Mesolithic to Neolithic in Greece. Continuity, discontinuity or change of course. *Documenta Praehistorica* 28:63–73.
9. Kotsakis K (2008) A Sea of Agency: Crete in the Context of the Earliest Neolithic in Greece in *Escaping the Labyrinth : The Cretan Neolithic in context*, eds. Isaakidou V, Tomkins PD. (Oxbow Books), pp. 49–72.
10. Broodbank C (2013) *The Making of the Middle Sea.* (Oxford University Press, Oxford).
11. Horejs B et al. (2015) The Aegean in the early 7th millennium BC: Maritime networks and colonization. *Journal of World Prehistory* 28(4):289–330.
12. Özdoğan M (2013) Neolithic Sites in the Marmara Region, Fikirtepe, Pendik, Yarımburgaz, Toptepe, Hocaçeşme, and Aşağı Pınar. *The Neolithic in Turkey. New Excavations & New Research* 5:167–269.
13. Özbal R, Gerritsen F (in press) Barcın Höyük in interregional perspective in *Social and Economic Changes in the Second Half of the Seventh Millennium in the Near East*, ed. Marciniak A. (Lockwood Press, Atlanta GA).
14. Özdoğan M (2007) Marmara Bölgesi Neolitik çağ kültürleri in *Türkiye’de Neolitik Dönem*, eds. Özdoğan M, Başgelen N. (Arkeolojive sanat yayinlari, Istanbul), pp. 401–426.
15. Özdoğan M (2007) Amidst Mesopotamia-centric and Euro-centric approaches: the changing role of the Anatolian peninsula between the East and the West. *Anatolian Studies* 57:17–24.
16. Özdoğan M (2010) Westward expansion of the Neolithic way of life: Sorting the Neolithic package into distinct packages in *Proceedings of the 6th International Congress on the Archaeology of the Ancient Near East. May, 5th-10th 2008. “Sapienza” Università di Roma*, eds. Matthiae P, Pinnock F, Nigro L, Marchetti N. (Harrassowitz Verlag, Wiesbaden), pp. 883–97.
17. Gerritsen F, Özbal R, Thissen L (2013) Barcın Höyük: the beginnings of farming in the Marmara Region in *The Neolithic in Turkey: New Excavations and New Research. Northwestern Turkey and Istanbul*, eds. Özdoğan M, Başgelen N, Kuniholm P. (Archaeology and Art Publications, Istanbul) Vol. 5, pp. 93–112.
18. Roodenberg J, Van As A, Jacobs L, Wijnen M (2003) Early settlement in the plain of Yenişehir (NW Anatolia). The basal occupation layers at Menteşe. *Anatolica* 29:17–60.

19. Karul N, Avcı MB (2013) Aktopraklık in *The Neolithic in Turkey: New Excavations and New Research. North-western Turkey and Istanbul*, eds. Özdoğan M, Başgelen N, Kuniholm P. (Archaeology and Art Publications, Istanbul) Vol. 5, pp. 45–68.
20. Çakırlar C (2013) Rethinking Neolithic subsistence at the gateway to Europe with new archaeozoological evidence from Istanbul in *Barely Surviving or More than Enough? - The Environmental Archaeology of Subsistence, Specialisation and Surplus Food Production*, eds. Groot M, Lentjes D, Zeiler J. (Sidestone Press), pp. 59–79.
21. Evershed RP et al. (2008) Earliest date for milk use in the Near East and southeastern Europe linked to cattle herding. *Nature* 455(7212):528–531.
22. Thissen L, Özbal H, Türkekul-Bıyık A, Gerritsen F, Özbal R (2010) The land of milk? Approaching dietary preferences of Late Neolithic communities in NW Anatolia. *Leiden Journal of Pottery Studies* 26:157–172.
23. Milić M (2014) PXRf characterisation of obsidian from central Anatolia, the Aegean and central Europe. *Journal of Archaeological Science* 41:285–296.
24. Krauß R (2009) Karanovo und das südosteuropäische Chronologiesystem aus heutiger Sicht. *Eurasia Antiqua* 14:115–147.
25. Lichardus-Itten M (1993) Zum Beginn des Neolithikums im Tal der Struma (Südwest-Bulgarien). *Anatolica* 19:99–116.
26. Karamitrou-Mentessidi G et al. (2013) New evidence on the beginning of farming in Greece: The Early Neolithic settlement of Mavropigi in Western Macedonia (Greece). *Antiquity* 87(336).
27. Karamitrou-Mentessidi G (2014) About prehistoric sites in western Macedonia: prefectures of Kozani and Grevena in *A century of Research in Prehistoric Macedonia. International Conference Proceedings 22-24 November 2012*, eds. Stefani, Merousis, Dimoula. (Archaeological Museum of Thessaloniki, Thessaloniki), pp. 233–249.
28. Maniatis Y (2014) Radiocarbon dating of the major cultural changes in Prehistoric Macedonia: recent developments in *A century of Research in Prehistoric Macedonia. International Conference Proceedings 22-24 November 2012*, eds. Stefani, Merousis, Dimoula. (Archaeological Museum of Thessaloniki, Thessaloniki) Vol. 10, pp. 205–222.
29. Maniatis Y, Kotsakis K, Halstead P (2012) New radiocarbon chronology of early Neolithic in Macedonia. *Archaeological Work in Macedonia and Thrace*. In press (in Greek).
30. Kotsakis K (2014) Domesticating the periphery. New research into the Neolithic of Greece. *Pharos* 20(1):41–73.
31. Marinova E (2006) Vergleichende paläoethnobotanische Untersuchung zur Vegetationsgeschichte und zur Entwicklung der prähistorischen Landnutzung in Bulgarien. *Dissertationes Botanicae* 401.
32. Colledge S, Conolly J, Shennan S (2004) Archaeobotanical evidence for the spread of farming in the eastern Mediterranean. *Current Anthropology* 45:35–58.
33. Krauß R et al. (2014) Beginnings of the Neolithic in Southeast Europe: the Early Neolithic sequence and absolute dates from Džuljunica-Smărdeš (Bulgaria). *Documenta Praehistorica* 41:51–77.
34. Craig OE et al. (2005) Did the first farmers of central and eastern Europe produce dairy foods? *Antiquity* 79(306):882–894.
35. Scheu A et al. (2015) The genetic prehistory of domesticated cattle from their origin to the spread across Europe. *BMC Genetics* 16(1):54.
36. Valamoti S, Kotsakis K (2007) Transitions to agriculture in the Aegean: the archaeobotanical evidence in *The origin and spread of domestic plants in southwest Asia and Europe*, eds. Colledge S, Conolly J. (Left Coast Press, London), pp. 76–92.

37. Trantalidou K (2011) From Mesolithic fishermen and bird hunters to Neolithic goat herders: The mammal and bird bone assemblages in *The Cave of the Cyclops: Mesolithic and Neolithic Networks in the Northern Aegean, Greece II: Bone Tool Industries, Dietary Resources, the Paleoenvironment, and Archaeometrical Studies*, ed. Sampson A. (INSTAP Academic Press) Vol. 2.
38. Oross K, Bánffy E (2009) Three successive waves of Neolithisation: LBK development in Transdanubia. *Documenta Praehistorica* 36:175–189.
39. Bramanti B et al. (2009) Genetic discontinuity between local hunter-gatherers and central Europe's first farmers. *Science (New York, N. Y.)* 326(5949):137–40.
40. Skoglund P et al. (2012) Origins and genetic legacy of Neolithic farmers and hunter-gatherers in Europe. *Science (New York, N. Y.)* 336(6080):466–9.
41. Lazaridis I et al. (2014) Ancient human genomes suggest three ancestral populations for present-day Europeans. *Nature* 513(7518):409–413.
42. Kaczanowska M, Kozłowski JK (2014) The origin and spread of the western Linear Pottery culture: Between forager and food producing lifeways in Central Europe. *Archaeologiai Értésítő* 139(1):293–318.
43. Gerritsen FA, Özbal R, Thissen L (2013) The earliest Neolithic levels at Barcın Höyük, Northwestern Turkey. *Anatolica* 39:53–92.
44. Gerritsen F, Özbal R (in press) Barcın Höyük and the pre-Fikirtepe Neolithization of the Eastern Marmara Region in *Anatolian Metal VII, Der Anschnitt, Beiheft, Bochum, 2016.*, ed. Yalcin Ü.
45. Groenhuizen MR, Kluiving SJ, Gerritsen FA, Künnel M (2015) Geoarchaeological research at Barcın Höyük: Implications for the initial Neolithic occupation of northwest Anatolia. *Quaternary International* 359–360:452–461.
46. Arbuckle BS et al. (2014) Data sharing reveals complexity in the westward spread of domestic animals across Neolithic Turkey. *PLoS ONE* 9(6).
47. Alpaslan Roodenberg S, Gerritsen F, Özbal R (2013) Neolithic burials from Barcın Höyük. *Anatolica* 39:93–111.
48. Facorellis Y (2003) Radiocarbon dating the Greek Mesolithic in *The Greek Mesolithic - Problems and Perspectives*, eds. Galanidou N, Perlès C. (British School at Athens Studies), pp. 51–67.
49. Galanidou N (2014) Advances in the Palaeolithic and Mesolithic archaeology of Greece for the new millennium. *Pharos* 20:1–40.
50. Reingruber A, Thissen L (2005) <sup>14</sup>C database for the Aegean catchment (eastern Greece, southern Balkans and western Turkey) 10,000–5500 cal BC in *How did farming reach Europe?*, ed. Lichter C. (BYZAS 2), pp. 295–327.
51. Andreou S, Fotiadis M, Kotsakis K (1996) Review of Aegean prehistory V: The Neolithic and Bronze Age of northern Greece. *American Journal of Archaeology* 100:537–597.
52. Galanidou N, Perlès C (2003) An introduction to the Greek Mesolithic in *The Greek Mesolithic - Problems and Perspectives*, eds. Galanidou N, Perlès C. (British School at Athens Studies), pp. 27–32.
53. Krahtopoulou A (2000) Holocene alluvial history of northern Pieria, Macedonia, Greece in *Landscape and land use in postglacial Greece*, eds. Halstead P, Frederick C. (Sheffield Academic Press, Sheffield), pp. 15–27.
54. Kotsos S, Urem-Kotsou D (2006) Filling in the Neolithic landscape of central Macedonia, Greece in *Homage to Milutin Garasanin*, eds. Tasic N, Grozdanov C. (Serbian Academy of Science and Arts, Macedonian Academy of Sciences and Arts, Belgrade), pp. 193–207.

55. Perlès C, Quiles A, Valladas H (2013) Early seventh-millennium AMS dates from domestic seeds in the Initial Neolithic at Franchthi Cave (Argolid, Greece). *Antiquity* 87(338):1001–1015.
56. Kyparissi-Apostolika N (2003) The Mesolithic in Theopetra Cave: new data on a debated period of Greek prehistory in *The Greek Mesolithic - Problems and Perspectives*, eds. Galanidou N, Perlès C. (British School at Athens Studies) Vol. 10, pp. 189–198.
57. Jacobsen TW (1973) Excavations in the Franchthi Cave, 1969-1971: Part II. *Hesperia* 42:253–283.
58. Sampson A, Kaczanowska M, Kozłowski J, Alexandrowicz S (2010) *The prehistory of the island of Kythnos (Cyclades, Greece) and the Mesolithic settlement at Maroulas*. (The Polish Academy of Arts and Sciences, The University of the Aegean, Krakow).
59. Gallis K (1996) Burial customs in *Neolithic Civilisation in Greece*, ed. Papathanasopoulos G. (Goulandri Foundation-Museum of Cycladic Art, Athens), pp. 171–174.
60. Triantaphyllou S (2001) A bioarchaeological approach to prehistoric cemetery populations from central and western Greek Macedonia. *British Archaeological Reports International Series* 976.
61. Gallis K (1982) Human cremations from Neolithic Thessaly. *Editions of the Archaeologikon Deltion* 30.
62. Papathanasopoulos G (1996) Taphika ethima tou dirou in *Neolithikos politismos stin ellada*, ed. Papathanasopoulos G. (NP Goulandris Foundation-Museum of Cycladic Art, Athens), pp. 175–177.
63. Papathanasiou A (2001) A bioarchaeological analysis of Neolithic Alepotrypa Cave, Greece. *British Archaeological Reports International Series* 961.
64. Coleman JE (1977) *Keos. Results of excavations conducted by the University of Cincinnati under the auspices of the American School of Classical Studies at Athens*. (American School of Classical Studies, Princeton) Vol. 1, pp. 44–97.
65. Triantaphyllou S (2008) Living with the Dead: a Re-consideration of Mortuary Practices in the Greek Neolithic in *Escaping the Labyrinth: The Cretan Neolithic in Context*, eds. Isaakidou V, Tomkins P. (Oxbow Monographs, Oxford) Vol. 8, pp. 139–157.
66. Kyparissi-Apostolika N, Kotzamani G (2005) Worlds in Transition: Mesolithic/Neolithic Lifestyles at the Cave of Theopetra, Thessaly, Greece in *How did farming reach Europe?*, ed. Lichter C. (BYZAS 2), pp. 173–182.
67. Kyparissi-Apostolika N (1999) The Palaeolithic deposits of Theopetra Cave in Thessaly (Greece) in *The paleolithic Archaeology of Greece and Adjacent Areas, Proceedings of the ICOPAG Conference held at Ioannina 1994*, eds. Bailey G, Adam E, Panagopoulou E, Zachos K. (British School at Athens Studies) Vol. 3, pp. 232–239.
68. Kyparissi-Apostolika N (1999) The Neolithic use of Theopetra cave in Thessaly in *Neolithic Society in Greece*, ed. Halstead P. (Sheffield Academic Press, Sheffield).
69. Valladas H et al. (2007) TL age-estimates for the Middle Palaeolithic layers at Theopetra Cave (Greece). *Quaternary Geochronology* 2(1):303–308.
70. Karkanas P et al. (2014) Tephra correlations and climatic events between the MIS6/5 transition and the beginning of MIS3 in Theopetra Cave, central Greece. *Quaternary Science Reviews* 30:1–12.
71. Manolis SK, Stravopodi HJ (2003) An assessment of the human skeletal remains in the Mesolithic deposits of Theopetra Cave: a case study in *The Greek Mesolithic - Problems and Perspectives*, eds. Galanidou N, Perlès C. (British School at Athens Studies) Vol. 10, pp. 207–216.
72. Facorellis Y, Kyparissi-Apostolika N, Maniatis Y (2001) The cave of Theopetra, Kalambaka: radiocarbon evidence for 50,000 years of human presence. *Radiocarbon* 43(2; VOL B):1029–1048.

73. Bessios M, Adaktylou F (2006) A Neolithic settlement at Revenia near Korinos. *AEMTH* 18(2004):357–366 (in Greek).
74. Urem-Kotsou D, Papaioanou A, Silva-Gracia T, Adaktylou F, Besios M (2011) Settlement of Early and Middle Neolithic in Revenia Korinou. Preliminary results of the ceramic analysis. *AEMTH* 25. In press (in Greek).
75. Triantaphyllou S, Adaktylou F (2014) The treatment of the dead during the Early Neolithic period in Macedonia a first insight in Revenia, Korinou of Northern Pieria. *AEMTH* 28. In press (in Greek).
76. Kotsakis K, Halstead P (2004) Anaskafi sta neolithika paliambela kolindrou. *AEMTH* 16(2002):407–415 (in Greek).
77. Urem-Kotsou D, Papaioannou A, Papadakou T, Saridaki N, Intze Z (2014) Early and Middle Neolithic Pottery in Macedonia in *A Century of Research in Prehistoric Macedonia 1912-2012*, eds. Stefani E, Merousis N, Dimoula A. (Archaeological Museum of Thessaloniki, Zitis, Thessaloniki), pp. 505–517.
78. Halstead P, Kotsakis K (2002) Paliambela in *Archaeological Reports for 2001-2002*. (Society for the Promotion of Hellenic Studies and the British School at Athens), p. 80.
79. Halstead P, Kotsakis K (2003) Paliambela in *Archaeological Reports for 2002-2003*. (Society for the Promotion of Hellenic Studies and the British School at Athens), pp. 49–66.
80. Ziota C (2011) The neolithic settlement of “Kleitos I” and the new updates for prehistoric research in *The Archaeological Work in Upper Macedonia in 2009*, ed. Karamitrou-Mentessidi G. (Archaeological Museum of Aiani), pp. 211–230. (in Greek).
81. Gamba C et al. (2014) Genome flux and stasis in a five millennium transect of European prehistory. *Nature Communications* 5:5257.
82. Kircher M, Sawyer S, Meyer M (2012) Double indexing overcomes inaccuracies in multiplex sequencing on the Illumina platform. *Nucleic acids research* 40(1):e3.
83. Meyer M, Kircher M (2010) Illumina sequencing library preparation for highly multiplexed target capture and sequencing. *Cold Spring Harbor protocols* 2010(6):pdb.prot5448.
84. Bunce M, Oskam CL, Allentoft ME (2012) Quantitative real-time PCR in aDNA research in *Ancient DNA*, eds. Shapiro B, Hofreiter M. (Springer), pp. 121–132.
85. Meyer M et al. (2012) A high-coverage genome sequence from an archaic Denisovan individual. *Science (New York, N.Y.)* 338(6104):222–6.
86. Jonsson H, Ginolhac A, Schubert M, Johnson PLF, Orlando L (2013) mapdamage2.0: fast approximate Bayesian estimates of ancient DNA damage parameters. *Bioinformatics* 29(13):1682–1684.
87. Der Sarkissian C et al. (2015) Evolutionary genomics and conservation of the endangered Przewalski’s horse. *Current Biology* 25(19):2577–2583.
88. Gnirke A et al. (2009) Solution hybrid selection with ultra-long oligonucleotides for massively parallel targeted sequencing. *Nature biotechnology* 27(2):182–9.
89. van Oven M, Kayser M (2009) Updated comprehensive phylogenetic tree of global human mitochondrial DNA variation. *Hum. Mutat.* 30(2):E386–E394.
90. Blumenstiel B et al. (2001) Targeted exon sequencing by in-solution hybrid selection. *Current Protocols in Human Genetics*.
91. Dabney J et al. (2013) Complete mitochondrial genome sequence of a middle pleistocene cave bear reconstructed from ultrashort DNA fragments. *Proceedings of the National Academy of Sciences* 110(39):15758–15763.

92. Kircher M (2011) Analysis of high-throughput ancient DNA sequencing data. *Ancient DNA* pp. 197–228.
93. Aronesty E (2013) Comparison of sequencing utility programs. *TOBIOIJ* 7(1):1–8.
94. Li H, Durbin R (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 25(14):1754–1760.
95. Li H et al. (2009) The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25(16):2078–2079.
96. Breese MR, Liu Y (2013) NGSUtils: a software suite for analyzing and manipulating next-generation sequencing datasets. *Bioinformatics* 29(4):494–496.
97. DePristo MA et al. (2011) A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nature Genetics* 43(5):491–498.
98. Vianello D et al. (2013) HAPLOFIND: A new method for high-throughput mtDNA haplogroup assignment. *Human Mutation* 34(9):1189–1194.
99. Fu Q et al. (2013) A revised timescale for human evolution based on ancient mitochondrial genomes. *Current Biology* 23(7):553–559.
100. Szécsényi-Nagy A et al. (2015) Tracing the genetic origin of Europe’s first farmers reveals insights into their social organization. *Proceedings of the Royal Society of London B: Biological Sciences* 282(1805):20150339.
101. Brandt G et al. (2013) Ancient DNA reveals key stages in the formation of central European mitochondrial genetic diversity. *Science (New York, N. Y.)* 342(6155):257–61.
102. Haak W et al. (2015) Massive migration from the steppe was a source for Indo-European languages in Europe. *Nature* 522(7555):207–211.
103. Ralf A, Oven M, Zhong K, Kayser M (2015) Simultaneous analysis of hundreds of Y-chromosomal SNPs for high-resolution paternal lineage classification using targeted semiconductor sequencing. *Human mutation* 36(1):151–159.
104. Keller A et al. (2012) New insights into the Tyrolean Iceman’s origin and phenotype as inferred by whole-genome sequencing. *Nature communications* 3:698.
105. Haak W et al. (2010) Ancient DNA from European early Neolithic farmers reveals their near Eastern affinities. *PLOS biology* 8(11):e1000536.
106. Lacan M et al. (2011) Ancient DNA reveals male diffusion through the Neolithic Mediterranean route. *Proceedings of the National Academy of Sciences U S A* 108(24):9788–9791.
107. Lacan M et al. (2011) Ancient DNA suggests the leading role played by men in the Neolithic dissemination. *Proceedings of the National Academy of Sciences of the United States of America* 108(45):18255–9.
108. Battaglia V et al. (2009) Y-chromosomal evidence of the cultural diffusion of agriculture in Southeast Europe. *European Journal of Human Genetics* 17(6):820–830.
109. Semino O et al. (2000) The genetic legacy of Paleolithic Homo sapiens sapiens in extant Europeans: Y chromosome perspective. *Science* 290(5494):1155–1159.
110. Cinnioglu C et al. (2004) Excavating Y-chromosome haplotype strata in Anatolia. *Human genetics* 114(2):127–148.
111. Korneliusen TS, Albrechtsen A, Nielsen R (2014) ANGSD: analysis of next generation sequencing data. *BMC bioinformatics* 15(1):356.

112. Rasmussen M et al. (2011) An Aboriginal Australian genome reveals separate human dispersals into Asia. *Science* 334(6052):94–98.
113. Li H (2011) A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics* 27(21):2987–2993.
114. Briggs A, Stenzel U (2007) Patterns of damage in genomic DNA sequences from a Neandertal. *Proceedings of the National Academy of Sciences* 104(37):14616–14621.
115. Skoglund P et al. (2014) Separating endogenous ancient DNA from modern day contamination in a Siberian Neandertal. *Proceedings of the National Academy of Sciences* 111(6):2229–34.
116. McKenna A et al. (2010) The Genome Analysis ToolKit: A MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Research* 20(9):1297–1303.
117. Karolchik D et al. (2004) The UCSC Table Browser data retrieval tool. *Nucleic Acids Research* 32(suppl 1):D493–D496.
118. Felsenstein J (1981) Evolutionary trees from DNA sequences: a maximum likelihood approach. *Journal of Molecular Evolution* 17(6):368–376.
119. Wang C, Zhan X, Liang L, Abecasis GR, Lin X (2015) Improved ancestry estimation for both genotyping and sequencing data using projection procrustes analysis and genotype imputation. *American Journal of Human Genetics* 96(6):926–937.
120. Wang C, Zhan X (2015) *LASER: Locating Ancestry from SEquence Reads*, 2.02 edition. Available at [http://csg.sph.umich.edu/chaolong/LASER/LASER\\_Manual.pdf](http://csg.sph.umich.edu/chaolong/LASER/LASER_Manual.pdf).
121. Hellenthal G et al. (2014) A genetic atlas of human admixture history. *Science* 343(6172):747–751.
122. Busby GBJ et al. (2015) The role of recent admixture in forming the contemporary west Eurasian genomic landscape. *Current Biology*.
123. Purcell S et al. (2007) PLINK: a toolset for whole-genome association and population-based linkage analysis. *American Journal of Human Genetics* 81(3):559–575.
124. Chang CC et al. (2015) Second-generation PLINK: rising to the challenge of larger and richer datasets. *GigaSci* 4(1).
125. Hinrichs AS et al. (2006) The UCSC Genome Browser Database: update 2006. *Nucleic Acids Research* 34:590–8.
126. (2014) PyLiftOver (<https://github.com/konstantint/pyliftover>). Accessed: 2015-09-30.
127. Fu Q et al. (2014) Genome sequence of a 45,000-year-old modern human from western Siberia. *Nature* 514(7523):445–9.
128. Seguin-Orlando A et al. (2014) Genomic structure in Europeans dating back at least 36,200 years. *Science (New York, N.Y.)*.
129. Jones ER et al. (2015) Upper Palaeolithic genomes reveal deep roots of modern Eurasians. *Nature communications* 6:8912.
130. Mathieson I et al. (2015) Genome-wide patterns of selection in 230 ancient Eurasians. *Nature* 528(7583):499–503.
131. Olalde I et al. (2014) Derived immune and ancestral pigmentation alleles in a 7,000-year-old Mesolithic European. *Nature* 507(7491):225–228.
132. Skoglund P et al. (2014) Genomic diversity and admixture differs for Stone-Age Scandinavian foragers and farmers. *Science* 344(6185):747–750.

133. Omrak A et al. (2016) Genomic Evidence Establishes Anatolia as the Source of the European Neolithic Gene Pool. *Current Biology* 26(2):270–275.
134. Patterson N et al. (2012) Ancient admixture in human history. *Genetics* 192(3):1065–93.
135. Sánchez-Quinto F et al. (2012) Genomic affinities of two 7,000-year-old Iberian hunter-gatherers. *Current biology : CB* 22(16):1494–9.
136. Raghavan M et al. (2014) Upper Palaeolithic Siberian genome reveals dual ancestry of Native Americans. *Nature* 505(7481):87–91.
137. Gronau I, Hubisz MJ, Gulko B, Danko CG, Siepel A (2011) Bayesian inference of ancient human demography from individual genome sequences. *Nature Genetics* 43(10):1031–1034.
138. Veeramah KR et al. (2011) An early divergence of KhoeSan ancestors from those of other modern humans is supported by an ABC-based analysis of autosomal resequencing data. *Molecular Biology and Evolution* 29(2):617–630.
139. Alexander DH, Novembre J, Lange K (2009) Fast model-based estimation of ancestry in unrelated individuals. *Genome Research* 19(9):1655–1664.
140. Jakobsson M, Rosenberg NA (2007) CLUMPP: a cluster matching and permutation program for dealing with label switching and multimodality in analysis of population structure. *Bioinformatics* 23(14):1801–1806.
141. Rosenberg NA (2003) distruct: a program for the graphical display of population structure. *Molecular Ecology Notes* 4(1):137–138.
142. Yang MA, Slatkin M (2015) Using Ancient Samples in Projection Analysis. *G3 (Bethesda, Md.)* 6(January):g3.115.023788–.
143. Jeffreys H (1946) An invariant form for the prior probability in estimation problems. *Proceedings of the Royal Society of London. Series A: Mathematical and physical sciences* 186(1007):453–461.
144. Voight BF et al. (2005) Interrogating multiple aspects of variation in a full resequencing data set to infer human population size changes. *Proceedings of the National Academy of Sciences of the United States of America* 102(51):18508–13.
145. Bollongino R et al. (2012) Modern taurine cattle descended from small number of near-eastern founders. *Molecular Biology and Evolution* 29(9):2101–2104.
146. Hudson R, Slatkin M, Maddison W (1992) Estimation of levels of gene flow from DNA sequence data. *Genetics* 132:583–589.
147. Lawson DJ, Hellenthal G, Myers S, Falush D (2012) Inference of population structure using dense haplotype data. *PLoS genetics* 8(1):e1002453.
148. Leslie S et al. (2015) The fine scale genetic structure of the British population. *Nature* 519:309–314.
149. van Dorp L et al. (2015) Evidence for a Common Origin of Blacksmiths and Cultivators in the Ethiopian Ari within the Last 4500 years: Lessons for Clustering-Based Inference. *PLoS Genetics* 11(8):e1005397.
150. Busing F, Meijer E, Van Der Leeden R (1999) Delete-m Jackknife for Unequal m. *Statistics and Computing* 9:3–8.
151. Kirin M et al. (2010) Genomic runs of homozygosity record population history and consanguinity. *PloS one* 5(11):e13996.
152. Pemberton TJ et al. (2012) Genomic patterns of homozygosity in worldwide human populations. *American Journal of Human Genetics* 91(2):275–292.



153. The 1000 Genomes Project Consortium (2015) A global reference for human genetic variation. *Nature* 526(7571):68–74.
154. Prüfer K et al. (2014) The complete genome sequence of a Neanderthal from the Altai Mountains. *Nature* 505(7481):43–49.
155. Sudmant PH et al. (2015) An integrated map of structural variation in 2,504 human genomes. *Nature* 526(7571):75–81.
156. Rajeevan H, Soundararajan U, Kidd JR, Pakstis AJ, Kidd KK (2012) ALFRED: an allele frequency resource for research and teaching. *Nucleic Acids Research* 40(D1):1010–1015.
157. Sabeti PC et al. (2007) Genome-wide detection and characterization of positive selection in human populations. *Nature* 449(7164):913–8.
158. Canfield VA et al. (2013) Molecular phylogeography of a human autosomal skin color locus under natural selection. *G3 (Bethesda)* 3(11):2059–67.
159. Wilde S et al. (2014) Direct evidence for positive selection of skin, hair, and eye pigmentation in Europeans during the last 5,000 years. *Proceedings of the National Academy of Sciences* 111(13):4832–4837.
160. McVean G (2012) An integrated map of genetic variation from 1,092 human genomes. *Nature* 491(7422):56–65. 10.1038/nature11632.
161. Yuasa I et al. (2006) Distribution of the F374 allele of the SLC45A2 (MATP) gene and founder-haplotype analysis. *Annals of Human Genetics* 70(6):802–811.
162. Donnelly MP et al. (2012) A global view of the OCA2-HERC2 region and pigmentation. *Human genetics* 131(5):683–96.
163. Eiberg H et al. (2008) 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. *Human genetics* 123(2):177–187.
164. Walsh S et al. (2014) Developmental validation of the HIrisPlex system: DNA-based eye and hair colour prediction for forensic and anthropological usage. *Forensic Science International: Genetics* 9:150–161.
165. Bersaglieri T et al. (2004) Genetic signatures of strong recent positive selection at the lactase gene. *American Journal of Human Genetics* 74(6):1111–1120.
166. Burger J, Kirchner M, Bramanti B, Haak W, Thomas MG (2007) Absence of the lactase-persistence-associated allele in early Neolithic Europeans. *Proceedings of the National Academy of Sciences of the United States of America* 104(10):3736–41.
167. Allentoft ME et al. (2015) Population genomics of Bronze Age Eurasia. *Nature* 522(7555):167–172.
168. Malmström H et al. (2010) High frequency of lactose intolerance in a prehistoric hunter-gatherer population in northern Europe. *BMC evolutionary biology* 10(1):1.
169. Plantinga TS et al. (2012) Low prevalence of lactase persistence in Neolithic south-west Europe. *European journal of human genetics* 20(7):778–782.
170. Cassidy LM et al. (2016) Neolithic and Bronze Age migration to Ireland and establishment of the insular Atlantic genome. *Proceedings of the National Academy of Sciences* 113(2):368–373.
171. Raj T et al. (2013) Common risk alleles for inflammatory diseases are targets of recent positive selection. *The American Journal of Human Genetics* 92(4):517–529.

172. Helgason A et al. (2007) Refining the impact of TCF7L2 gene variants on type 2 diabetes and adaptive evolution. *Nature genetics* 39(2):218–225.
173. Huff CD et al. (2012) Crohn’s disease and genetic hitchhiking at IBD5. *Molecular biology and evolution* 29(1):101–111.
174. Peltekova VD et al. (2004) Functional variants of OCTN cation transporter genes are associated with Crohn disease. *Nature genetics* 36(5):471–475.
175. Peng Y et al. (2010) The ADH1B Arg47His polymorphism in east Asian populations and expansion of rice domestication in history. *BMC Evol Biol* 10:15.
176. Oota H et al. (2004) The evolution and population genetics of the ALDH2 locus: random genetic drift, selection, and low levels of recombination. *Annals of Human Genetics* 68(2):93–109.
177. Kuznetsov IB, McDuffie M, Moslehi R (2009) A web server for inferring the human N-acetyltransferase-2 (NAT2) enzymatic phenotype from NAT2 genotype. *Bioinformatics* 25(9):1185–6.
178. Patillon B et al. (2014) A homogenizing process of selection has maintained an “ultra-slow” acetylation NAT2 variant in humans. *Human Biology* 86(3):185–214.
179. Relethford J, Sabbagh A, Darlu P, Crouau-Roy B, Poloni ES (2011) Arylamine N-acetyltransferase 2 (NAT2) genetic diversity and traditional subsistence: A worldwide population survey. *PLoS ONE* 6(4):e18507.
180. Magalon H et al. (2008) Population genetic diversity of the NAT2 gene supports a role of acetylation in human adaptation to farming in Central Asia. *European Journal of Human Genetics* 16(2):243–51.
181. García-Closas M et al. (2011) A single nucleotide polymorphism tags variation in the arylamine N-acetyltransferase 2 phenotype in populations of European background. *Pharmacogenetics and genomics* 21(4):231.
182. Selinski S et al. (2011) Genotyping NAT2 with only two SNPs (rs1041983 and rs1801280) outperforms the tagging SNP rs1495741 and is equivalent to the conventional 7-SNP NAT2 genotype. *Pharmacogenetics and Genomics* 21(10):673–678.
183. Cao A, Galanello R (2010) Beta-thalassemia. *Genetics in Medicine* 12(2):61–76.
184. Tadmouri G et al. (1998) Molecular and population genetic analyses of  $\beta$ -Thalassemia in Turkey. *American journal of hematology* 57(3):215–220.
185. Itan Y, Jones BL, Ingram CJ, Swallow DM, Thomas MG (2010) A worldwide correlation of lactase persistence phenotype and genotypes. *BMC evolutionary biology* 10(1):1.