

1 *Research paper*

2 Ancient ancestry informative markers for identifying fine-scale ancient population

3 structure in Eurasians

4

5 Umberto Esposito¹, Ranajit Das², Mehdi Pirooznia³, and Eran Elhaik^{1*}

6

7 ¹ Department of Animal and Plant Sciences, University of Sheffield, Sheffield, UK S10

8 2TN

9 ² Manipal University, Manipal Centre for Natural Sciences (MCNS), Manipal, Karnataka,

10 India

11 ³ Johns Hopkins University, Department of Psychiatry and Behavioral Sciences,

12 Baltimore, MD, USA 21205

13

14

15

16

17

18

19

20 * Please address all correspondence to Eran Elhaik at e.elhaik@sheffield.ac.uk

21

22 **Keywords:** ancient DNA, ancient ancestry informative markers, population structure,

23 PCA, admixture mapping

24 **Abstract**

25

26 The rapid accumulation of ancient human genomes from various places and time periods,
27 mainly from the past 15,000 years, allows us to probe the past with an unparalleled
28 accuracy and reconstruct trends in human biodiversity. Alongside providing novel
29 insights into the population history, population structure permits correcting for population
30 stratification, a practical concern in gene mapping in association studies. However, it
31 remains unclear which markers best capture ancient population structure as not all
32 markers are equally informative. Moreover, the high missingness rates in ancient,
33 oftentimes haploid, DNA, may distort the population structure and prohibit genomic
34 comparisons. In past studies, ancestry informative markers (AIMs) were harnessed to
35 address such problems, yet whether AIMs finding methods are applicable to aDNA
36 remains unclear. Here, we define ancient AIM (aAIMs) and develop a framework to
37 evaluate established and novel AIMs-finding methods. We show that a novel principal
38 component analysis (PCA)-based method outperforms all methods in capturing ancient
39 population structure and identifying admixed individuals. Our results highlight important
40 features of the genetic structure of ancient Eurasians and the choice of strategies to
41 identify informative markers. This work can inform the design and interpretation of
42 population and medical studies employing ancient DNA.

43

44

45

46 **Author summary**

47 Ancient DNA studies aim to identify geographical origin, migration routes, and disease
48 susceptibility genes through the analysis of genetic markers such as single nucleotide
49 polymorphisms (SNPs) in growing cohorts of ancient data. In addition to the existence of
50 sub-structure in the studied population (i.e., differences in ancestry), ancient DNA suffers
51 from high missingness rates and is oftentimes haploid, which may distort the inferred
52 population structure and lead to spurious results. It is thereby imperative to address this
53 possible bias by identifying the most accurate population structure. Due to the success of
54 past studies in addressing similar problems using ancestry informative markers (AIMs),
55 we defined ancient ancestry informative markers (aAIMs) that like AIMs can be used to
56 interrogate ancient population structure. To find aAIMs, we designed a framework to
57 evaluate established and novel AIMs-finding methods. We developed a database of
58 150,278 autosomal SNPs from 302 ancient genomes and 21 populations recovered from
59 Europe, the Middle East, and North Eurasia dated to time periods from 14,000 to 1,500
60 years ago. We then applied two existing and three novel AIMs-finding methods and
61 compared their performances against the complete dataset. We found that a novel
62 principal component analysis (PCA)-based method captured the ancient population
63 structure most accurately. Importantly, we introduce here a novel concept of aAIMs, a
64 novel method that effectively identifies aAIMs, and a framework to compare the
65 performances of AIMs. The outcome of our studies can improve the accuracy of genetic
66 studies employing ancient DNA.

67

68 **Introduction**

69

70 Population stratification or geographic variation are a major concern in population,
71 biomedical, and evolutionary studies. In genetic association studies, mismatching cases
72 and controls introduces genetic heterogeneity that can lead to spurious associations and
73 obscure the true association [1, 2]. In large groups the stratification bias may be less
74 pronounced, but it is practically unavoidable in the case of rare diseases due to the
75 difficulties in recruiting genetically homogeneous participants [3]. These problems are
76 particularly challenging since the human population structure itself remain contentious.
77 Nonetheless, is now clear that conquering population structure requires considering
78 ancient DNA (aDNA) [4-6].

79

80 The advent of next-generation sequencing and the availability of large-scale genomic
81 data and genotyping techniques have facilitated investigations of genomic variability that
82 are central to understanding our evolutionary history and genomic origins. Over the last
83 decade a plethora of ancient human genome sequencing projects have been
84 accomplished, generating more than a thousand ancient genomes [7]. The revolution in
85 aDNA sequencing has aided in investigations of ancient human migration, human
86 adaptation, agricultural lifestyle, and disease co-evolution [7]. Notwithstanding its
87 usefulness in delineating the evolutionary history of mankind, aDNA data can be
88 problematic due to its haploidity and high missingness [6], which require having a large
89 number of SNPs to infer population structure. However, SNPs are not equally
90 informative and may distort the population structure. The plethora of mismatching

91 markers sequenced in different genomes has revived the “of AIMs matrix” problem and
92 the difficulty of comparing genomes. These problems are not new and rather reminiscent
93 of the early stages of human population genetics. Then, one of the most successful
94 solutions was using ancestry informative markers (AIMs).
95
96 AIMs are SNPs which exhibit large variation in minor allele frequencies (MAF) among
97 populations. Over the past two centuries DNA studies scour genomes for these genetic
98 patterns and produced numerous AIM sets for various purposes including determining an
99 individual’s ancestry, detecting stratification in biomedical studies, inferring geographic
100 structure, and localizing biogeographical origins [e.g., 8, 9-12]. AIM panels can delineate
101 population structure in a cost effective manner by identifying population specific
102 markers, which in turn help in detecting and correcting for variation in individual
103 ancestry that can confound methods like admixture mapping, Mendelian Randomization
104 trials, association studies, and forensics by increasing false positive results and/or
105 reducing Power [e.g., 13, 14, 15]. In the case of genetic association studies, AIMs-based
106 solution has been preferred over methods like genomic control (GC) correction, which is
107 only applicable in genome-wide scale data [16]. However, it remains uncertain which
108 AIMs to use since all AIMs panels have limitations [17] and their applicability to ancient
109 genomes was never tested. The characteristics of ideal AIMs are remain contentious with
110 some authors preferring common SNPs (minor allele frequency >1%) [16], SNPs with
111 high F_{ST} [18], SNPs with high pairwise MAF between populations [17], or SNPs that
112 satisfy several criteria. Consequently, AIMs may not overlap across studies that focus on
113 particular populations and even those reported in global studies do not necessarily

114 overlap. Finally, studies typically show that AIMS can separate populations or broadly
115 classify individuals into subcontinental populations, rather than capture the population
116 structure of the complete SNP set or allow fine-population mapping. Given the
117 uncertainties surrounding AIMS, their potential incompatibility to capture ancient
118 structure and admixtures, and the challenges imposed by aDNA data, it is unclear
119 whether, if at all, AIMS-finding methods or AIMS can be utilized to study ancient
120 population structure.

121

122 In this study, we defined ancient ancestry informative markers (aAIMs) as SNPs that vary
123 in their MAF across ancient populations (Figure 1) and attempted to identify and validate
124 the first autosomal aAIMs. Since AIMS-finding tools were never tested on aDNA, it was
125 necessary to compare their ability in finding aAIMs. For that, we interrogated a
126 comprehensive dataset of 302 ancient genomes classified to 21 populations from Europe,
127 the Middle East, and North Eurasia. This dataset was used to compare two existing AIMS
128 finding algorithms: Infocalc [19] and Wright's F_{ST} [20, 21], three novel Admixture- and
129 PCA-based algorithms, and two random SNP sets in identifying aAIMs that can capture
130 the population structure and identify admixed individuals. Our study offers a
131 methodological framework to evaluate AIMS, contrasts different strategies to find aAIMs,
132 and reports the first set of aAIMs.

133

134

135

136

137 **Results**

138

139 **Ancient genomic data**

140

141 We developed a framework to identify and evaluate the efficacy of aAIM candidates in
142 capturing ancient population structure and allow admixture mapping (Figure 2). We
143 constructed a dataset of 150,278 autosomal SNPs from 302 ancient genomes and 21
144 populations recovered from Europe, the Middle East, and North Eurasia dated to time
145 periods spanning the past 14,000 years till 1,500 years ago (Figure 3, Table S1). Due to
146 the limited availability of ancient genomes, our dataset was inconsistent over time and
147 space. For instance, there were 57 Central European genomes from the Late Neolithic to
148 the Bronze Age, but populations, such as, Central Western Mesolithic Europeans, Bronze
149 Age Jordanians, Chalcolithic, and Mesolithic Russians, comprised of three genomes each.

150

151 Missingness largely varied between the samples and markers. The sample-based
152 missingness ranged from 0.05% (KK1) to 99.2% (I1951) with an average of 54%. The
153 sample-based missingness varied among populations with Levantine Epipaleolithic
154 Neolithic genomes having the highest missingness ($n=19$, $\mu=90\pm 16\%$) and Mesolithic
155 Swedish genomes having the lowest one ($n=8$, $\mu=29\pm 27\%$). The variant-based
156 missingness ranged from 30% to 98% with an average of 54%.

157

158 Principal component analysis (PCA) of the ancient genomes substantiated previous
159 observations of a Europe–Middle East contrast along the vertical principal component

160 (PC1) and parallel clines (PC2) in both Europe and the Middle East (Figure 4). Genomes
161 from the Epipaleolithic and Neolithic Levantine clustered at one extreme of the Near
162 East-Europe cline with some overlapping with Neolithic Turkish and Central European
163 genomes. Neolithic Iranians clustered between Central European genomes. While ancient
164 Spanish, Armenian, Central EU, and British genomes occupied the intermediate position
165 of Near Eastern and North Eurasian genomes, Russian and Swedish genomes clustered at
166 the end of the Near East-Europe cline.

167

168 We next applied an unsupervised ADMIXTURE analysis to the dataset. Analyzing the
169 results generated with various number of splits (K) (Figure S1), no choice of K
170 minimized the cross-validation error (CVE) (Figure S2), likely because the high noise
171 and missingness in the data prevented the CVE from stabilizing. We observed that at
172 $K=10$, multiple genomes (e.g., Britain Iron Saxon, Mesolithic Neolithic Caucasus
173 population, Bronze Age Jordanian and Epipaleolithic Levantine, Chalcolithic, Mesolithic
174 and Early Mid Bronze Russian, Early Neolithic Spanish, Mesolithic and Mid Neolithic
175 Swedish, and Neolithic Turkish) appeared homogeneous in relation to their population
176 and assigned to a distinct allele frequency profile or admixture components (Figure 5). In
177 these figure, putative ancient ancestral components, like the *Early Neolithic European*
178 (brown) and *Russia Mid Late Bronze* (magenta), predominantly found among European
179 genomes, may be identified. Except their predominance in Neolithic Turkish genomes,
180 these components also exist in most Neolithic Central Europeans. Some 20-30% of
181 Central European genomes have discernible fractions of *Europe Late Neolithic-Early*
182 *Bronze* (navy-blue) and *Russia Mid-Late Bronze* (deep-pink) components, respectively.

183 Two components (cyan and dark purple) appeared sporadically in a few populations,
184 likely due to noise.

185

186 Identifying aAIM candidates

187

188 To identify aAIM candidates, we employed Infocalc and F_{ST} , commonly used to detect
189 AIMS. We also implemented three novel Admixture- and PCA-based methods
190 (Admixture₁, Admixture₂, and PCA-derived [PD]). Finally, we selected two random SNP
191 sets of 10,000 and 15,000 markers, which approximated the number of AIMS identified
192 by the various methods (Rand_{10k} and Rand_{15k}). Four criteria were adopted to evaluate
193 how the candidate aAIMs capture the population structure depicted by the complete SNP
194 set (CSS): first, by qualitatively comparing the dispersal of genomes obtained from a
195 PCA to that of the CSS. Second, by comparing the Euclidean distances between the
196 admixture proportions of each genome and those obtained from the CSS. To avoid
197 inconsistencies between the SNP sets, we used admixture components obtained through a
198 *supervised* ADMIXTURE (see *methods*). Third, by testing which aAIMs classify
199 individuals to populations most accurately. The abilities to identify admixed individuals
200 and evaluated for the top performing method.

201 As with the CSS, genomes with over 90% missingness were removed, leaving each
202 dataset with 223-263 genomes (Table S2). 310 SNPs without data were removed from the
203 Rand_{10k} dataset. The final number of aAIM candidates identified using each method is
204 shown in Table S3. Overlapping aAIMs between the methods are remarkably small and
205 range from 560 (Rand_{10k} and Admixture₁) to 2,160 (Admixture₁ and Admixture₂).

206 Interestingly, Infocalc and F_{ST} , oftentimes used together share only ~10% of their aAIM
207 candidates. The PD method shares 13.7% of its aAIMs with F_{ST} and ~10% with Infocalc.

208 Comparing the sequence properties of the aAIM candidates, we found that for ancient
209 populations (Figure 6a) Infocalc's aAIMs mirrored the MAF of the CSS with most
210 variants having low MAF (45% of the aAIMs have $MAF < 0.1$). The F_{ST} aAIM also had
211 high frequency of low-mid MAF values. By contrast, the PD and Admixture-based
212 methods exhibited higher frequencies of high MAF SNPs with Admixture₂ having the
213 highest proportion of high MAF aAIMs (75% of the aAIMs have $MAF > 0.4$).

214 Interestingly, the MAF distributions exhibited similar distributions in modern populations
215 (Figure 6b), though with fewer alleles in lowest MAF bins for all the methods.

216 Unsurprisingly, most of the aAIM variants were non-functional (94.6-96.3%) and vary
217 little from the CSS's annotation (Table S4).

218

219 Comparative testing of aAIM candidates

220

221 We compared the performances of aAIMs candidates to each other and to the CSS in
222 capturing the population structure and classifying individuals to populations through
223 three analyses. First, we calculated the PCA for each SNP set and compared the
224 population dispersion along the primary two axis. Similarly to the CSS (Figure 4), all the
225 methods depicted the Europe–Middle East contrast (PC1) and parallel clines (PC2) in the
226 European genomes so that ancient Jordanian, Levantine, Turkic, and Spanish genomes
227 clustered at one extreme of the Near East-Europe cline, whereas the genomes from
228 Russia and Sweden clustered at the other end (Figure S3). However, much like the

229 random sets, Infocalc and F_{ST} did not separate Levantine and Turkish individuals from
230 each other. Infocalc also merged the Caucasus individuals with central Europeans. The
231 admixture-based methods and PD separated all the ancient populations, similarly to the
232 CSS and better, in the case of Scandinavians and Russians.

233

234 We next quantitatively assessed which dataset produced the closest admixture signature
235 to that of the CSS (Figures 5). For that, we calculated the admixture proportions in
236 relation to ten putatively ancient ancestral populations (Figures S4-5) and then computed
237 their Euclidean distances to their counterparts obtained by the CSS (Figure 7). The PD
238 aAIMs had significantly short Euclidean distances ($\mu=0.13$, $\sigma=0.1$, $n=302$) compared to
239 all other aAIMs (Welch t -test p -values: Infocalc 0.002, F_{ST} 8.5×10^{-13} , Admixture₁ 2.2×10^{-16} ,
240 Admixture₁ 2×10^{-16} , Rand_{10k} 5×10^{-6} , and Rand_{15k} 0.001). Infocalc's aAIMs produced
241 the second shortest distances from the CSS ($\mu=0.17$, $\sigma=0.15$), however they were not
242 statistically shorter than the distances obtained by the two random datasets (Welch t -test
243 p -values: Rand_{10k} 0.12 and Rand_{15k} 0.77 respectively), suggesting that Infocalc was
244 unable to capture the population structure. F_{ST} -derived AIMs ($\mu=0.2$, $\sigma=0.13$) performed
245 worse than the Rand_{15k} aAIMs (Welch t -test p -value 0.004), and similarly to the Rand_{10k}
246 aAIMs (Welch t -test, p -value=0.13). The admixture-based datasets performed worst of all
247 aAIMs ($\mu_1=0.22$, $\sigma_1=0.15$ and $\mu_2=0.24$, $\sigma_2=0.16$) and significantly worse than the two
248 random datasets (Welch t -test: Admixture₁ [Rand_{10k} p -value=0.002] and [Rand_{15k} p -
249 value= 1.6×10^{-5}]; Admixture₂ [Rand_{10k} p -value= 1.7×10^{-5}] and [Rand_{15k} p -value= 2.5×10^{-8}]).

250

251 We last assessed which aAIMs dataset allows classifying individuals to population
252 groups most accurately. For that, an admixture-based population classifier was applied to
253 the admixture proportions produced by all the datasets and their accuracy was compared
254 to that of the CSS ($76\pm 25\%$) and the known population classification ([Table S1](#)). The
255 mean classification accuracy per population ranged from 3% (F_{ST}) to 61% (PD) with the
256 PD outperforming all other methods ([Table 1](#)). In other words, ~13k (8%) of the SNPs
257 are sufficiently informative to classify individuals to populations with 80% of the
258 accuracy of the CSS. In nine out of 21 population groups (22% of the individuals) PD-
259 based classification was similar or more accurate than the CSS. All other methods
260 performed similarly or worse than the random SNP sets ($42\pm 22\%$ and $50\pm 23\%$) with
261 Infocalc ($50\pm 23\%$) outperforming the remaining methods. Of note are the poor
262 performances of F_{ST} aAIMs, likely due to the high sensitivity of F_{ST} to aDNA data. As
263 expected, high missingness was associated with incorrect predictions ([Figure S6](#)). For
264 example, the low-coverage low-quality Britain Anglo-Saxon genomes proved
265 challenging for all the methods (0-40%) but predicted correctly with the CSS (100%).
266 Due to the high accuracy of the PD aAIMs compared to the alternative datasets, we
267 continued to analyze its aAIMs,

268

269 [Inference of admixed samples](#)

270

271 Admixture mapping is a powerful method of gene mapping to map phenotypic variation
272 or diseases that show differential risk by ancestry and takes advantage of higher densities
273 of genetic variants and extensions to admixed populations [22]. Thereby a large number

274 of markers throughout the genome is necessary to allow inference of local chromosomal
275 ancestry blocks. [Figure 8](#) illustrates the genome-wide distribution of PD aAIMs. To test
276 whether these aAIMs can identify admixture in hybrid individuals, ancient individuals
277 were hybridized to form 120 mixed individuals, each associated with three datasets: CSS,
278 PD aAIMs, and a random SNP set of the size of PD aAIMs ([Table 2](#)).

279

280 The genetic distances between the CSS and PD aAIMs were significantly smaller
281 ($\mu=0.05$, $\sigma=0.04$) than the distances between the CSS and the random SNP sets ($\mu=0.45$,
282 $\sigma=0.15$, Welch *t*-test *p*-values= 2.2×10^{-8}) as well as between the OD and the random SNP
283 sets ($\mu=0.43$, $\sigma=0.15$, Welch *t*-test *p*-values= 1.9×10^{-8}). We, thus, demonstrated that PD
284 aAIMs can be used to infer admixed individuals and be used in future admixture mapping
285 involving aDNA.

286

287

288

289

290 **Discussion**

291

292 The use of ancient genomes in research is at its infancy and expected to intensify as data
293 are becoming available. It is reasonable to expect that many of the tools employed to
294 study modern-day genomes will need to be adapted to the ancient DNA environment.
295 Some of the most useful tools in addressing population, biomedical, and evolutionary
296 questions were ancestry informative markers (AIMs), however it is unclear whether they

297 are applicable to ancient genomic data, which not only represent populations with
298 different population structure, but has some unique characteristics like high missingness
299 and haploid genomes [6].

300

301 In this study, we defined ancient ancestry informative markers (aAIMs) (Figure 1) and
302 sought to identify those using various methods. The number of aAIMs identified by each
303 method ranges from 9 to 15 thousands. These numbers of the same magnitude of large
304 AIMs studies [e.g., 23, 24] and reasonable provided the potential relatedness of the
305 ancient populations and the near absence of heterozygote markers in the data. To find
306 which of the aAIMs candidates produced by each method best represent the true
307 population structure, we used the complete SNP set as a benchmark for qualitatively and
308 quantitatively comparisons.

309

310 Identifying the ideal AIMs set that would be both small and include redundancies (in case
311 of sequencing failure), capture the population structure, and allow identifying admixed
312 individuals remains one of the challenges of population genetics. We showed that aAIMs
313 identified through a PCA-derived (PD) method outperformed all other methods in
314 agreement with previous studies that tested PCA-based methods [16]. Some
315 classifications made by the PD were more accurate than those made using the CSS,
316 which highlights the negative influence of ancestry *uninformative* markers. To the best of
317 our knowledge, such markers and their influence were never explored. Infocalc and F_{ST}
318 aAIMs, typically used in conjunction to identify AIMs [10] and have been reported to
319 perform well in admixed populations [25] have oftentimes underperformed random

320 SNPs. Not only was F_{ST} already shown to be particularly small within continental
321 populations [26], but these methods may be particularly sensitive to ancient DNA data
322 that is both haploid and has high missingness (Figure S6). We also found no relationships
323 between MAF and aAIMs performances (Figure 5). Enrichment for high or low MAF
324 SNPs did not guarantee success, although the PD harbored more common SNPs than
325 most underperforming methods.

326

327 The applicability of the PD aAIMs for admixture mapping combined with tools that can
328 homogenize cases and controls [e.g., 27] enable future association studies to be carried
329 out on ancient DNA samples. Indeed, Cassidy et al. [28] provided evidence for the
330 existence of Hemochromatosis alleles in ancient genomes and point at the association of
331 hemochromatosis alleles in ancient Irish. Due to the nature of the ancient data and to
332 enable admixture mapping studies we refrained from optimizing the number of aAIMs.
333 Further investigations with additional data may identify formerly common markers
334 associated with the disease that with time became rare and undetectable.

335

336 Our study has several limitations. We studied an uneven number of Eurasian populations
337 from various times and locations, causing a skew towards markers that predict central
338 European populations from the Late-Neolithic Bronze Age. A modest attempt to reduce
339 this bias was made by including modern-day African and Asian populations, however a
340 more comprehensive analyses should be made when more global genomes are available.
341 Second, the aAIMs were calculated independently by each method with individual
342 populations considered independent, although the PCA and ADMIXTURE plots indicate

343 that central European populations may not be independent. Finally, due to high
344 missingness of the data, it is likely that our study missed informative markers that could
345 improve the classification accuracy in newly sequenced populations. We thereby advise
346 applying our method to more comprehensive aDNA datasets when such will be available.

347

348 In summary, AIMs are some of the most effective tools that spear-headed population
349 genetics over the past two decades and ancillary to the challenge of understanding
350 population structure. We defined ancient AIMs (aAIMs), proposed a framework to
351 evaluate AIMs-finding methods, demonstrated the accuracy of a novel aAIMs-finding
352 method, and reported the most successful set of aAIMs. Future analyses may benefit from
353 using our method to uncover powerful aAIMs and using our aAIMs to refine ancient
354 population structure models.

355

356 **Methods**

357

358 **Sample collection**

359

360 Ancient DNA genomic data were obtained from 11 publications depicting 207 ancient
361 genomes ([Table S1](#)). In the case of sequence data, sequence reads were aligned to the
362 human reference assembly (UCSC hg19-<http://genome.ucsc.edu/>) using the Burrows
363 Wheeler Aligner (BWA version 0.7.15) [29], allowing two mismatches in the 30-base
364 seed. Alignments were then imported to binary (bam) format, sorted, and indexed using
365 SAMtools (version 1.3.1) [30]. Picard (version 2.1.1) (<http://picard.sourceforge.net/>) was

366 then used for MarkDuplicates to remove reads with identical outer mapping coordinates
367 (which are likely PCA artifacts). The Genome Analysis Toolkit RealignerTargetCreator
368 module (GATK version 3.6) [31, 32] was used to generate SNP and small InDel calls for
369 the data within the targeted regions only. GATK InDelRealigner/BaseRecalibrator was
370 then used for local read realignment around known InDels and for base quality score
371 recalibration of predicted variant sites based on dbSNP 138 and 1000 Genomes known
372 sites, resulting in corrections for base reported quality. The recalibration was followed by
373 SNP/InDel calling with the GATK HaplotypeCaller. Variants were filtered for a
374 minimum confidence score of 30 and minimum mapping quality of 40. At the genotype
375 level, all genotypes that had a genotype depth less than 4 ($GD < 4$) or a genotype quality
376 score less than 10 ($GQ < 10$) were removed from the dataset by setting them to missing in
377 the VCF. GATK DepthofCoverage was then used to re-examine coverage following the
378 realignment. VCFtools (version 0.1.14) [33] were used to convert the VCF file to PLINK
379 format [34]. The final dataset comprised of 150,278 autosomal SNPs from 302 ancient
380 DNA (aDNA) genomes ([Table S1; Additional file 1](#)). Eight aDNA genomes (I0247,
381 I1584, I1955, ATP9, IR1, Kostenki14, MA1, and Ust_Ishim) without any country/region
382 designation were omitted in the closest population determination calculations. For
383 coherency, the genomes were divided into 21 populations, based on the sampling
384 country/region and their era.

385

386 [Data analyses](#)

387

388 [The genetic structure canvas of ancient Eurasian genomes](#). The population structure of
389 the ancient genomes was described using principal component analysis (PCA)
390 implemented in PLINK v1.9 (<https://www.cog-genomics.org/plink/1.9/>). Individuals with
391 high SNP missingness were removed using --mind 0.9 flag alongside the --pca command
392 for all the aAIMs datasets. We also applied the model-based clustering methods
393 implemented in ADMIXTURE v1.3 [35]. All PCA and Admixture plots were generated
394 in R v3.2.3. Minor allele frequency (MAF) was calculated using PLINK (--maf
395 command) for ancient populations and for modern ones, MAF was calculated from the
396 1000 Genomes populations (ALL.2of4intersection.20100804.genotypes) [36]. Percentage
397 of rare and novel variants and other functional information were obtained through VEP
398 (McLaren et al. 2016).

399

400 [Identifying aAIMs via five methods](#). aAIMs were considered markers that can infer the
401 ancestry of ancient DNA (aDNA) genomes in a similar accuracy to the complete SNP set
402 (CSS). We compared methods to detect candidate AIMs, three of which are novel:

403

404 **1.** Infocalc (Rosenberg et al. 2003), which determines the amount of information
405 multiallelic markers provide about an individual's ancestry by calculating the
406 informativeness (I) of each of each markers separately and ranks the SNPs by
407 their informativeness. Infocalc determines I based on the mathematical expression
408 described in Rosenberg et al. (2003). We compared the performances of four
409 choices of the top 5,000, 10,000, 15,000, and 20,000 most informative markers in

410 the Infocalc v1.1 output file (results not shown). The 15,000 dataset outperformed
411 all other datasets and was selected for further analyses.

412 **2.** *F_{ST}*. Wright's fixation indices (*F_{ST}*) [21] measures the degree of differentiation
413 among populations potentially arising due to genetic structure within populations.
414 Given a set of populations (Table S1), we employed PLINK [34] to estimate *F_{ST}*
415 separately for all the markers using `-fst` command alongside `--within` flag, that
416 defines population IDs of the genomes. Due to the high fragmentation of the data,
417 *F_{ST}* values could only be calculated for 46% of the dataset. We compared the
418 performances of four choices of the highest 5,000, 10,000, 15,000, and 20,000 *F_{st}*
419 values. The 15,000 dataset outperformed all other datasets and was selected for
420 further analyses.

421 **3.** Admixture₁. This method assumes that AIMs have high allelic frequencies in
422 certain subpopulations and drive the differentiation of admixture components.
423 Analyzing ADMIXTURE's output file (P file) for *K* of 10, we identified the
424 markers (rows) that had high allele frequency (>0.9) in only one admixture
425 component (columns). We identified 9,309 from the five columns with the highest
426 number of such markers.

427 **4.** Admixture₂. This method assumes that AIMs embody both high allelic
428 frequencies in certain subpopulations and high variance between these allelic
429 frequencies that differentiate of admixture components. Analyzing
430 ADMIXTURE's output file (P file) for *K* of 10, we identified 11,418 SNPs that
431 for each SNP (rows) had high variance (≥ 0.04) and high allele frequency range
432 (maxima - minima ≥ 0.65) between the admixture component (columns).

433 **5. PC-based (PD) approach.** This methods assumes that AIMs can replicate the
434 population structure of subpopulations represented by the variation in the first two
435 PCs. This is an interactive PC-based approach that identifies the smallest set of
436 markers necessary to capture the population structure of a group of individuals as
437 captured by the CSS. More specifically, for each population group ([Table S1](#)) in
438 which at least 100 SNPs were available, we calculated PCA and used PC1 and
439 PC2 to plot the individuals after all SNPs with high missingness (>0.05) were
440 excised. If the population group had insufficient SNPs we relaxed the missingness
441 threshold by additional 0.05, though 0.05 were sufficient for almost all groups.
442 We then scored the SNPs by their informativeness as in [37] and visually
443 compared the plot to that obtained from the CSS ([Figure S7](#)). If the plots were
444 dissimilar, we repeated the analysis using additional 100 top scored SNPs until
445 either the plots exhibited high similarity or a threshold of 2000 SNPs was reached.
446 We were unable to complete the analyses for 3 populations due to the small
447 number of individuals. The PD method is available on
448 <https://github.com/eelhaik/PCA-derived-aAIMs>. On average 861 SNPs were
449 found per population group. Overall, the dataset comprised of 13,027 SNPs.

450

451 To compare the prediction accuracy of the aAIMs subsets, two control datasets (Rand_{10k}
452 and Rand_{15k}) were generated by randomly sampling 10,000 and 15,000 SNPs from the
453 CSS, respectively. aAIMs identified by all methods are available as [Additional file 2](#).

454

455 **Classifying individuals to populations from genomic data. Identifying ancient**
456 **admixture components.** We selected a random hundred ancient genomes and removed
457 six for insufficient data (>95% missingness). To those, we added 20 Han Chinese and 20
458 Yoruba modern genomes from the 1000 Genomes Project (Durbin et al. 2010). We then
459 applied *supervised* ADMIXTURE with various K 's ranging from 8 to 13. While we were
460 unable to find a single K where culturally related genomes exhibited homogeneous
461 admixture patterns, the most robust population substructure was found for K of 10. Two
462 more components were obtained by analyzing Spanish and German genomes that
463 appeared indistinguishable along with five Yoruba genomes separately. We observed
464 very little admixture with the Han and Yoruba. Overall, we identified 10 admixture
465 components in ancient genomes, corresponding to allele frequencies of hypothetical
466 populations. Similarly to Elhaik et al. [9], we simulated 15 samples for each hypothetical
467 population, by generating 30 alleles whose average corresponds to the mean allele
468 frequency of that population and assigning those genotypes to the simulated individuals.
469 The putative ancestral ancient populations are available in [Additional file 3. Relabeling](#)
470 **populations.** Initially, the labels from the corresponding papers were used to classify
471 individuals to population. The consistency of these labels with data was evaluated by
472 carrying out a *supervised* ADMIXTURE analysis on the genomic data combined with the
473 150 putative ancient ancestral individuals. Due to the high similarity of the admixture
474 patterns between individuals of different groups living in similar periods or entire groups
475 (e.g., Neolithic individuals from Hungary and those from Germany), we re-labeled some
476 of the population to reduce the number of populations and create more genomically
477 homogeneous populations,. For instance, Natufian and Neolithic samples from Jordan are

478 grouped into the label Levant Epipaleolithic Neolithic. Overall, we identified 23
479 populations, whose labels are all of the form “area_time period.” In the case of the
480 Caucasus labelling, all the samples from Iran (except Iran_HotuIIIb) were excavated in
481 the Zagros Mountains, south of the Caucasus. Given their admixture similarity with
482 Armenians and Georgians from the same periods and their proximity to the Caucasus,
483 this area was labelled as Caucasus. Iran_HotuIIIb was found in a more eastern region,
484 just below the southeastern edge of the Caspian Sea, and given its similarity to Georgians
485 and other Iranians it was included in the group Caucasus Mesolithic Neolithic.

486 **Genomically defining reference populations.** For each population with $N_p > 4$, where
487 N_p is the number of individuals in the population, individuals were grouped in clusters
488 through an iterative process that uses a k -means clustering technique paired with multiple
489 pairwise F -tests. Iterations ran over the number of k clusters $[2, N_p/2]$. At each iteration i ,
490 k -means was used to identify the k clusters, then the F -test was applied on each pair of
491 clusters to test whether they are significantly ($P < 0.05$) different. If the two clusters are
492 different from all the pairs at iteration i , the process advances to $i+1$ until at least one pair
493 violates the condition, in which case $k_{op}=i-1$ is the optimal number of clusters or
494 reference populations within that population. **Assigning individuals to populations.** We
495 developed an admixture-based classifier, which is not sensitive to exclusion of random
496 groups of individuals nor inclusion of large numbers of individuals from admixed groups
497 and was trained on a third of the data. Using *supervised* ADMIXTURE, we calculated the
498 admixture proportions of the individuals in relation to the putative ancient ancestral
499 populations. Population assignment was then made based on the minimal Euclidean
500 distance between the admixture components of each genome and those of the reference

501 populations. The assignment accuracy was calculated based on the known classification
502 ([Table S1](#)).

503

504 [Assessing admixture mapping](#). **Creating hybrid individuals**. We selected 15 individuals
505 from five populations that showed homogeneity in their admixture components ([Figure 5](#))
506 and randomly sampled 120 pairs. Since selecting random alleles from each parent was
507 problematic due to the high missingness of the data, we randomly selected half the
508 genotypes of each parent to form 120 “offspring” or hybrid genomes. Each hybrid had
509 three SNP sets: the CSS, PD aAIMs, and a random SNP set of the size of PD aAIMs with
510 SNPs selected randomly for each hybrid. **Assessing admixture accuracy**. We defined
511 genetic distances (d) as the Euclidean distance between two set of admixture proportions.
512 We applied a *supervised* admixture to the three SNP sets of each hybrid and calculated
513 their distances d from each another.

514

515 [Graphics](#). Maps were drawn using the ‘*rworldmap*’ package implemented in R v3.2.3.

516

517 [Availability of data and materials](#). The dataset supporting the conclusions of this article
518 are included within the article and its additional files.

519

520 **Acknowledgment**

521 This study was partially supported by the MRC Confidence in Concept Scheme award
522 2014-University of Sheffield to E.E. (Ref: MC_PC_14115). We thank Grace Holland

523 who was partially supported by the UK EPSRC Doctoral Training Partnership Grant
524 EP/N509735/1 as a Vacation Bursary Training Project.

525

526 **References**

527

- 528 1. Marchini J, Cardon LR, Phillips MS, Donnelly P: **The effects of human**
529 **population structure on large genetic association studies.** *Nat Genet*
530 2004, **36**:512-517.
- 531 2. Scannell JW, Blanckley A, Boldon H, Warrington B: **Diagnosing the decline**
532 **in pharmaceutical R&D efficiency.** *Nat Rev Drug Discov* 2012, **11**:191-200.
- 533 3. Yusuf S, Wittes J: **Interpreting geographic variations in results of**
534 **randomized, controlled trials.** *N Engl J Med* 2016, **375**:2263-2271.
- 535 4. Rasmussen S, Allentoft ME, Nielsen K, Orlando L, Sikora M, Sjogren KG,
536 Pedersen AG, Schubert M, Van Dam A, Kapel CM, et al: **Early divergent**
537 **strains of *Yersinia pestis* in Eurasia 5,000 years ago.** *Cell* 2015, **163**:571-
538 582.
- 539 5. D'Costa VM, King CE, Kalan L, Morar M, Sung WWL, Schwarz C, Froese D,
540 Zazula G, Calmels F, Debruyne R, et al: **Antibiotic resistance is ancient.**
541 *Nature* 2011, **477**:457.
- 542 6. Morozova I, Flegontov P, Mikheyev AS, Bruskin S, Asgharian H, Ponomarenko
543 P, Klyuchnikov V, ArunKumar G, Prokhortchouk E, Gankin Y, et al: **Toward**
544 **high-resolution population genomics using archaeological samples.**
545 *DNA Res* 2016, **23**:295-310.
- 546 7. Marciniak S, Perry GH: **Harnessing ancient genomes to study the history**
547 **of human adaptation.** *Nat Rev Genet* 2017, **advance online publication.**
- 548 8. Elhaik E, Greenspan E, Staats S, Krahn T, Tyler-Smith C, Xue Y, Tofanelli S,
549 Francalacci P, Cucca F, Pagani L, et al: **The GenoChip: a new tool for genetic**
550 **anthropology.** *Genome Biol Evol* 2013, **5**:1021-1031.
- 551 9. Elhaik E, Tatarinova T, Chebotarev D, Piras IS, Maria Calò C, De Montis A,
552 Atzori M, Marini M, Tofanelli S, Francalacci P, et al: **Geographic population**
553 **structure analysis of worldwide human populations infers their**
554 **biogeographical origins.** *Nat Commun* 2014, **5**:1-12.
- 555 10. Phillips C, Parson W, Lundsberg B, Santos C, Freire-Aradas A, Torres M,
556 Eduardoff M, Borsting C, Johansen P, Fondevila M, et al: **Building a forensic**
557 **ancestry panel from the ground up: The EUROFORGEN Global AIM-SNP**
558 **set.** *Forensic Sci Int Genet* 2014, **11**:13-25.
- 559 11. Mao X, Bigham AW, Mei R, Gutierrez G, Weiss KM, Brutsaert TD, Leon-Velarde
560 F, Moore LG, Vargas E, McKeigue PM, et al: **A genomewide admixture**
561 **mapping panel for Hispanic/Latino populations.** *Am J Hum Genet* 2007,
562 **80**:1171-1178.

- 563 12. Elhaik E, Yusuf L, Anderson AIJ, Pirooznia M, Arnellos D, Vilshansky G, Ercal
564 G, Lu Y, Webster T, Baird ML, Esposito U: **The Diversity of REcent and**
565 **Ancient huMan (DREAM): a new microarray for genetic anthropology**
566 **and genealogy, forensics, and personalized medicine.** *Genome Biol Evol*
567 2017, **9**:3225-3237.
- 568 13. Qin H, Zhu X: **Power comparison of admixture mapping and direct**
569 **association analysis in genome-wide association studies.** *Genet Epidemiol*
570 2012, **36**:235-243.
- 571 14. Seldin MF, Price AL: **Application of ancestry informative markers to**
572 **association studies in European Americans.** *PLoS Genet* 2008, **4**:e5.
- 573 15. Halder I, Shriver M, Thomas M, Fernandez JR, Frudakis T: **A panel of**
574 **ancestry informative markers for estimating individual**
575 **biogeographical ancestry and admixture from four continents: utility**
576 **and applications.** *Hum Mutat* 2008, **29**:648-658.
- 577 16. Huckins LM, Boraska V, Franklin CS, Floyd JAB, Southam L, Gcan, Wtccc,
578 Sullivan PF, Bulik CM, Collier DA, et al: **Using ancestry-informative**
579 **markers to identify fine structure across 15 populations of European**
580 **origin.** *Eur J Hum Genet* 2014, **22**:1190.
- 581 17. Kidd KK, Speed WC, Pakstis AJ, Furtado MR, Fang R, Madbouly A, Maiers M,
582 Middha M, Friedlaender FR, Kidd JR: **Progress toward an efficient panel of**
583 **SNPs for ancestry inference.** *Forensic Sci Int Genet* 2014, **10**:23-32.
- 584 18. Xu S, Huang W, Qian J, Jin L: **Analysis of genomic admixture in Uyghur and**
585 **its implication in mapping strategy.** *Am J Hum Genet* 2008, **82**:883-894.
- 586 19. Rosenberg NA, Li LM, Ward R, Pritchard JK: **Informativeness of genetic**
587 **markers for inference of ancestry.** *Am J Hum Genet* 2003, **73**:1402-1422.
- 588 20. Kidd JR, Friedlaender FR, Speed WC, Pakstis AJ, De La Vega FM, Kidd KK:
589 **Analyses of a set of 128 ancestry informative single-nucleotide**
590 **polymorphisms in a global set of 119 population samples.** *Investig Genet*
591 2011, **2**:1.
- 592 21. Wright S: *Evolution and the genetics of populations. A treatise in three*
593 *volumes.* Chicago, IL: University of Chicago Press; 1968.
- 594 22. Shriner D: **Overview of admixture mapping.** *Curr Protoc Hum Genet* 2013,
595 **Chapter 1**:Unit 1 23.
- 596 23. Tian C, Hinds DA, Shigeta R, Adler SG, Lee A, Pahl MV, Silva G, Belmont JW,
597 Hanson RL, Knowler WC, et al: **A genomewide single-nucleotide-**
598 **polymorphism panel for Mexican American admixture mapping.** *Am J*
599 *Hum Genet* 2007, **80**:1014-1023.
- 600 24. Paschou P, Lewis J, Javed A, Drineas P: **Ancestry informative markers for**
601 **fine-scale individual assignment to worldwide populations.** *J Med Genet*
602 2010, **47**:835-847.
- 603 25. Ding L, Wiener H, Abebe T, Altaye M, Go RC, Kericsmar C, Grabowski G, Martin
604 LJ, Hershey GK, Chakorborty R, Baye TM: **Comparison of measures of**
605 **marker informativeness for ancestry and admixture mapping.** *BMC*
606 *Genomics* 2011, **12**:622.
- 607 26. Elhaik E: **Empirical distributions of F_{ST} from large-scale Human**
608 **polymorphism data.** *PLoS One* 2012, **7**:e49837.

- 609 27. Elhaik E, Ryan D: **Pair Matcher (PaM): fast model-based optimisation of**
610 **treatment/case-control matches using demographic and genetic data.**
611 *bioRxiv* 2017.
- 612 28. Cassidy LM, Martiniano R, Murphy EM, Teasdale MD, Mallory J, Hartwell B,
613 Bradley DG: **Neolithic and Bronze Age migration to Ireland and**
614 **establishment of the insular Atlantic genome.** *Proc Natl Acad Sci U S A*
615 2016, **113**:368-373.
- 616 29. Li H, Durbin R: **Fast and accurate short read alignment with Burrows-**
617 **Wheeler transform.** *Bioinformatics* 2009, **25**:1754-1760.
- 618 30. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis
619 G, Durbin R: **The Sequence Alignment/Map format and SAMtools.**
620 *Bioinformatics* 2009, **25**:2078-2079.
- 621 31. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A,
622 Garimella K, Altshuler D, Gabriel S, Daly M, DePristo MA: **The genome**
623 **analysis toolkit: a MapReduce framework for analyzing next-generation**
624 **DNA sequencing data.** *Genome Res* 2010, **20**:1297-1303.
- 625 32. DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, Philippakis
626 AA, del Angel G, Rivas MA, Hanna M, et al: **A framework for variation**
627 **discovery and genotyping using next-generation DNA sequencing data.**
628 *Nat Genet* 2011, **43**:491-498.
- 629 33. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker
630 RE, Lunter G, Marth GT, Sherry ST, et al: **The variant call format and**
631 **VCftools.** *Bioinformatics* 2011, **27**:2156-2158.
- 632 34. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, Maller J,
633 Sklar P, de Bakker PI, Daly MJ, Sham PC: **PLINK: a tool set for whole-**
634 **genome association and population-based linkage analyses.** *Am J Hum*
635 *Genet* 2007, **81**:559-575.
- 636 35. Alexander DH, Novembre J, Lange K: **Fast model-based estimation of**
637 **ancestry in unrelated individuals.** *Genome Res* 2009, **19**:1655-1664.
- 638 36. Durbin RM, Abecasis GR, Altshuler DL, Auton A, Brooks LD, Gibbs RA, Hurles
639 ME, McVean GA: **A map of human genome variation from population-**
640 **scale sequencing.** *Nature* 2010, **467**:1061-1073.
- 641 37. Paschou P, Ziv E, Burchard EG, Choudhry S, Rodriguez-Cintron W, Mahoney
642 MW, Drineas P: **PCA-correlated SNPs for structure identification in**
643 **worldwide human populations.** *PLoS Genet* 2007, **3**:1672-1686.
- 644 38. Marcus JH, Novembre J: **Visualizing the geography of genetic variants.**
645 *Bioinformatics* 2017, **33**:594-595.

646

647

648 **Figure Legends**

649

650 [Figure 1](#)

651 **Geographic distribution of the highly differentiated rs7896530 in ancient (A) and**

652 **modern-day (B) populations.** The geographic distributions of the T (black) and G

653 (yellow) alleles in ancient and modern-day populations were obtained from our dataset

654 ([Table S1](#)) and the Geography of Genetic Variants Browser [38], respectively.

655

656

657 Figure 2

658 **A scheme to identify and evaluate aAIMs.**

659

660

661 **Figure 3**

662 **Geographical locations of the ancient genomes.** The geographical coordinates of the
663 ancient genomes. The shapes plotted in the map designate the country of origin of the
664 genomes and their colors designate the era. The total number of ancient genomes from a
665 specific era is shown.

666

667

668 **Figure 4**

669 **Scatter plot of all ancient populations along the first two principal components.**

670 Symbols corresponding to individuals and their color and shape correspond to the

671 location map and the era table, respectively.

672

673

674 **Figure 5**

675 **Ancient population structure inferred by ADMIXTURE analysis.** Each individual is

676 represented by a vertical (100%) stacked column of genetic components proportions

677 shown in color for $K=10$.

678

679

680 **Figure 6**

681 **Minor allele frequency distributions for aAIMs identified with various methods.**

682 MAF frequencies were calculated for ancient (A) and modern-day (B) populations. To

683 avoid confusion, the distributions represent the frequency of the minor allele in each

684 datasets, which was the same one in 91.5% of the genotypes.

685

686

687 **Figure 7**

688 **Violin plots comparing the Euclidean distances between the admixture proportions**

689 **of the ancient genomes obtained from the CSS and those obtained from the aAIM**

690 **sets.**

691

692

693 Figure 8

694 **Genome wide distribution of SNPs in the CSS (dots) and PD (red bars) datasets.**

695

696

697 **Supporting Information Legends**

698 Elhaik et al 2018 – Supp – Figures S1-S7 and Tables S1-S4

699 Additional file 1.zip – Genotype data of the aDNA samples

700 Additional file 2.zip – aAIMs candidates used in all analyses

701 Additional file 3.zip – Genotype file of the putative ancient ancestral populations

702 **Tables**

703

704 **Table 1**

705 **Accuracy in classifying individuals to populations using the aAIM candidates. Mean**

706 and standard deviation for each SNP set are provided in the last row.

707

Population	<i>n</i>	CSS	PD	F_{ST}	Infocalc	Admixture ₁	Admixture ₂	Rand _{10k}	Rand _{15k}
Britain Iron Saxon	10	10 (100)	4 (40)	0 (0)	0 (0)	0 (0)	0 (0)	1 (10)	3 (30)
Caucasus Chalcolithic Bronze	22	21 (95)	8 (36)	0 (0)	12 (55)	6 (27)	4 (18)	13 (59)	9 (41)
Caucasus Mesolithic Neolithic	9	6 (67)	7 (78)	0 (0)	6 (67)	1 (11)	7 (78)	4 (44)	4 (44)
Central EU Early Neolithic	26	17 (65)	14 (54)	4 (15)	18 (69)	4 (15)	5 (19)	14 (54)	18 (69)
Central EU Late Neolithic Bronze	57	16 (28)	17 (30)	19 (33)	19 (33)	13 (23)	21 (37)	25 (44)	21 (37)
Central EU Mid Neolithic Chalcolithic	6	2 (33)	3 (50)	0 (0)	3 (50)	3 (50)	3 (50)	2 (33)	2 (33)
Central Northern EU Late Neolithic Bronze	20	18 (90)	9 (45)	0 (0)	6 (30)	0 (0)	5 (25)	4 (20)	6 (30)
Central Western EU Mesolithic	3	3 (100)	2 (67)	0 (0)	3 (100)	0 (0)	0 (0)	1 (33)	3 (100)
Italy Mid Neolithic Chalcolithic	4	4 (100)	3 (75)	0 (0)	1 (25)	1 (25)	0 (0)	1 (25)	1 (25)
Jordan Bronze	3	3 (100)	2 (67)	0 (0)	0 (0)	2 (67)	3 (100)	1 (33)	2 (67)
Levant Epipaleolithic Neolithic	19	7 (37)	6 (32)	0 (0)	9 (47)	8 (42)	7 (37)	4 (21)	7 (37)
Russia Chalcolithic	3	2 (67)	3 (100)	0 (0)	1 (33)	0 (0)	2 (67)	1 (33)	1 (33)
Russia Early Mid Bronze	19	19 (100)	15 (79)	0 (0)	10 (53)	0 (0)	18 (95)	10 (53)	14 (74)
Russia Late Chalcolithic	9	6 (67)	6 (67)	0 (0)	5 (56)	0 (0)	1 (11)	3 (33)	3 (33)
Russia Mesolithic	3	2 (67)	2 (67)	0 (0)	2 (67)	0 (0)	1 (33)	2 (67)	2 (67)
Russia Mid Late Bronze	22	15 (68)	16 (73)	0 (0)	7 (32)	0 (0)	0 (0)	4 (18)	6 (27)
Spain Early Neolithic	6	4 (67)	5 (83)	0 (0)	6 (100)	4 (67)	4 (67)	4 (67)	5 (83)
Spain Mid Neolithic Chalcolithic	18	7 (39)	6 (33)	0 (0)	7 (39)	5 (28)	3 (17)	5 (28)	5 (28)
Sweden Mesolithic	8	8 (100)	8 (100)	0 (0)	7 (88)	4 (50)	1 (13)	6 (75)	7 (88)
Sweden Mid Neolithic	4	4 (100)	1 (25)	1 (25)	2 (50)	1 (25)	0 (0)	4 (100)	2 (50)
Turkey Neolithic	24	23 (96)	18 (75)	0 (0)	12 (50)	3 (13)	6 (25)	8 (33)	11 (46)
		76±25	61±23	3±9	50±27	21±23	33±32	42±22	50±23

708

709

710 **Table 2**

711 **Accuracy of inferring hybrid individuals using the PD aAIMs.** The parental
 712 populations and the number of hybrids generated from them are shown. Each hybrid was
 713 represented by three datasets: CSS, PD aAIMs, and a random SNP set. The average
 714 genetic distances (d) between the admixture components of these datasets per population
 715 are shown.

716

Parental population A	Parental population B	# Hybrids	$\bar{d}(\text{CSS, PD})$	$\bar{d}(\text{CSS, random set})$	$\bar{d}(\text{PD, random set})$
Britain Iron Saxon	Britain Iron Saxon	6	0.026	0.212	0.208
Britain Iron Saxon	Russia Late Chalcolithic	9	0.009	0.610	0.601
Britain Iron Saxon	Sweden Mesolithic	9	0.051	0.344	0.337
Britain Iron Saxon	Turkey Neolithic	9	0.003	0.428	0.431
Britain Iron Saxon	Spain Early Neolithic	9	0.108	0.221	0.241
Russia Late Chalcolithic	Russia Late Chalcolithic	6	0.009	0.443	0.448
Russia Late Chalcolithic	Sweden Mesolithic	9	0.062	0.578	0.561
Russia Late Chalcolithic	Turkey Neolithic	9	0.063	0.661	0.633
Russia Late Chalcolithic	Spain Early Neolithic	9	0.101	0.520	0.491
Sweden Mesolithic	Sweden Mesolithic	6	0.000	0.384	0.384
Sweden Mesolithic	Turkey Neolithic	9	0.055	0.567	0.522
Spain Early Neolithic	Sweden Mesolithic	9	0.108	0.402	0.377
Turkey Neolithic	Turkey Neolithic	6	0.001	0.627	0.626
Spain Early Neolithic	Turkey Neolithic	9	0.092	0.483	0.493
Spain Early Neolithic	Spain Early Neolithic	6	0.041	0.197	0.172

717

718

719

720

721

722

723

724 **Competing interests**

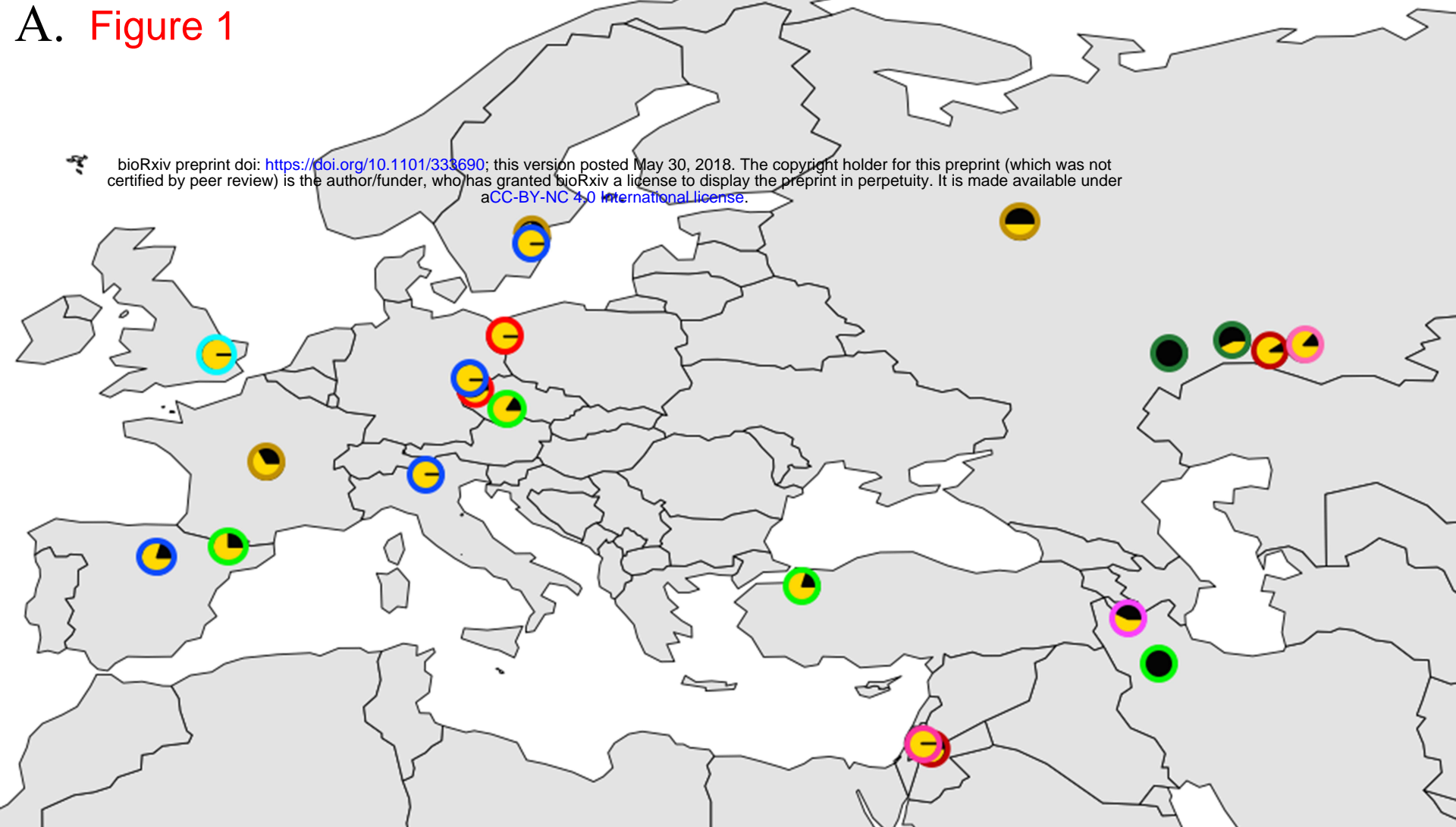
725 EE is a consultant to DNA Diagnostic Centre. The funders had no role in study design,

726 data collection and analysis, decision to publish, or preparation of the manuscript.

727

A. Figure 1

bioRxiv preprint doi: <https://doi.org/10.1101/333690>; this version posted May 30, 2018. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC 4.0 International license.



Mesolithic

Late-Neolithic Bronze

Early-Mid Bronze

Early-Neolithic/Neolithic

Chalcolithic/Late Chalcolithic

Mid-Late Bronze

Mid-Neolithic/Mid-Neolithic Chalcolithic

Early Bronze/Chalcolithic Bronze

Epipaleolithic

Anglo-Saxon

B.

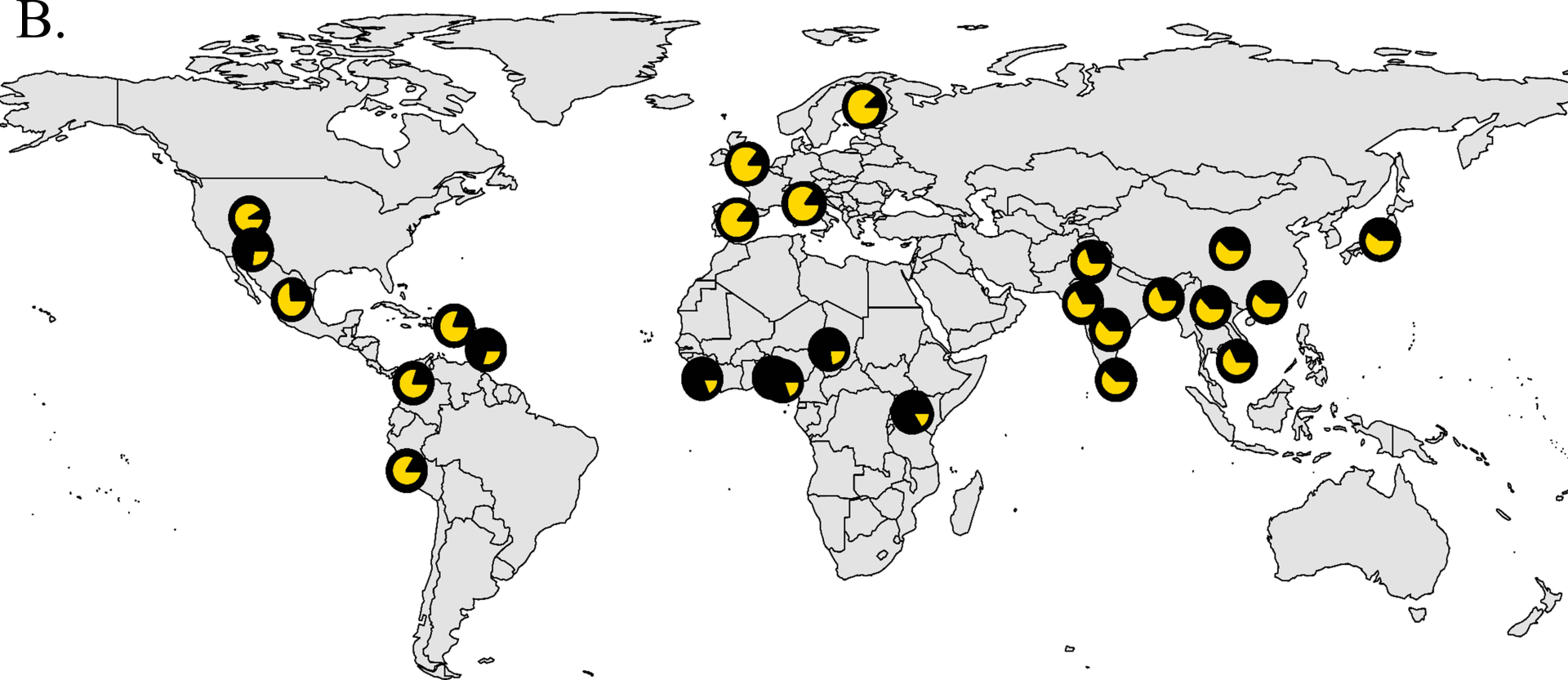
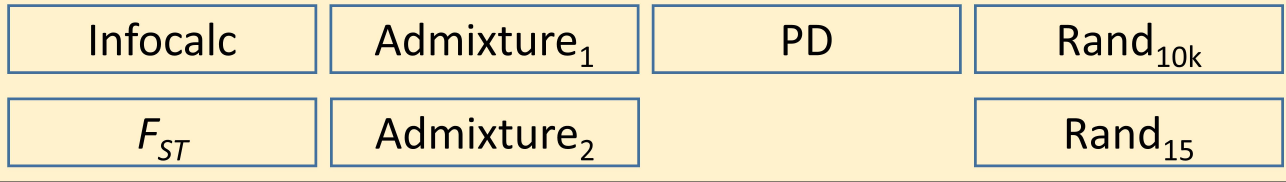


Figure 2

Ancient DNA dataset

Apply aAIMs methods to the aDNA dataset

Initial aAIM candidate sets



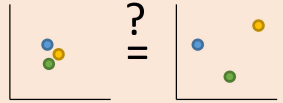
Filter samples/SNPs for missingness

Final aAIM candidate sets

Calculate overlap between the aAIM sets
Calculate aAIMs overlap
Calculate aAIMs MAF distributions
Functional annotation

Calculate PCA for each aAIM set

1) Qualitative assessment of PCA plots



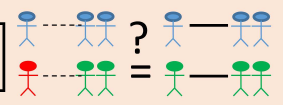
Calculate admixture for each aAIM set

2) Comparison of admixture plots



Classify individuals to populations

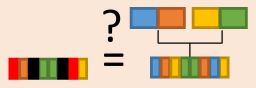
3) Comparison of individual classification



Further analyses with the best aAIMs

Hybrids dataset

4) Test hybrid inference



Calculate genome wide distribution of aAIMs

Ancient genomic data

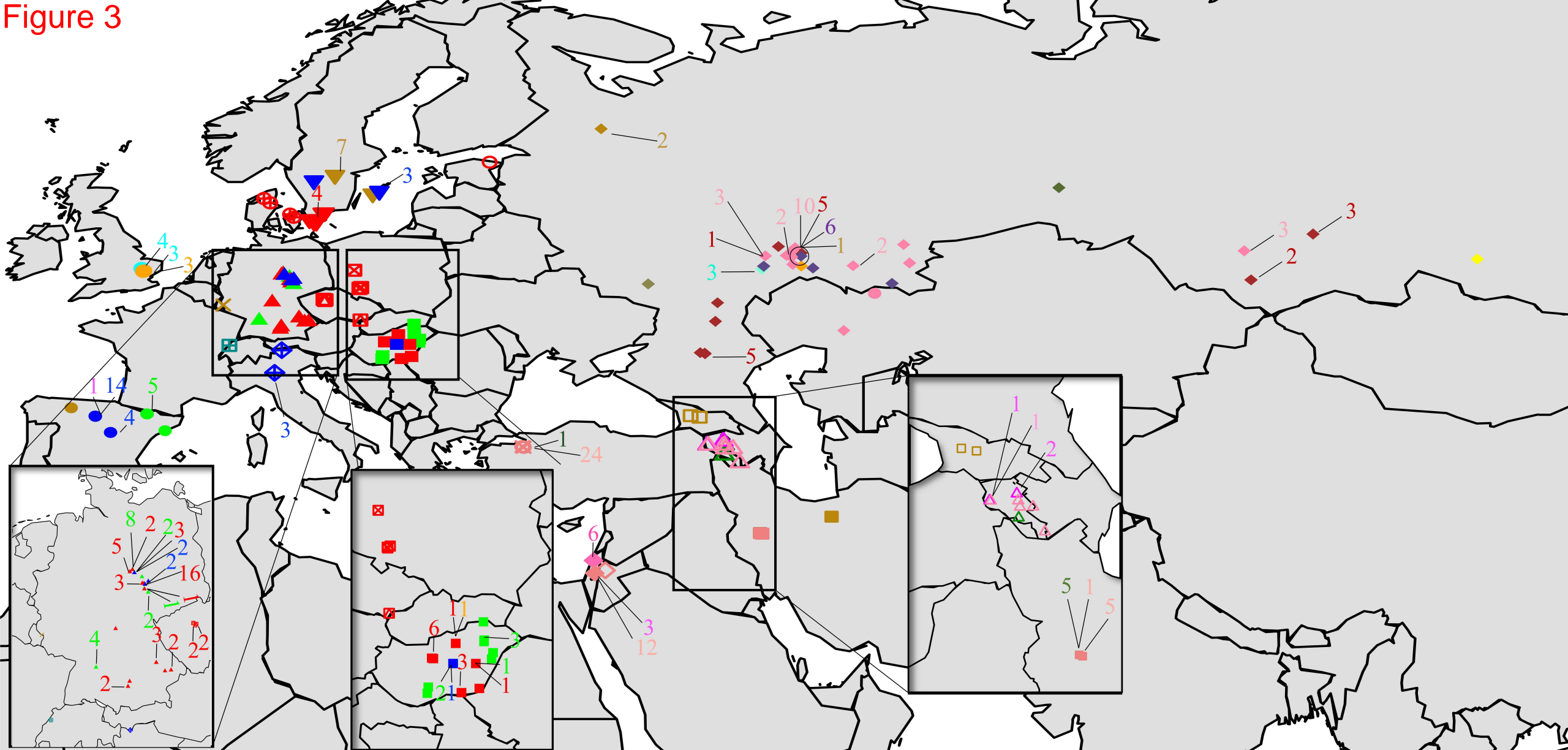
Identifying aAIM candidates

Criteria to evaluate the candidate aAIMs:

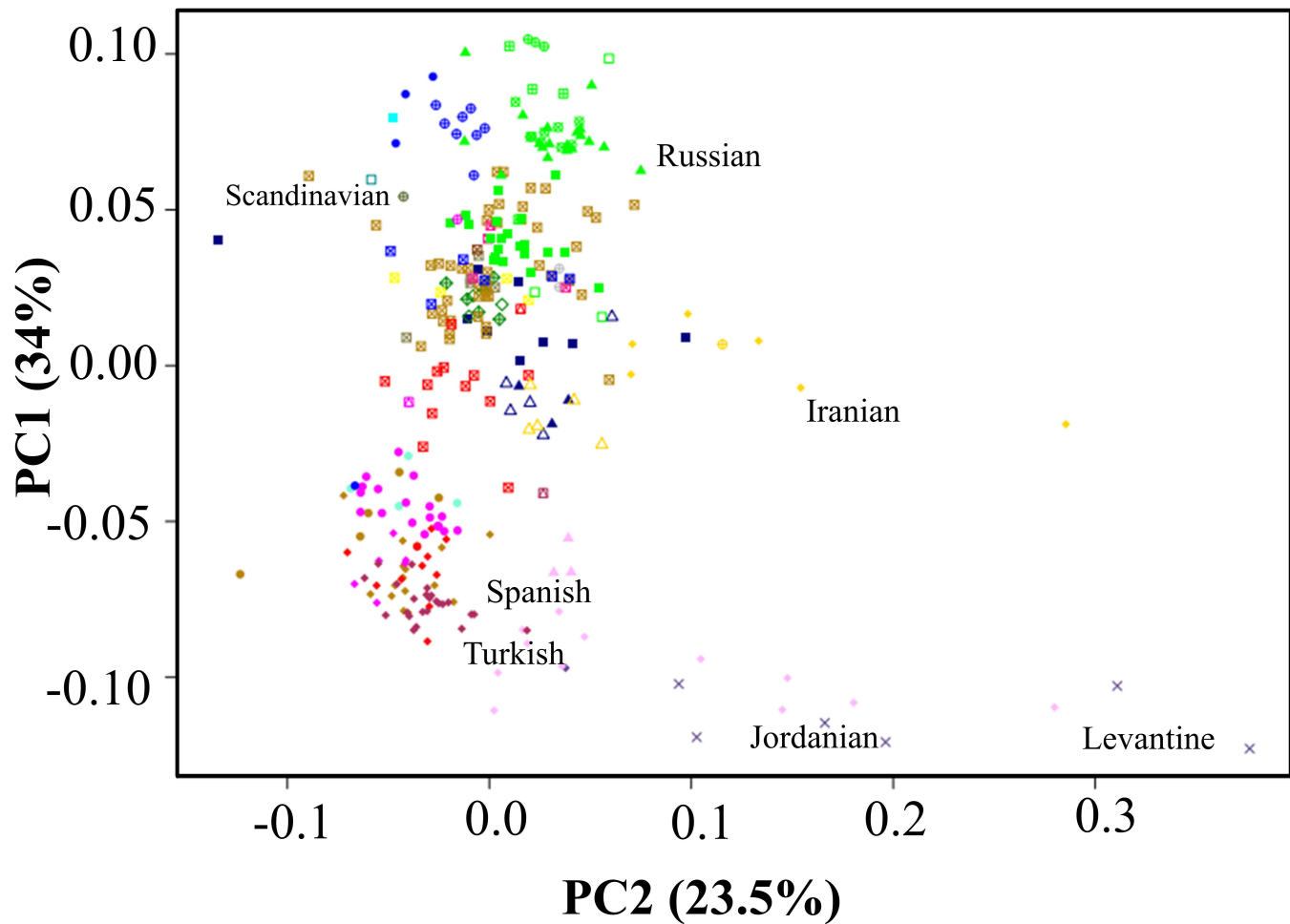
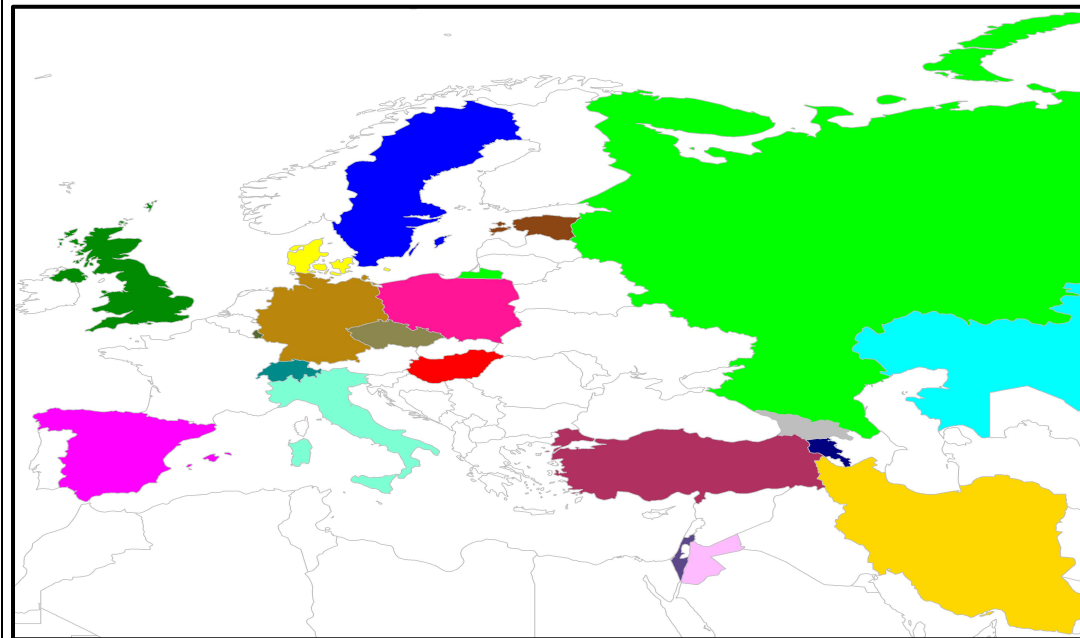
Can aAIMs capture the population structure?

Can aAIMs be used for admixture mapping?

Figure 3



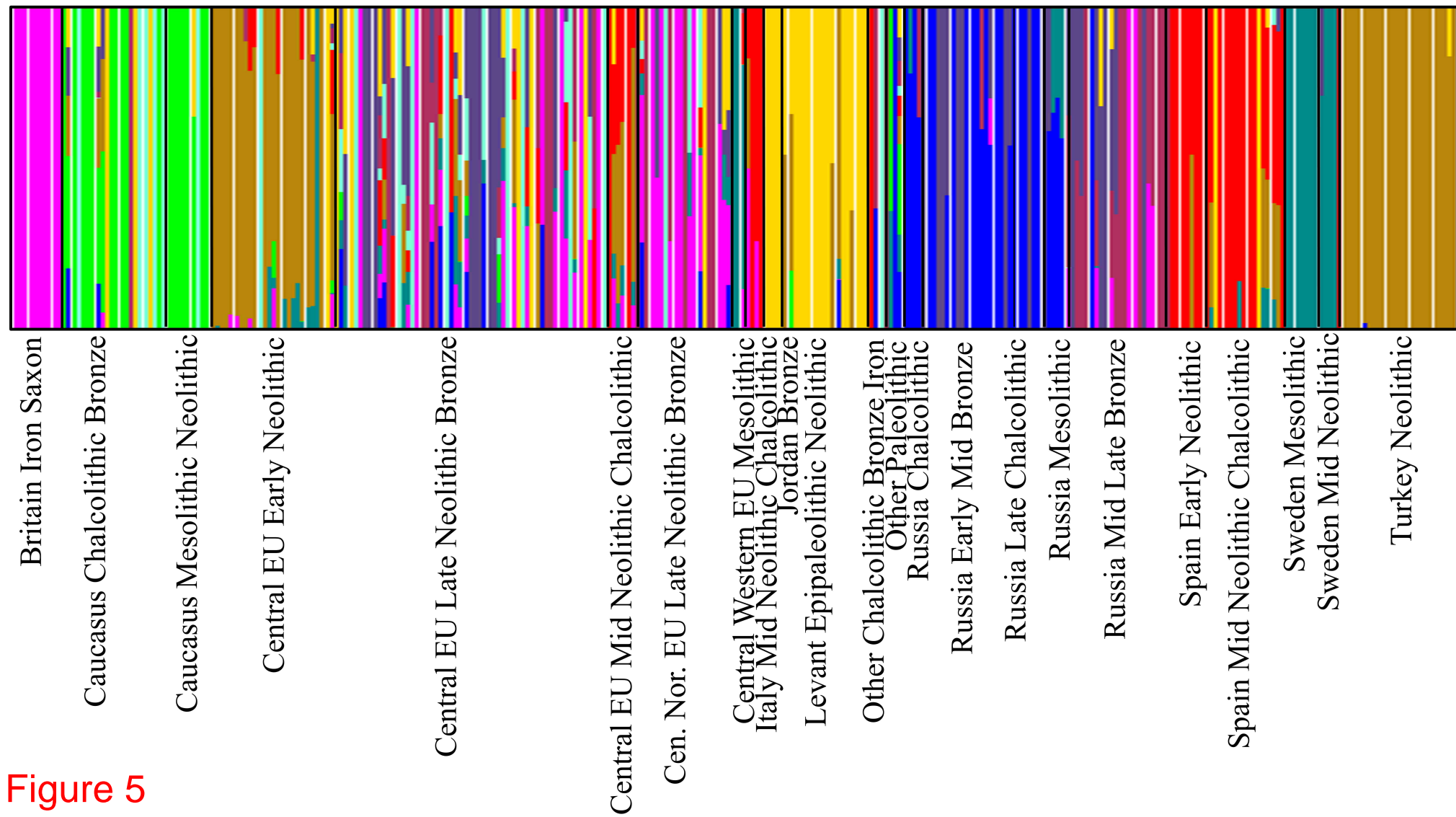
- Upper Paleolithic (1)
- Mesolithic (16)
- Early-Neolithic (32)
- Mid-Neolithic (32)
- Late-Neolithic Bronze (77)
- Neolithic (43)
- Late Copper (9)
- Chalcolithic (11)
- Early Bronze (7)
- Early-Mid Bronze (19)
- Mid-Late Bronze (31)
- Iron/Late Iron (5)
- Russia Eneolithic (3)
- Russia Kostenki (1)
- Russia MA1 (1)
- Russia Ust Ishim (1)
- Israel Epipaleolithic (6)
- Britain Anglo-Saxon (7)



Upper Paleolithic (30,000 - 12,000 BP)	Epipaleolithic (12,000 BP - 8,300 BC)	Mesolithic (10,000 - 5,000 BC)
□	×	⊕
Early Neolithic (8,300 - 5,500 BC)	Neolithic (8,300 - 5,500 BC)	Mid Neolithic Chalcolithic (5,500 - 4,500 BC)
◆	◆	●
Chalcolithic (4,500 - 4,000 BC)	Late Copper (4,000 - 3,300 BC)	Late Neolithic Bronze (3,300 - 3,000 BC)
△	⊗	⊗
Early Bronze (3,000 - 2,700 BC)	Early Mid Bronze (2,000 - 1,750 BC)	Mid Late Bronze (1,400 - 1,200 BC)
▲	▲	■
Iron (1,200 - 586 BC)	Anglo-Saxons (450 - 1,066 AD)	Other-Chalc. Bronze & Iron
◇	⊕	⊗

Figure 4

Figure 5



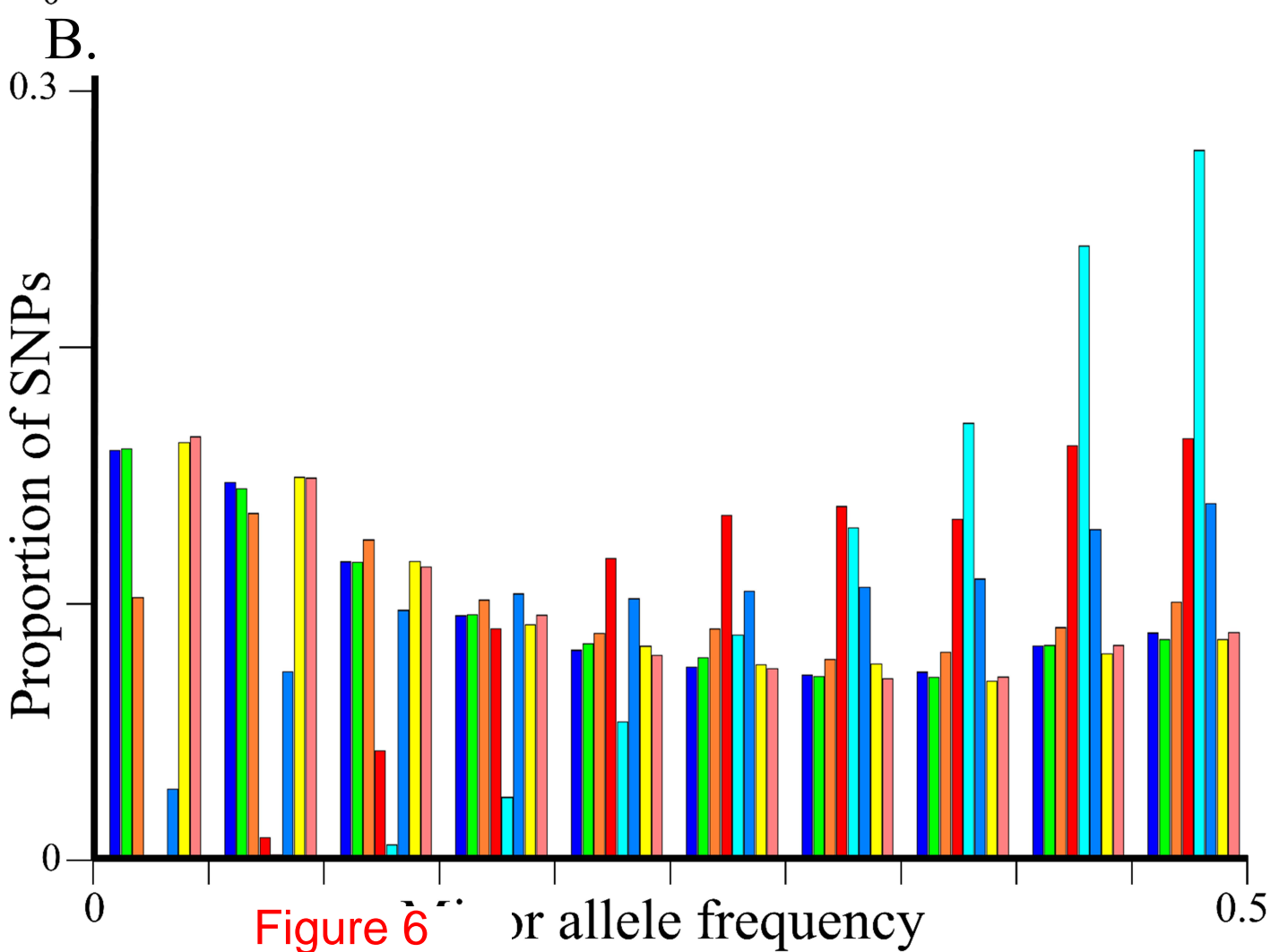
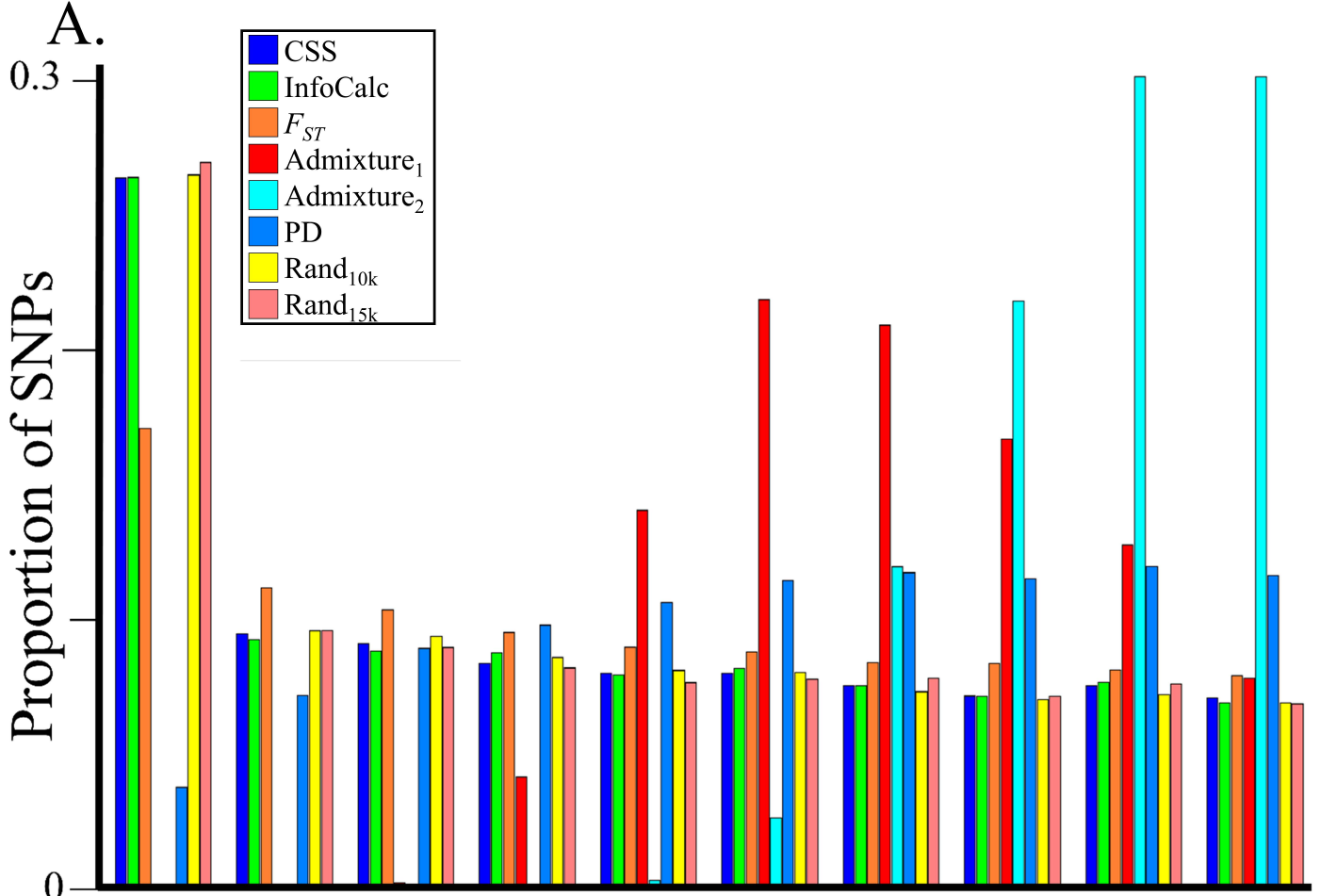


Figure 6

Figure 7

Total Euclidean distance

