

1 Population Genomics of Stone Age
2 Eurasia:
3 **Supplementary Information (Part 1)**

5	1) Data Generation and Authentication	7
6	Sampling, lab work and sequencing	7
7	Basic bioinformatics	8
8	DNA authentication	9
9	DNA contamination - results and implications	10
10	References	11
11	2) Imputation of ancient DNA	14
12	Introduction	14
13	Methods	14
14	Results	15
15	Effects of Low coverage	17
16	References	29
17	3) Demographic inference	31
18	3a) Phylogenetic analysis of mtDNA sequences	31
19	Methods	31
20	Results	31
21	Conclusion	32
22	References	35
23	3b) Y chromosome / Sex determination	35
24	Methods	35
25	Results	36
26	Sex determination	36
27	Phylogenetic placement	37
28	Sub-haplogroup analyses of newly reported samples	38
29	References	45
30	3c) Relatedness	45
31	Methods	46
32	Results	46
33	References	49
34	3d) Pop structure general, PCA/Admixture (Martin)	49
35	Methods	49
36	Results	52
37	Comparison of imputed and pseudo-haploid genotypes in PCA space	52
38	PCA position of samples flagged as contaminated	56
39	Genetic ancestry of newly reported samples	57

40	References	64
41	3e) Inferring the spatiotemporal spread of population movements in the past 13 millennia	
42		65
43	Introduction	65
44	Methods	65
45	Results	65
46	Figures	67
47	References	71
48	3f) HBD/ IBD sharing/ROH/clustering	73
49	Methods	73
50	Results	74
51	IBD-based hierarchical graph clustering	74
52	Runs of homozygosity and IBD sharing within clusters	104
53	Mixture models	107
54	References	113
55	3g) Selecting non-British individuals from the UK Biobank	114
56	Introduction	114
57	Methods	114
58	Results	115
59	Discussion	115
60	References	121
61	3h) Painting the UK BioBank	121
62	Introduction	121
63	Methods	121
64	Painting pipeline introduction	121
65	Reference/donor panel formation	122
66	Target/recipient panel formation	122
67	SNP selection and merging of the panels	123
68	Painting process	123
69	Painting at biobank scale	124
70	Results	125
71	Ancestry-PCs relationship	125
72	Ancestry-geographic variation	125
73	Discussion	127
74	Figures	128
75	References	145
76	3i) Building a population history model of Europeans and using it to assign local ancestry	
77	to all haplotypes in modern and ancient European samples	147
78	Introduction	147
79	Method and Results	147
80	Model of population structure	147
81	Local Ancestry Using Tree Sequences:	148
82	Figures:	151
83	References	153

84	4) Natural selection and trait evolution	155
85	4a) Estimating allele frequency trajectories of trait-associated variants	155
86	Introduction	155
87	Methods	156
88	SNP Ascertainment	156
89	GWAS SNPs	156
90	Control SNPs	156
91	Simulated Neutral SNPs	156
92	1000G ARG	157
93	Data pre-processing	157
94	ARG Inference	157
95	Modifications to CLUES	157
96	ARG sampling using Relate	157
97	Ancient DNA samples	158
98	Selection Analysis	158
99	CLUES with Modern 1000G data	158
100	CLUES with aDNA Time Series	158
101	CLUES with aDNA Ancestral Paintings	159
102	Reference and mapping bias filters	159
103	Genome-wide selection	160
104	Results	162
105	Selection in 1000G EUR	162
106	Selection in aDNA Time Series	163
107	Peak 1: LINC02797, AL583808.1	165
108	Peak 2: MCM6	166
109	Peak 3: PTH1R - AC109583.3	167
110	Peak 4: KLF3 - RNA5SP158	168
111	Peak 5: P4HA2, SLC22A4	169
112	Peak 6: FLT4	170
113	Peak 7: HLA-A - HLA-W	171
114	Peak 8: FADS2	172
115	Peak 9: HECTD4	173
116	Peak 10: AC092143.1, MC1R	174
117	Peak 11: MAPT	175
118	Selection in simulations with Ancestral Paintings	176
119	Selection in aDNA with Ancestral Paintings	178
120	Peak 1: AL591122.1	180
121	Peak 2: NEGR1	181
122	Peak 3: MCM6	182
123	Peak 4: CACNA2D2	183
124	Peak 5: KLF3-AS1	184
125	Peak 6: BANK1	185
126	Peak 7: SLC45A2	186
127	Peak 8: SLC22A4	187

128	Peak 9: GRK6	188
129	Peak 10: HLA	189
130	Peak 11: MSRA	190
131	Peak 12: ABO	191
132	Peak 13: MYBPC3	192
133	Peak 14: FADS2	193
134	Peak 15: HECTD4	194
135	Peak 16: RFLNA	195
136	Peak 17: MARK3	196
137	Peak 18: SEMA6D	197
138	Peak 19: CSK	198
139	Peak 20: DPEP1	199
140	Peak 21: MAPT	200
141	Discussion	201
142	Pan-ancestry selection	201
143	Ancestry stratified selection trajectories	203
144	References	204
145	4b) Detangling Direct and Indirect impacts of sample age from the Mesolithic-Neolithic	
146	data on genotype imputation	221
147	References	225
148	4c) Over-dispersion in polygenic scores across ancient populations	225
149	Introduction	226
150	Methods	226
151	Results	228
152	Polygenic scores across ancient populations	228
153	Individual scores across time and space	228
154	Differentiation among ancient populations and GBR	228
155	Discussion	229
156	Figures	232
157	References	242
158	4d) Identifying candidates for positive selection using patterns of ancient population	
159	differentiation	244
160	Introduction	244
161	Methods	244
162	Results	245
163	Eurasian scan	245
164	West Eurasian scan	246
165	Neolithic vs. hunter-gatherer scan	248
166	Figures	250
167	References	310
168	4e) Correlation between components of variation in population structure and components	
169	of variation in SNP-trait association	320
170	Introduction	321
171	Methods	321
172	Results	322

173	Discussion	323
174	Figures	324
175	References	336
176	4f) Polygenic prediction for height, eye colour and hair colour in ancient Danish samples	
177		336
178	References	342
179	4g) Calling chr17q21.31 KANSL1 Duplications in Ancient Genomes	343
180	4h) Calculating ancestral contributions to modern complex phenotypes	345
181	Introduction	345
182	Methods	346
183	Results	347
184	Discussion	349
185	Figures/tables	350
186	References	352
187	4i) Pathogenic structural variants in ancient vs. modern-day humans	355
188	References	378
189		
190		
191		
192		
193		

194 Key to Supplementary Data Tables

No.	Title
Supplementary Table I.	Basic overview of samples and genetic data
Supplementary Table II.	Samples, dates and isotopic data
Supplementary Table III.	Reservoir correction calculations for radiocarbon dates
Supplementary Table IV.	Full sample metadata (dates and isotopic data non-cleaned and non-combined)
Supplementary Table V.	DNA contamination estimation
Supplementary Table VI.	Relatedness Estimates
Supplementary Table VII.	Ancient Genomes dataset (incl. all imputed)
Supplementary Table VIII.	IBD mixture model sets
Supplementary Table IX.	IBD Ancestry proportions for set "deep"
Supplementary Table X.	IBD Ancestry proportions for set "postNeol"
Supplementary Table XI.	IBD Ancestry proportions for set "postBA"
Supplementary Table XII.	IBD Ancestry proportions for set "hgEur"
Supplementary Table XIII.	IBD Ancestry proportions for set "fEur"
Supplementary Table XIV.	IBD Ancestry proportions for set "postNeolScand"
Supplementary Table XV.	CLUES modern
Supplementary Table XVI.	CLUES ancient

195

196

197

1) Data Generation and Authentication

198

199

200 Morten E. Allentoft^{1,2}, Simon Rasmussen³, Gabriel Renaud⁴, Abigail D. Ramsøe², Thorfinn

201

Korneliusson², Martin Sikora²

202

203 ¹Trace and Environmental DNA (TrEnD) Lab, Curtin University, Perth, Australia

204 ²Lundbeck Foundation GeoGenetics Centre, GLOBE Institute, University of Copenhagen,

205

Copenhagen, Denmark

206

³Novo Nordisk Foundation Center for Protein Research, University of Copenhagen,

207

Copenhagen, Denmark

208

⁴Department of Health Technology, Section of Bioinformatics, Technical University of

209

Denmark, Kongens Lyngby, Denmark

210 Sampling, lab work and sequencing

211 The lab work component of this project followed the same procedures outlined in Allentoft et

212 al. ¹ and Damgaard et al. ², also including sampling of the petrous part of the temporal bone

213 following the discovery of exceptional DNA preservation in these bones ^{3,4}. While new

214 ancient DNA (aDNA) extraction and library methods are continually optimised and presented

215 in the literature, we prioritised method consistency throughout the project period to avoid the

216 risk of introducing batch effects in the data.

217

218 A total of 962 Stone Age and early Bronze Age human skeletons from across Eurasia were

219 sampled for this project. An initial molecular 'screening' to assess the endogenous DNA

220 content (proportion of DNA sequences identified as human) was performed by shallow

221 shotgun sequencing resulting in 317 samples (**Supplementary Table I**) being selected for

222 deeper sequencing. We applied a threshold of <1% endogenous DNA for rejecting samples

223 in the project, except for a few Danish skeletons that were prioritised despite displaying even

224 lower contents. Of the 317 samples, 211 were teeth, 91 were petrous bones, and 15 were

225 pieces of other types of bones (long bones, ribs, cranial bones).

226

227 Following some AMS-dates which were conducted late in the project, two of the 317 samples

228 (NEO901 and NEO902, **Supplementary Table I**) proved too young to be relevant for this

229 project. These two samples have only been included for imputation purposes together with

230 the data that are released as part of this project, but they have not been included in

231 downstream analyses. Hence the final number of new samples that were sequenced,

232 analysed and discussed as part of this project is 315, and this is the number we refer to

233 throughout the study.

234

235 All the pre-PCR-amplification lab work was conducted in dedicated clean laboratories at the
236 Lundbeck Foundation GeoGenetics Centre (GLOBE Institute, University of Copenhagen),
237 according to strict aDNA guidelines⁵⁻⁷. To reduce the amount of non-target DNA in the
238 extracts, the outermost surfaces of the samples were first removed using a sterile cutting
239 disc. Teeth were processed by separating the crown from the root by a cutting disc and the
240 inner dentine was then removed from the root with a pointy drilling bit. By this procedure
241 each root sample was proportionally enriched for the outer cementum layer which is known
242 for its high endogenous DNA content^{2,8,9}. Petrous bones were sampled by cutting off slices
243 (with a cutting disc) until reaching the dense otic capsule which was used for DNA extraction.
244 The samples were crushed into smaller pieces before the lysis step.

245

246 To further increase the endogenous DNA yield, we performed a brief 'pre-digestion' step
247 prior to the extraction protocol following Damgaard et al.². After this pre-digestion, we added
248 3.5mL of fresh digestion buffer to each sample and incubated them for 24h before the DNA
249 was purified with silica-in-solution similar to Rohland & Hofreiter¹⁰ but using the optimised
250 binding buffer from Allentoft et al.¹. Double-stranded blunt-end libraries were constructed
251 from the extracted DNA using NEBNext DNA Prep Master Mix Set E6070 (New England
252 Biolabs Inc.) with protocol modifications¹¹ and then amplified with indexed Illumina-specific
253 adapters prepared as in¹². The DNA concentration of each amplified library was quantified
254 on an Agilent 2200 TapeStation and sequencing (80bp and 100bp single read) was
255 performed on Illumina HiSeq 2500 and Illumina HiSeq 4000 platforms at the Danish National
256 High-throughput DNA Sequencing Centre.

257 Basic bioinformatics

258 The Illumina data was base-called using Illumina software *CASAVA (v.1.8.2)*¹³ and
259 sequences were de-multiplexed with the requirement of full matching of the six nucleotide
260 index which was used for library preparation. Adapter sequences and leading/trailing
261 stretches of Ns were trimmed from the reads and bases with quality 2 or less were removed
262 using *AdapterRemoval (v.2.1.3)*. Trimmed reads of at least 30bp were mapped using *bwa*
263 *(v.0.7.10)*¹⁴ with the seed disabled to allow for higher sensitivity¹⁵. Reads were mapped to
264 the human reference genome build 37 including mitochondrial DNA (rCRS) and to
265 mitochondrial DNA alone. Mapped reads were filtered for mapping quality 30 and sorted
266 using *Picard (v.1.127)* (<http://picard.sourceforge.net>) and *samtools*¹⁶. Data was merged to
267 library level and duplicates removed using *Picard MarkDuplicates (v.1.127)* and hereafter
268 merged to sample level. Sample level BAMs were re-aligned using *GATK (v.3.3.0)* and

269 hereafter had the md-tag updated and extended BAQs calculated using *samtools calmd*
270 (*v.1.10*)¹⁶. Read depth and coverage were determined using *pysam*
271 (<https://github.com/pysam-developers/pysam>) and *BEDtools* (*v.2.23.0*)¹⁷.

272
273

274 DNA authentication

275 To investigate authenticity of the ancient DNA molecules, post-mortem DNA damage
276 patterns were determined using *mapDamage2.0*¹⁸. Cytosine deamination was recorded for
277 each sample as the fraction of C-to-T transitions at the first 5' position of the DNA reads when
278 compared to the reference genome. For the 317 samples included here, we observed C-to-T
279 deamination fractions ranging from 10.4% to 67.8%, with an average of 38.3% across all 315
280 samples (**Supplementary Table I**). These numbers generally reflect highly damaged
281 molecules as expected for DNA that is thousands of years old.

282

283 Next, we applied three different methods to estimate levels of DNA contamination; two based
284 on mitochondrial genome data and one method investigating X-chromosomal data in males.
285 All contamination estimates are reported in **Supplementary Table V** with summary values
286 provided in **Supplementary Table I**.

287

288 First, estimates of present-day human contamination for the mitochondrial genome were
289 performed using the iterative Bayesian framework implemented in *Schmutzi*¹⁹. Briefly,
290 ancient DNA sequences were realigned to the mitochondrial genome of the revised
291 Cambridge Reference Sequence (rCRS, NCBI Reference Sequence: NC_012920.1) using
292 *BWA*¹⁴ with parameters for increased sensitivity (*-n 0.01 -o 2 -l 16500*). The mapping was
293 performed exclusively to the mitochondrial genome to mitigate the impacts of nuclear
294 NUMTs on the mitochondrial alignments. The resulting BAM file was used as input for
295 *Schmutzi* using five iterations and by subsampling samples with coverage to 500X, should
296 they exceed that¹⁹. The iterative approach was run using a database of Eurasian
297 mitogenomes. The point estimate for the final contamination rate with the maximum *a*
298 *posteriori* probability is reported in the sample summary **Supplementary Table I**, whereas the
299 95% confidence interval (lower and upper bound as well as the point estimate) are reported
300 in **Supplementary Table V**.

301

302 Next, we applied *ContamMix* in order to estimate the fraction of non-endogenous reads in
303 the mitochondrial genome by comparing the reconstructed mtDNA consensus sequence to

304 311 possible contaminate genomes²⁰. For each sample, an in-house perl script was used to
305 construct two different versions of the endogenous mitochondrial genome. The first approach
306 (CONTAMIX_APPROX_1Xdif05) used sites with at least 1x coverage, and at each position a
307 base was only called if it was observed in at least 50% of reads covering the site. The
308 second approach (CONTAMIX_PRECISE_5Xdif07) only considered sites with at least 5x
309 coverage and 70% of reads agreeing. Both approaches used reads with a base quality of
310 ≥ 20 and mapping quality of ≥ 30 .

311

312 Lastly, we applied *ANGSD*²¹ on X-chromosomal data in males. This approach quantifies
313 heterozygosity on the X chromosome. As males only have one copy of the X chromosome,
314 any heterozygosity is expected to arise from either contamination or sequencing error. As
315 heterozygosity due to contamination is expected to be restricted to mainly known diagnostic
316 polymorphic sites, *ANGSD* quantifies the heterozygosity in these sites in. It then compares it
317 to adjacent sites in order to ascertain the level of background sequencing error, and thus
318 estimates the extent of contamination. For each sample, we removed the pseudoautosomal
319 regions on the X chromosome and filtered out reads with a base quality < 20 and mapping
320 quality < 30 .

321

322 DNA contamination - results and implications

323 **Supplementary Table V** lists all the contamination estimation results for the 317 samples
324 across the three applied methods. The vast majority of the samples show very low levels of
325 contamination ($\leq 5\%$) across all methods. A total of 33 samples, however, display
326 contamination estimates $> 5\%$ in one or more of the methods, but there are considerable
327 inconsistencies between the methods (**Supplementary Table V**). It is well known that
328 contamination estimates are not reliable for very low coverage genomic data (refs) and this is
329 further complicated by DNA damage in the sequences. Indeed, we observe that the 33
330 potentially problematic samples have an average coverage of 0.11X and a median of 0.06X
331 which is considerably lower than the full dataset with its average of coverage of 0.75X and
332 median of 0.26X. So, instead of simply excluding data from these precious samples based
333 on estimates that are likely imprecise, we “flagged” samples as potentially contaminated in
334 downstream analyses and took a more analytical approach in the evaluation. Flagging was
335 applied as follows:

336

337 **Samples with nuclear coverage $< 0.1X$ and MT coverage $< 10X$:** Not flagged, since both
338 estimates are likely unreliable

339

340 **Samples with nuclear coverage <0.1X and MT coverage ≥10X:** Flagged as contaminated
341 if any MT estimate is >5%; ignore nuclear estimate as likely unreliable

342

343 **Male samples with nuclear coverage ≥0.1X:** Flagged as contaminated if any nuclear (X-
344 chromosomal) estimate is > 5%; ignore MT as only nuclear data are relevant for genome-
345 wide analyses

346

347 **Female samples with nuclear coverage ≥0.1X:** Flagged as contaminated if any MT
348 estimate is >5% as no nuclear estimate is available

349

350 Based on this approach we have a total of 15 samples (NEO1, NEO3, NEO76, NEO77,
351 NEO158, NEO162, NEO168, NEO221, NEO226, NEO537, NEO657, NEO671, NEO677,
352 NEO746, NEO815) that we have flagged as “possibly contaminated” (See Table S2) in our
353 downstream analyses. The further analytical evaluation of these samples is described in
354 Supplement Note **S3d**.

355

356 References

- 357 1. Allentoft, M. E. *et al.* Population genomics of Bronze Age Eurasia. *Nature* **522**, 167–
358 172 (2015).
- 359 2. Damgaard, P. B. *et al.* Improving access to endogenous DNA in ancient bones and
360 teeth. *Sci. Rep.* **5**, 11184 (2015).
- 361 3. Edson, S. M. *et al.* Sampling of the cranium for mitochondrial DNA analysis of human
362 skeletal remains. *Forensic Science International: Genetics Supplement Series* **2**, 269–
363 270 (2009).
- 364 4. Gamba, C. *et al.* Genome flux and stasis in a five millennium transect of European
365 prehistory. *Nat. Commun.* **5**, 5257 (2014).
- 366 5. Knapp, M., Clarke, A. C., Horsburgh, K. A. & Matisoo-Smith, E. A. Setting the stage--
367 Building and working in an ancient DNA laboratory. *Annals of Anatomy-Anatomischer*
368 *Anzeiger* **194**, 3–6 (2012).
- 369 6. Fulton, T. L. & Shapiro, B. Setting Up an Ancient DNA Laboratory. *Methods Mol. Biol.*

370 **1963**, 1–13 (2019).

371 7. Orlando, L. *et al.* Ancient DNA analysis. *Nature Reviews Methods Primers* **1**, 1–26
372 (2021).

373 8. Higgins, D., Kaidonis, J., Townsend, G., Hughes, T. & Austin, J. J. Targeted sampling
374 of cementum for recovery of nuclear DNA from human teeth and the impact of common
375 decontamination measures. *Investig. Genet.* **4**, 18 (2013).

376 9. Hansen, H. B. *et al.* Comparing Ancient DNA Preservation in Petrous Bone and
377 Tooth Cementum. *PLoS One* **12**, e0170940 (2017).

378 10. Rohland, N. & Hofreiter, M. Comparison and optimization of ancient DNA
379 extraction. *Biotechniques* **42**, 343–352 (2007).

380 11. Margaryan, A. *et al.* Population genomics of the Viking world. *Nature* **585**,
381 390–396 (2020).

382 12. Meyer, M. & Kircher, M. Illumina sequencing library preparation for highly
383 multiplexed target capture and sequencing. *Cold Spring Harb. Protoc.* **2010**,
384 db.prot5448 (2010).

385 13. Hosseini, P., Tremblay, A., Matthews, B. F. & Alkharouf, N. W. An efficient
386 annotation and gene-expression derivation tool for Illumina Solexa datasets. *BMC Res.*
387 *Notes* **3**, 183 (2010).

388 14. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–
389 Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).

390 15. Schubert, M. *et al.* Improving ancient DNA read mapping against modern
391 reference genomes. *BMC Genomics* **13**, 178 (2012).

392 16. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools.
393 *Bioinformatics* **25**, 2078–2079 (2009).

394 17. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing
395 genomic features. *Bioinformatics* **26**, 841–842 (2010).

396 18. Jónsson, H., Ginolhac, A., Schubert, M., Johnson, P. L. F. & Orlando, L.
397 mapDamage2.0: fast approximate Bayesian estimates of ancient DNA damage

- 398 parameters. *Bioinformatics* **29**, 1682–1684 (2013).
- 399 19. Renaud, G., Slon, V., Duggan, A. T. & Kelso, J. Schmutzi: estimation of
400 contamination and endogenous mitochondrial consensus calling for ancient DNA.
401 *Genome Biol.* **16**, 224 (2015).
- 402 20. Fu, Q. *et al.* A revised timescale for human evolution based on ancient
403 mitochondrial genomes. *Curr. Biol.* **23**, 553–559 (2013).
- 404 21. Korneliussen, T. S., Albrechtsen, A. & Nielsen, R. ANGSD: Analysis of Next
405 Generation Sequencing Data. *BMC Bioinformatics* **15**, 356 (2014).
- 406

2) Imputation of ancient DNA

Bárbara Sousa da Mota^{1,2}, Andrew Vaughn³ & Olivier Delaneau^{1,2}

¹ Department of Computational Biology, University of Lausanne, Switzerland

² Swiss Institute of Bioinformatics, University of Lausanne, Switzerland

³ Center for Computational Biology, University of California, Berkeley, USA

Introduction

GLIMPSE is a statistical method developed to impute low-coverage human genomes. It has been shown that GLIMPSE efficiently produces accurate results when employed to impute low-coverage present-day genomes. Here we seek to demonstrate that GLIMPSE is also a suitable imputation tool of low-coverage ancient genomes. To benchmark it, we used a subset of 42 previously published ancient genomes with a mean depth of coverage above 10x (sample list is given in [Table S2.1](#)) that we downsampled to lower coverages in order to match the coverage we observed for the 317 genomes sequenced in the present study. Then, we imputed the resulting genomes and assessed the accuracy of the imputed calls by comparing with the original high-coverage genomes. Specifically, we examined the imputation performance regarding (i) depth of coverage, (ii) minor allele frequency, (iii) ancestry and living period of target samples. We also imputed the samples with the present 'gold standard' method, Beagle4.1¹, to show how it compares with GLIMPSE v1.0.1² (<https://odelaneau.github.io/GLIMPSE/>).

Methods

We first prepared all necessary files for the imputation step. We used samtools 1.10 to downsample the 42 high-coverage genomes to coverages 0.1x, 0.2x, 0.4x, 0.8x, 1.0x, 2.0x and 4.0x. Then, we prepared a list of candidate variant sites at which imputation was performed by retaining all sites in 1000 Genomes version 3 that were (i) bi-allelic SNPs and (ii) non-singleton (i.e. informative for imputation). For each of the seven tested coverages, we computed genotype likelihoods (VCF/PL field) at all candidate variant sites across all target samples using bcftools 1.10. To minimise computation time, we restricted this data generation procedure to chromosome 20.

Then, we performed imputation of all the resulting VCF files. We first divided chromosome 20 into 35 chunks with size between 1Mb and 2Mb. To prevent edge effects, we also added

442 additional buffer regions of 200kb on each side. Splitting chromosome 20 according to these
443 parameters was done using GLIMPSE_chunk v1.0.1. We then performed the imputation with
444 GLIMPSE_phase using the reference panel 1000 Genomes version 3, a cosmopolitan
445 collection of whole genome sequenced modern samples that we feel well adapted to the
446 various ancestries included in our data set. GLIMPSE_phase was run using the following
447 parameters: 10 burn-in iterations (*--burn 10*), 15 main iterations (*--main 15*) and a depth of 2
448 for the conditional state selection based on Positional Burrows-Wheeler transform (*--pbwt-*
449 *depth 2*). Finally, we ligated all imputed chunks back together into chromosome-wide VCF
450 files using GLIMPSE_ligate v1.0.1.

451

452 In addition to this, we also performed imputation with Beagle 4.1 with exactly the same input
453 data and reference panel and set its parameters *modelscale* to 2 and *niterations* to 0. These
454 parameter settings allow Beagle v4.1 to run with tractable running times on the data while
455 retaining good accuracy. The chunks of data imputed with Beagle v4.1 were ligated together
456 with bcftools concat 1.10.

457

458 To evaluate imputation performance, we employed GLIMPSE_concordance v1.0.1 using as
459 validation set all genotypes in high-coverage data that were covered by at least eight
460 sequencing reads and at which the most likely genotype was at least 1,000 times more likely
461 than the second best given the genotype likelihoods reported in the VCF/PL fields.

462 Specifically, we computed (i) the squared correlation and (ii) the concordance between
463 imputed and validation genotypes. For (i), we compared minor allele dosages (VCF/DS field)
464 within multiple minor allele frequency (MAF) bins. For (ii), we compared best guess
465 genotypes (VCF/GT field) and stratified the results depending on the type of validation
466 genotype: homozygous reference allele, heterozygous and homozygous alternative allele. As
467 further validation, we increased genomic coverage to 27.5X, 18.9X and 5.4X by deep
468 sequencing a previously published family trio (mother, father, son) from the Late Neolithic
469 mass burial at Koszyce in Poland³. This presented an opportunity to validate imputed
470 genotypes and haplotypes on the basis of Mendel's rules of inheritance⁴.

471 Results

472 In Fig. S2.1 and S2.2, we present the imputation accuracy per downsampled genome to 1.0x
473 for chromosome 20. We divided samples into eight classes based on expected genetic
474 proximity and plotted each class separately. In Fig. S2.1, the minor allele frequency (MAF)
475 for each of these groups was estimated from the 1000G reference panel, using European,
476 African, South East Asian, East Asian or American allele frequencies according to the place

477 of origin of samples. As expected, imputation accuracy decreases as minor allele frequency
478 decreases too. For common variants ($MAF \geq 5\%$), imputation accuracy is remarkably high
479 (>0.9) and closely matches what is usually obtained for modern samples. An exception to
480 this are African samples which exhibit lower accuracy in some cases (especially for baa01).
481 This likely results from the reference panel we used that does not represent well the
482 underlying ancestries of these samples. In Fig. S2.2, we present imputation accuracy per
483 genome as genotype discordance: the fraction of validation genotypes incorrectly imputed
484 stratified by homozygous and heterozygous genotypes. As expected, homozygous
485 reference alleles exhibit lower error rates than heterozygous and homozygous alternative
486 alleles. Error rates are remarkably low: less than 1% overall and less than 5% for the most
487 challenging genotypes to impute (RA and AA). Again, African samples exhibit much higher
488 error rates.

489

490 In Fig. S2.3, we present how imputation accuracy varies for all 42 samples depending on
491 multiple factors expected to affect accuracy. First, we look at coverage and find that 0.4x
492 coverage is enough to get 0.9 imputation accuracy at common variants ($MAF \geq 10\%$). Of note,
493 even 0.1x allows reaching 0.8 at common variants. Second, we considered whether imputing
494 the 42 jointly with the remaining 1,622 low coverage samples could decrease imputation
495 accuracy and did not find evidence of this happening: we get very similar accuracy results.
496 Finally, we check how GLIMPSE imputation does compare to Beagle 4.1 in case of ancient
497 low coverage samples and find that GLIMPSE brings substantial accuracy boost across the
498 entire frequency range. In Fig. S2.4, we show the phasing and imputation performance we
499 obtained across the entire genome for multiple coverages (0.1x to 4x). We obtained
500 Mendelian error rates from 0.1% at 4X to 0.55% at 0.1X (Extended Figure 1E,
501 Supplementary figure S2.4A). When looking only at sites at which at least one sample is
502 heterozygote (i.e. excluding triple homozygotes), we find that Mendel error rate ranges from
503 less than 2% at 4x and up to 8-10% for 0.1x (Supplementary figure S2.4B). Similarly, we
504 obtained switch error rates between 2% and 6%. Altogether, our validation analysis showed
505 that ancient European genomes can be imputed confidently from coverages above 0.4x and
506 highly valuable data can still be obtained with coverages as low as 0.1X when using specific
507 QC on the imputed data.

508

509 *Imputation of the full dataset*

510 Given the outcome of the benchmarking described above, we then proceeded with the
511 imputation of the full dataset. In total, we retained 1,664 samples with at least 0.1x mean
512 coverage. Similarly as before, we extracted all variable positions in 1000 Genomes version 3
513 that correspond to non-singleton bi-allelic SNPs and call genotype likelihoods at all these

514 variants for all samples using bcftools v1.10, thereby resulting in a VCF containing data for
515 1,664 samples across 43,285,119 SNPs. Of note, the reference genome used in this
516 analysis was hg19, b37. We then used GLIMPSE_chunk to split all the data into 1,841
517 chunks of 1Mb to 2Mb with overlapping 200kb buffers on each side. All these chunks of data
518 were imputed using GLIMPSE_phase v1.0.1 with 1000 Genomes version 3 as a reference
519 panel of haplotypes. Imputed chunks were ligated back together using GLIMPSE_ligate
520 v1.0.1, resulting in chromosome-wide VCF files containing the following information: (i) best
521 genotype guesses (VCF/GT field), (ii) expected non-reference allele dosage (VCF/DS field),
522 (iii) genotype posterior probabilities /VCF/GP field) and (iv) haplotype pairs sampled during
523 imputation (VCF/HS field). Finally, we used GLIMPSE_sample v1.0.1 to produce consensus
524 haplotype calls at all variants for all samples from the VCF/HS field.
525

526 Effects of Low coverage

527 It should be noted, however, that certain issues may arise when using this imputation on very
528 low coverage data. Specifically, the imputation errors GLIMPSE makes with low coverage
529 data tend to predominantly occur by incorrectly filling in the major allele at a given SNP. This
530 is not an issue specific to GLIMPSE itself but is instead inherent to any kind of Bayesian
531 approach to imputation. In the absence of informative data about the allele at a particular
532 SNP, imputation methods will fall back on the reference panel, or a set of confidently imputed
533 genomes, for imputation, which will tend to fill in missing data with the major allele at each
534 SNP.

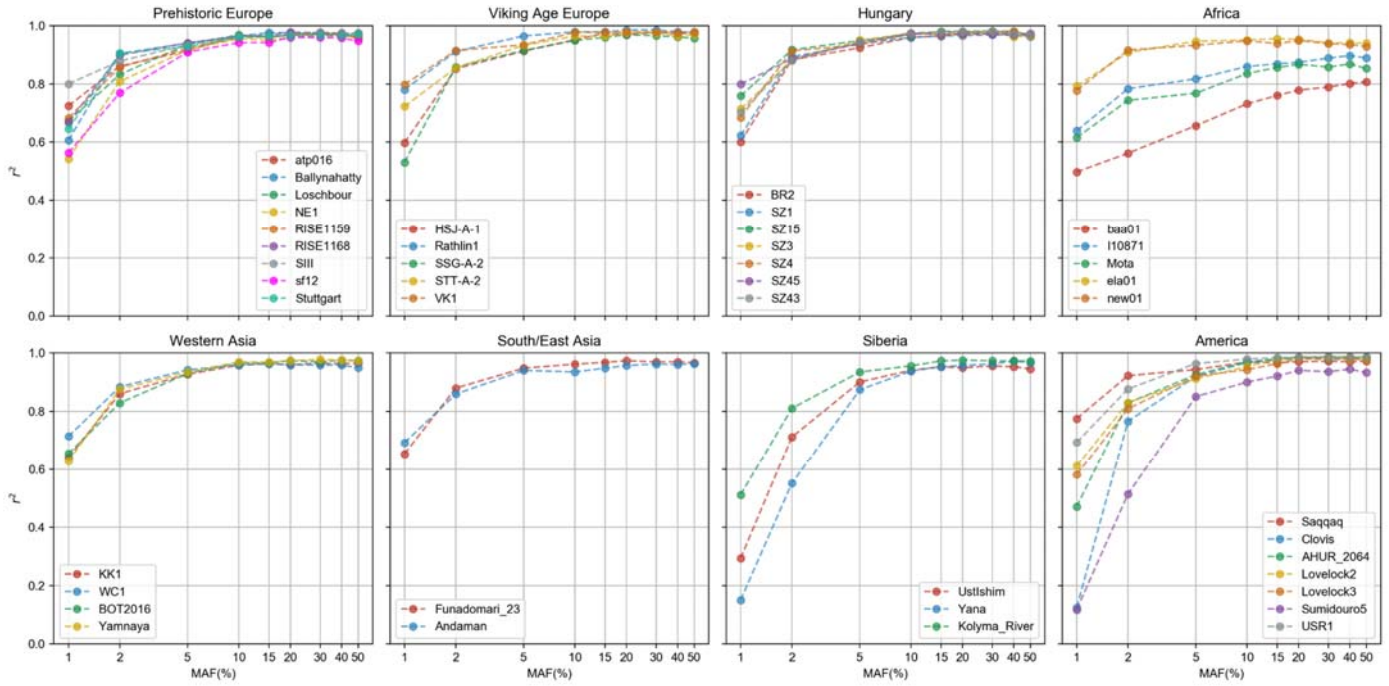
535

536 To illustrate this phenomenon, we took the 1,492 samples that passed all filters (described in
537 detail in Supplementary Note 3d) and retained only SNPs passing the 1000 Genomes strict
538 mask. For computational considerations, we considered only chromosome 8, which gave
539 1,139,150 total SNPs. In Figure S2.5, we plotted, for each of the 2,984 haplotypes, the total
540 number of allele differences between that genome and the reference sequence against the
541 coverage of that sample. Computing Spearman's ρ showed a substantial correlation, which
542 persisted after filtering for SNPs on INFO > 0.5, leaving 584,280 SNPs, and INFO > 0.8,
543 leaving 336,842 SNPs (Figures S2.6 and S2.7 respectively). Visually, it appeared that this
544 correlation was driven by a reduced number of differences to the reference sequence among
545 very low coverage genomes, and we confirmed this by considering SNPs of all INFO scores
546 and noticing that the correlation decreases sharply when samples below 0.3x are dropped
547 (Figure S2.9) and can be decreased even further by dropping samples below 2x coverage
548 (Figure S2.10). We also confirmed that the correlation is very small when only retaining very

549 high INFO score SNPs, specifically $\text{INFO} > 0.97$, which retained 56,925 SNPs (Figure S2.8).
550 This is to be expected, as the imputation for high-confidence SNPs should show no
551 significant biases.

552

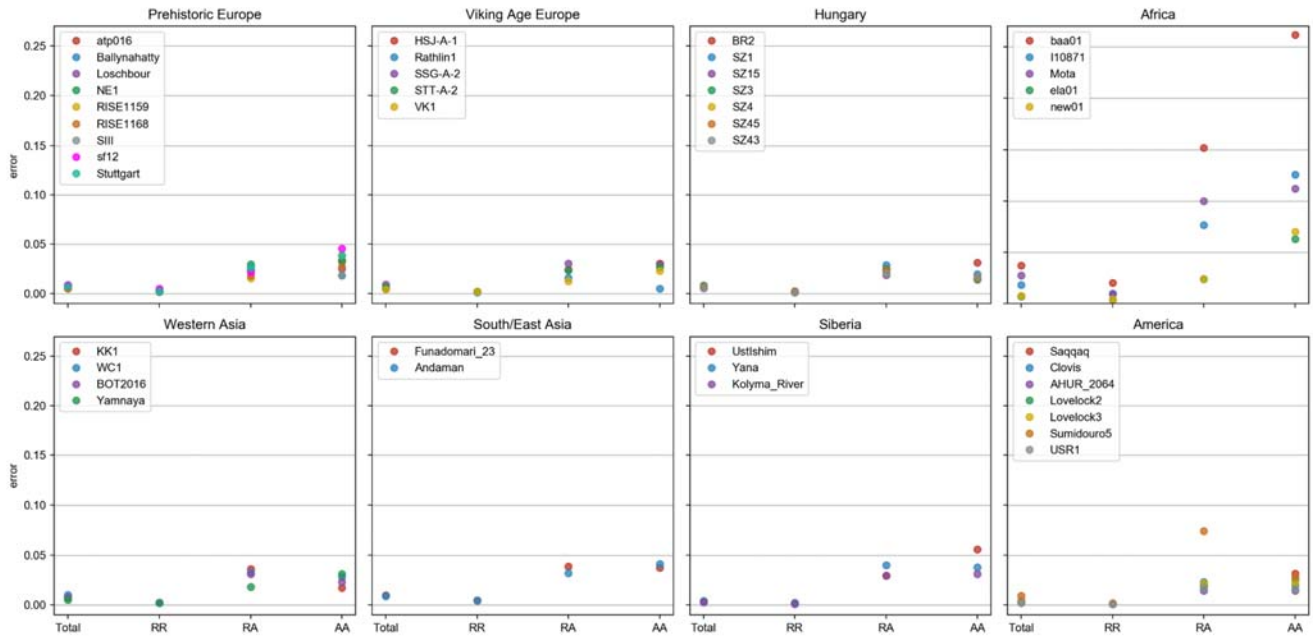
553 This phenomenon, where coverage can be predictive of sequence, is important to keep in
554 mind when running certain types of analyses on these data. Analyses such as PCA and
555 admixture modelling, which mainly rely on common SNPs that are shared among many
556 individuals, are not expected to be significantly affected, as imputation is quite accurate for
557 SNPs with high MAF (see Supplementary Note 3d for a thorough analysis of how coverage
558 affects PCA). However, this observation has important implications for genealogy
559 reconstruction and other analyses that rely on overall sequence similarity or otherwise utilise
560 rare SNPs. We recommend that researchers using these imputed data carefully consider
561 what effect the inclusion of low coverage samples might have on their analyses and then
562 utilise appropriate MAF/INFO filters on SNPs and/or coverage filters on samples as
563 necessary.



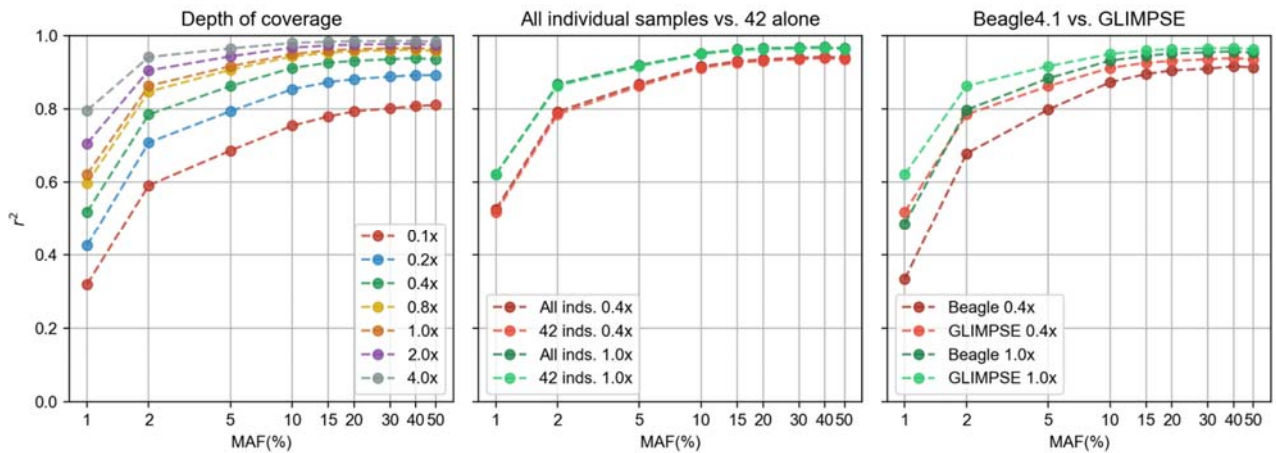
564

565 **Figure S2.1: Per-sample imputation accuracy (1).** Imputation accuracy as squared
 566 correlation between imputed and validation genotypes (y-axis). Samples were grouped into
 567 eight broad categories based on genetic proximity. Each of the plots corresponds to one of
 568 such categories. Minor allele frequencies (MAF; x-axis) were estimated from 1000 Genomes
 569 version 3 for matched continental groups (see Table S2.1 for details).

570



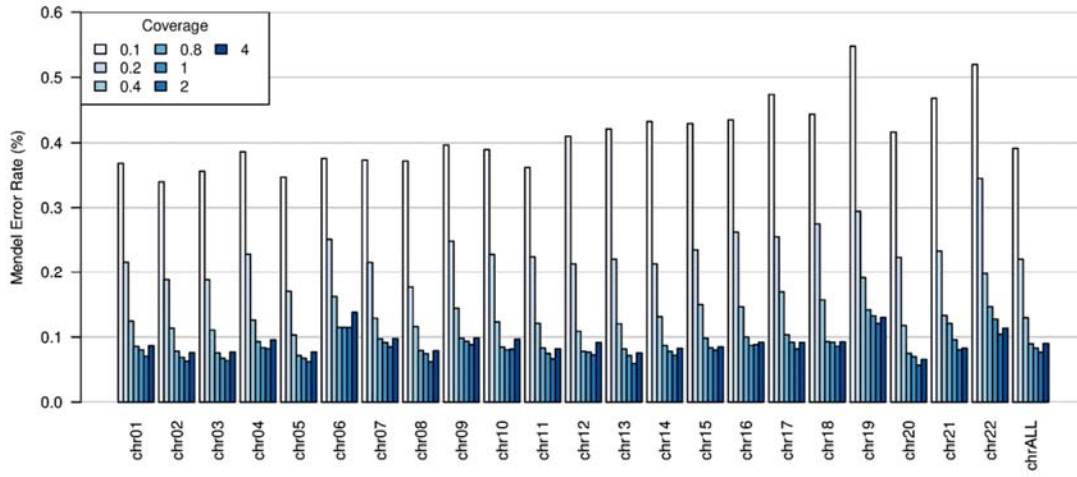
571 **Figure S2.2: Per-sample imputation accuracy (2).** Imputation accuracy as discordance
 572 between imputed and validation genotypes (y-axis). Samples were grouped into eight broad
 573 categories based on genetic proximity. Each of the plots corresponds to one of such
 574 categories. We report results across four types of genotypes: (i) Total; all genotypes
 575 together, (ii) RR; validation genotypes with two copies of the reference allele, (iii) RA;
 576 heterozygous genotypes (iv) AA; validation genotypes with two copies of the alternative
 577 allele.
 578



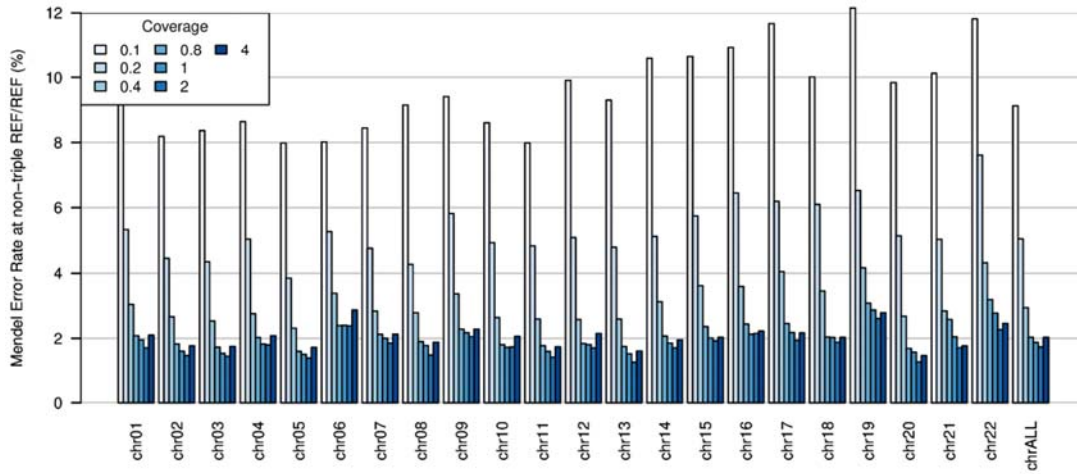
579 **Figure S2.3: Main parameters affecting imputation accuracy.** Imputation accuracy of
 580 GLIMPSE across the 42 samples regarding, from left to right: (i) sequencing coverage; (ii)
 581 effect of jointly imputing all 1.6K individual samples compared to imputing only the 42
 582 downsampled genomes, (iii) imputation done with Beagle4.1.
 583
 584

586 **Figure S2.4: Imputation and phasing accuracy for the Koszyce trio.** (A) Mendel error
587 rate across the 22 autosomes. A Mendel error is counted when the parental and offspring
588 genotypes violate Mendel transmission rules. (B) Same as before, excluding sites at which
589 all three samples are REF/REF in the high coverage data. (C) Switch error rates averaged
590 over the 3 samples. A switch error is counted between two consecutive heterozygous

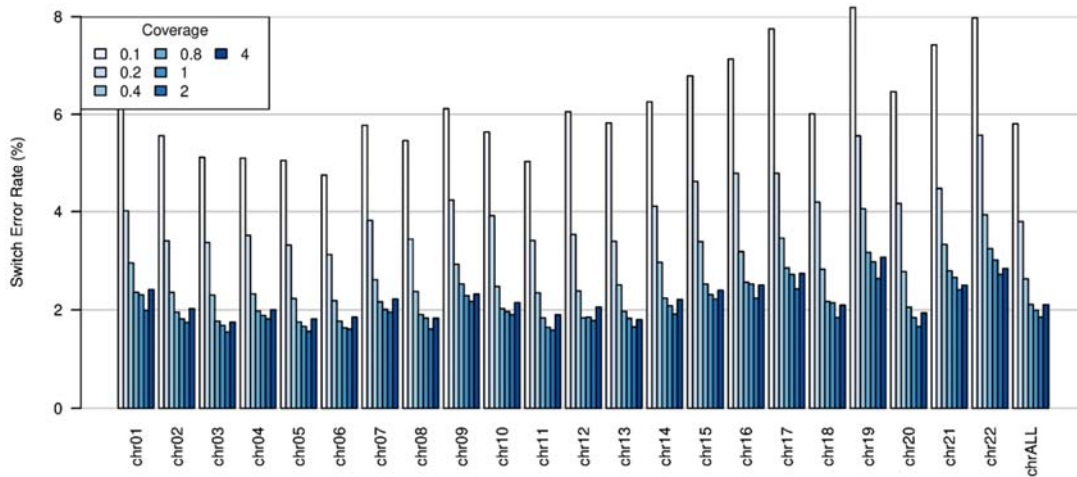
A. Mendel error rate at all sites



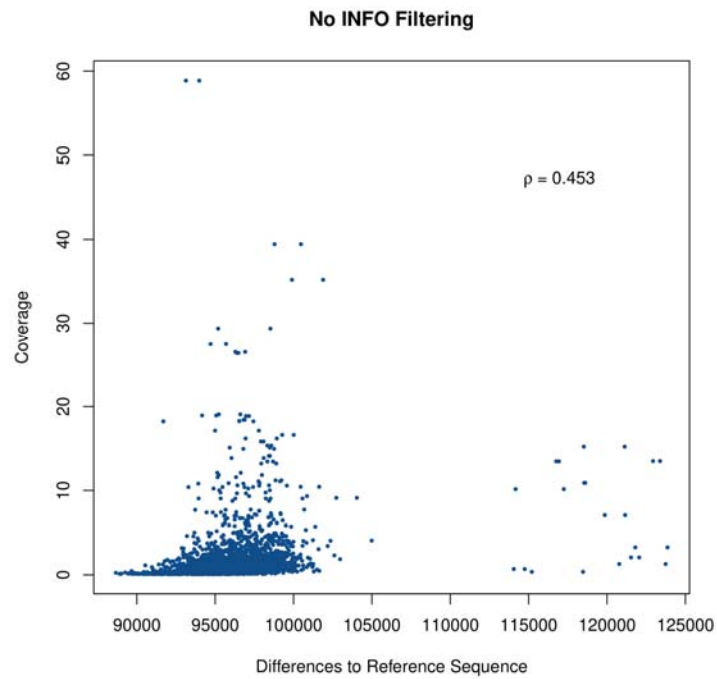
B. Mendel error rate at variable sites



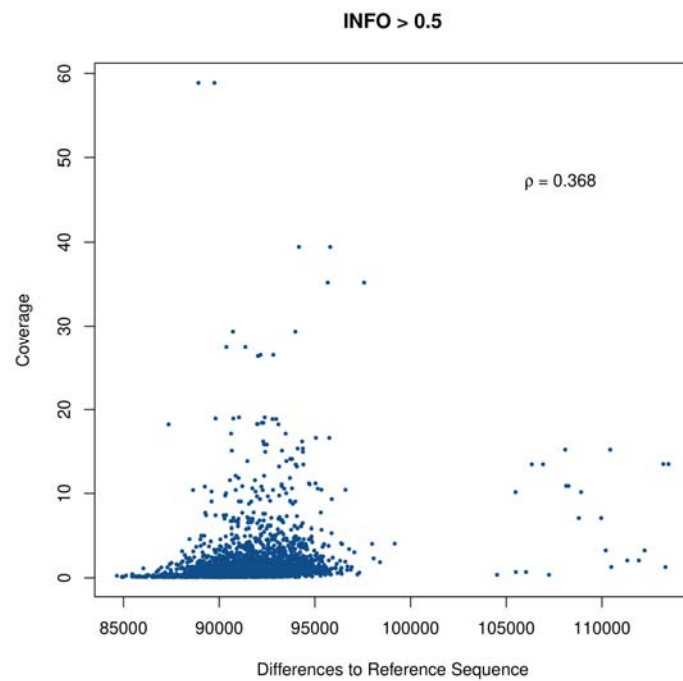
C. Switch error rate



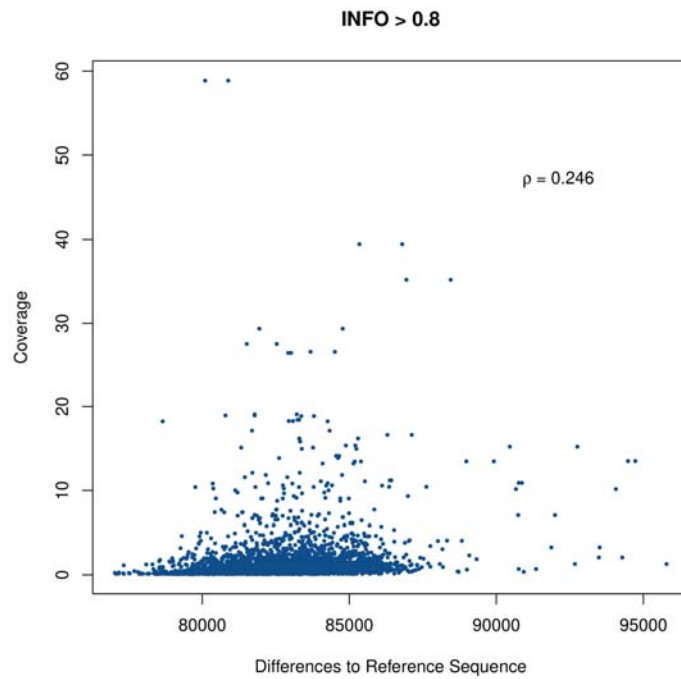
591 genotypes when the reported haplotypes are not consistent with those derived from the trio.



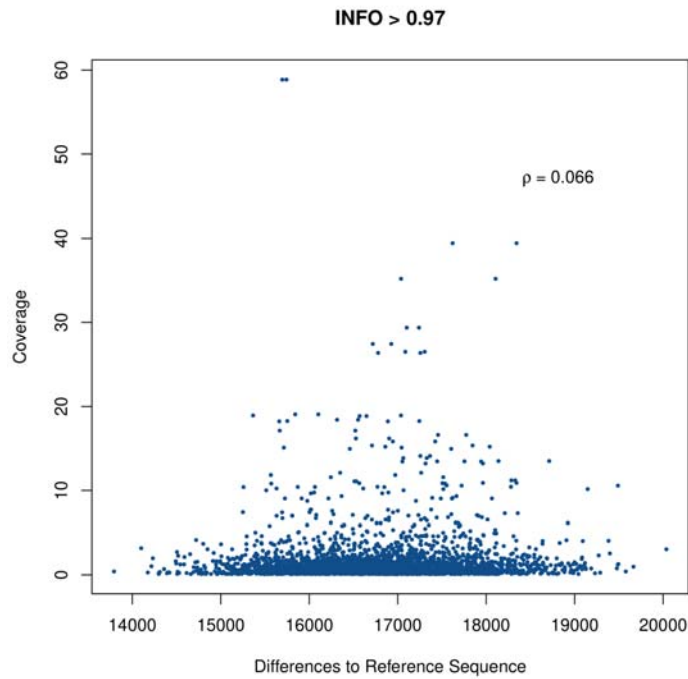
592 **Figure S2.5.** Total number of allele differences to the reference sequence for each of the
 593 2,984 haplotypes against coverage, without INFO filtering.



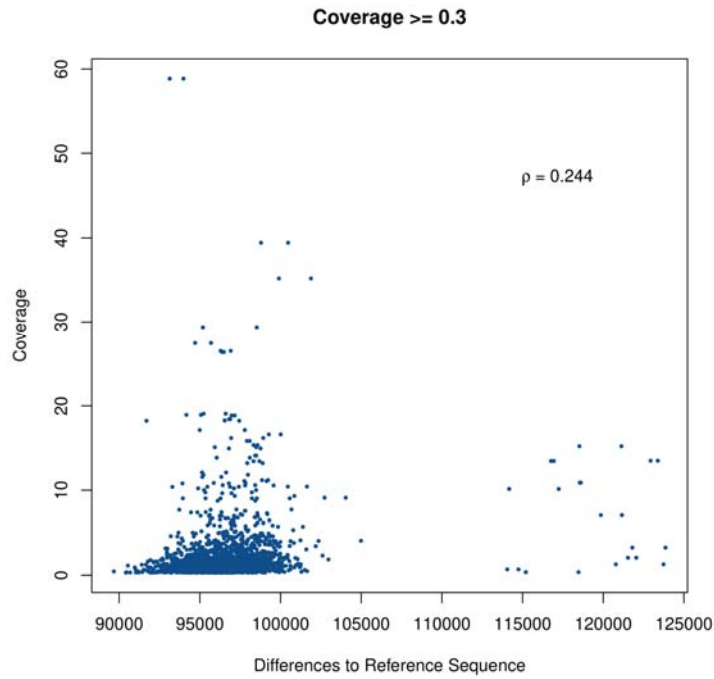
594 **Figure S2.6.** Total number of allele differences to the reference sequence for each of the
 595 2,984 haplotypes against coverage, with filtering for SNPs on INFO > 0.5.



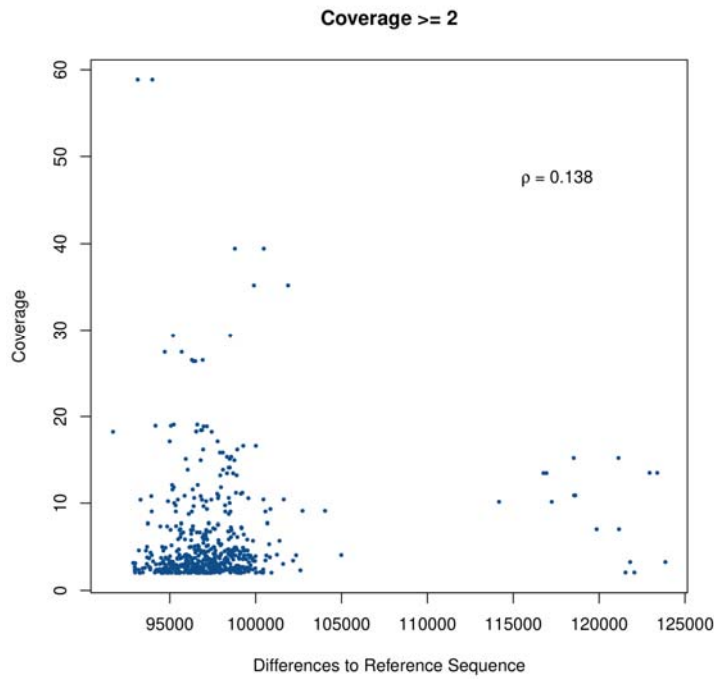
596 **Figure S2.7.** Total number of allele differences to the reference sequence for each of the
 597 2,984 haplotypes against coverage, with filtering for SNPs on INFO > 0.8.



598 **Figure S2.8.** Total number of allele differences to the reference sequence for each of the
 599 2,984 haplotypes against coverage, with filtering for SNPs on INFO > 0.97.



600 **Figure S2.9.** Total number of allele differences to the reference sequence for each of the
 601 2,984 haplotypes against coverage, with filtering for samples on coverage ≥ 0.3 .



602 **Figure S2.10.** Total number of allele differences to the reference sequence for each of the
 603 2,984 haplotypes against coverage, with filtering for samples on coverage ≥ 2 .
 604

605
606
607
608

SampleID	Country	Age (yBP)	MAF 1000G group	Coverage	Reference
atp16	Spain	4867 – 5212	EUR	13X	Günther et al. 2015 ⁵
Stuttgart	Germany	7020 – 7260	EUR	16X	Lazaridis et al. 2014 ⁶
Loschbour	Luxembourg	7940 – 8160	EUR	18X	Lazaridis et al. 2014 ⁶
Ballynahatty	Ireland	4970 – 5293	EUR	10X	Cassidy et al. 2016 ⁷
sf12	Sweden	8757 – 9033	EUR	59X	Günther et al. 2018 ⁸
NE1	Hungary	7021 – 7256	EUR	18X	Gamba et al. 2014 ⁹
Sunghir III	Russia	33031 – 35154	EUR	11X	Sikora et al. 2017 ¹⁰
Rathlin1	Ireland	3835 – 3976	EUR	11X	Cassidy et al. 2016 ⁷
SSG-A2	Iceland	950 – 1100	EUR	10X	Ebenesersdóttir et al. 2018 ¹¹
HSJ-A1	Iceland	950 – 1080	EUR	29X	Ebenesersdóttir et al. 2018 ¹¹
STT-A2	Iceland	950 – 1050	EUR	14X	Ebenesersdóttir et al. 2018 ¹¹
VK1	Greenland	750 – 950	EUR	12X	Margaryan et al. 2020 ¹²
BR2	Hungary	3060 – 3220	EUR	18X	Gamba et al. 2014 ⁹
SZ15	Hungary	1346 – 1538	EUR	11X	Amorim et al. 2018 ¹³
SZ3	Hungary	1346 – 1538	EUR	11X	Amorim et al. 2018 ¹³
SZ4	Hungary	1347 – 1538	EUR	10X	Amorim et al. 2018 ¹³
SZ45	Hungary	1348 – 1538	EUR	10X	Amorim et al. 2018 ¹³
SZ43	Hungary	1349 – 1538	EUR	12X	Amorim et al. 2018 ¹³
SZ1	Hungary	1150 – 1350	EUR	11X	Amorim et al. 2018 ¹³
baa01	South Africa	1831 – 1986	AFR	14X	Schlebusch et al. 2017 ¹⁴
ela01	South Africa	453 – 533	AFR	13X	Schlebusch et al. 2017 ¹⁴
new01	South Africa	327 – 508	AFR	11X	Schlebusch et al. 2017 ¹⁴
I10871	Cameroon	7800 – 7970	AFR	15X	Lipson et al. 2020 ¹⁵
Mota	Ethiopia	4419 – 4525	AFR	10X	Gallego Llorente et al. 2015 ¹⁶

KK1	Georgia	9550 – 9890	EUR	12X	Broushaki et al. 2016 ¹⁷
WC1	Iran	9032 – 9405	EUR	10X	Broushaki et al. 2016 ¹⁷
BOT2016	Kazakhstan	5318 – 5582	EUR	14X	de Barros Damgaard et al. 2018 ¹⁸
Yamnaya Karagash	Kazakhstan	4837 – 4968	EUR	26X	de Barros Damgaard et al. 2018 ¹⁸
Andaman	India	30 – 150	SAS	17X	Moreno-Mayar et al. 2018 ¹⁹
Funadomari 23	Japan	3550 – 3960	EAS	39X	Kanzawa-Kiriyama et al. 2019 ²⁰
Ust'-Ishim	Russia	42560 – 47480	ALL	35X	Fu et al. 2014 ²¹
Yana 1	Russia	30950 – 32950	ALL	27X	Sikora et al. 2019 ²²
Kolyma 1	Russia	9665 – 9906	ALL	15X	Sikora et al. 2019 ²²
USR1	USA	11270 – 11600	ALL	17X	Moreno-Mayar et al. 2018 ²³
AHUR_2064	USA	10770 – 11170	AMR	19X	Moreno-Mayar et al. 2018 ¹⁹
Lovelock2	USA	1818 – 1942	AMR	15X	Moreno-Mayar et al. 2018 ¹⁹
Lovelock3	USA	567 – 687	AMR	19X	Moreno-Mayar et al. 2018 ¹⁹
Saqqaq	Greenland	3600 – 4170	AMR	13X	Rasmussen et al. 2010 ²⁴
Clovis	USA	12572 – 12726	AMR	15X	Rasmussen et al. 2014 ²⁵
Sumidouro5	Brazil	10258 – 10552	AMR	16X	Moreno-Mayar et al. 2018 ¹⁹
RISE1159 *	Poland	4840 – 4709	EUR	27X	Schroeder et al. 2019 ³
RISE1168 *	Poland	4840 – 4709	EUR	19X	Schroeder et al. 2019 ³
RISE1160 *	Poland	4840 – 4709	EUR	5X	Schroeder et al. 2019 ³

609
610
611
612
613
614
615
616
617
618
619

Table S2.1: Detailed list of high coverage ancient validation genomes. From left to right: (1) original sample ID, (2) country of origin, (3) estimated age, (4) best matching continental group in 1000 Genomes used to stratify imputation accuracy results and (5) original coverage from which down-sampling has been performed. *indicates the Neolithic Koszyce trio that was first published in Schroeder et al. ³ but now sequenced to higher depth in this study and used for imputation validation purposes. RISE1160 is not counted among the 42 high coverage genomes.

620 References

- 621 1. Browning, B. L. & Browning, S. R. Genotype Imputation with Millions of Reference
622 Samples. *Am. J. Hum. Genet.* **98**, 116–126 (2016).
- 623 2. Rubinacci, S., Ribeiro, D. M., Hofmeister, R. J. & Delaneau, O. Efficient phasing and
624 imputation of low-coverage sequencing data using large reference panels. *Nat. Genet.*
625 **53**, 412 (2021).
- 626 3. Schroeder, H. *et al.* Unraveling ancestry, kinship, and violence in a Late Neolithic
627 mass grave. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 10705–10710 (2019).
- 628 4. Bateson, W. *Mendel's Principles of Heredity.* (Cambridge University Press, 1902).
- 629 5. Günther, T. *et al.* Ancient genomes link early farmers from Atapuerca in Spain to
630 modern-day Basques. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 11917–11922 (2015).
- 631 6. Lazaridis, I. *et al.* Ancient human genomes suggest three ancestral populations for
632 present-day Europeans. *Nature* **513**, 409–413 (2014).
- 633 7. Cassidy, L. M. *et al.* Neolithic and Bronze Age migration to Ireland and establishment
634 of the insular Atlantic genome. *Proc. Natl. Acad. Sci. U. S. A.* **113**, 368–373 (2016).
- 635 8. Günther, T. *et al.* Population genomics of Mesolithic Scandinavia: Investigating early
636 postglacial migration routes and high-latitude adaptation. *PLoS Biol.* **16**, e2003703
637 (2018).
- 638 9. Gamba, C. *et al.* Genome flux and stasis in a five millennium transect of European
639 prehistory. *Nat. Commun.* **5**, 5257 (2014).
- 640 10. Sikora, M. *et al.* Ancient genomes show social and reproductive behavior of
641 early Upper Paleolithic foragers. *Science* **358**, 659–662 (2017).
- 642 11. Ebenesersdóttir, S. S. *et al.* Ancient genomes from Iceland reveal the making
643 of a human population. *Science* **360**, 1028–1032 (2018).
- 644 12. Margaryan, A. *et al.* Population genomics of the Viking world. *Nature* **585**,
645 390–396 (2020).
- 646 13. Amorim, C. E. G. *et al.* Understanding 6th-century barbarian social

- 647 organization and migration through paleogenomics. *Nat. Commun.* **9**, 3547 (2018).
- 648 14. Schlebusch, C. M. *et al.* Southern African ancient genomes estimate modern
649 human divergence to 350,000 to 260,000 years ago. *Science* **358**, 652–655 (2017).
- 650 15. Lipson, M. *et al.* Ancient West African foragers in the context of African
651 population history. *Nature* **577**, 665–670 (2020).
- 652 16. Gallego Llorente, M. *et al.* Ancient Ethiopian genome reveals extensive
653 Eurasian admixture in Eastern Africa. *Science* **350**, 820–822 (2015).
- 654 17. Broushaki, F. *et al.* Early Neolithic genomes from the eastern Fertile Crescent.
655 *Science* **353**, 499–503 (2016).
- 656 18. de Barros Damgaard, P. *et al.* The first horse herders and the impact of early
657 Bronze Age steppe expansions into Asia. *Science* **360**, (2018).
- 658 19. Moreno-Mayar, J. V. *et al.* Early human dispersals within the Americas.
659 *Science* **362**, (2018).
- 660 20. Kanzawa-Kiriyama, H. *et al.* Late Jomon male and female genome sequences
661 from the Funadomari site in Hokkaido, Japan. *Anthropol. Sci.* **127**, (2019).
- 662 21. Fu, Q. *et al.* Genome sequence of a 45,000-year-old modern human from
663 western Siberia. *Nature* **514**, 445–449 (2014).
- 664 22. Sikora, M. *et al.* The population history of northeastern Siberia since the
665 Pleistocene. *Nature* **570**, 182–188 (2019).
- 666 23. Moreno-Mayar, J. V. *et al.* Terminal Pleistocene Alaskan genome reveals first
667 founding population of Native Americans. *Nature* **553**, 203–207 (2018).
- 668 24. Rasmussen, M. *et al.* Ancient human genome sequence of an extinct Palaeo-
669 Eskimo. *Nature* **463**, 757–762 (2010).
- 670 25. Rasmussen, M., Anzick, S. L., Waters, M. R. & Skoglund, P. The genome of a
671 Late Pleistocene human from a Clovis burial site in western Montana. *Nature* (2014).
- 672

3) Demographic inference

3a) Phylogenetic analysis of mtDNA sequences

Tharsika Vimala¹, Martin Sikora¹

¹Lundbeck Foundation GeoGenetics Centre, GLOBE Institute, University of Copenhagen,
Copenhagen, Denmark

Methods

We carried out a phylogenetic tree analysis using the reconstructed mitochondrial genomes from the human remains presented in this study. Only sequences with less than 2000 missing sites were included in the analysis to avoid biases caused by missing data. Sequences were aligned using *mafft*¹ and subsequently inputted to the maximum likelihood (ML) based phylogenetic tree inference tool *RAxML-ng*². The analysis was carried out by using the 'all-in-one'- option performing both an ML search and bootstrapping (--all --bs-trees 100) along with the substitution model GTR+I+G4.

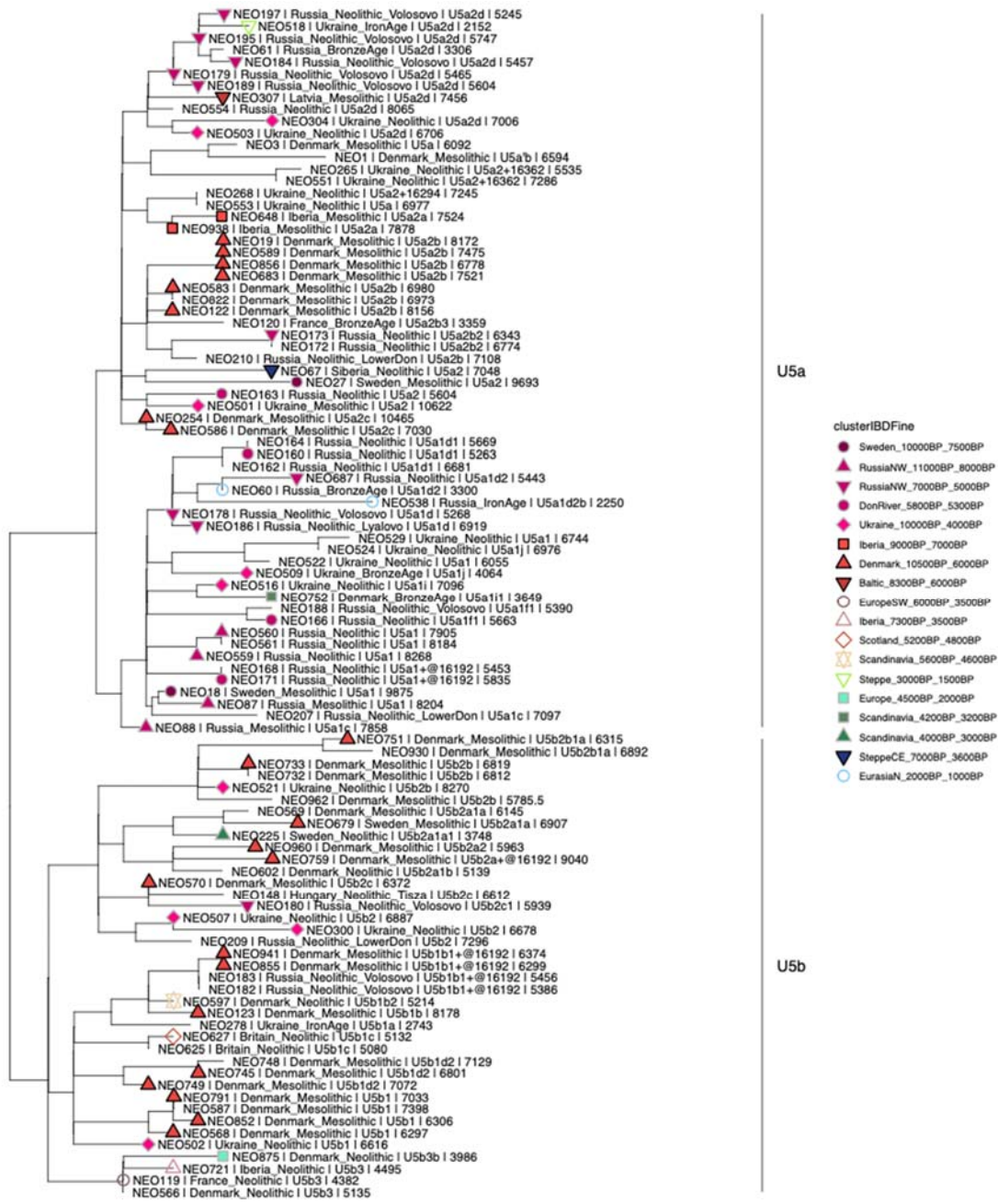
Results

From the resulting phylogenetic tree, we obtain an overview of how the remains are distributed across the haplogroups in our dataset. We find haplogroup U to be the most common haplogroup in the analysed set of individuals. Especially, subclade U5 is commonly observed among European hunter-gatherers. Focusing on the subclade U5a, we find the remains distributed into two main sub-haplogroups classified as U5a1 and U5a2. U5a1 is mainly influenced by Eastern European hunter-gatherers from Russia and Ukraine, while we also find the Scandinavian remains, NEO752 and NEO18, represented in this clade. In particular, the 9.8 kya old remains of NEO18 are interesting as the genetic structure analysis of the autosomes of NEO18 show evidence of Ukrainian hunter-gatherer ancestry (Figure 2, main). Haplogroup U5a2 shows a higher representation of Danish hunter-gatherers, specifically in subclade U5a2a, in which we also identify two Mesolithic Iberian individuals, NEO648 and NEO938. We likewise observe a Mesolithic Latvian individual (NEO307) within the U5a2d subclade, which is primarily dominated by Ukrainian and Russian hunter-gatherers (Figure S3a.1). These observations are congruent with the autosomal structure analysis displaying Ukrainian hunter-gatherer ancestry in these individuals. The clade

706 representing haplogroup U5b is mainly influenced by Danish hunter-gatherers along with a
707 few Western European hunter-gatherers from Britain, France, and Iberia. We do, however,
708 also observe a few Ukrainian hunter-gatherers clustering closely to the Danish hunter-
709 gatherer individuals, which could indicate a continuous level of gene flow from Eastern
710 Europe into Scandinavia. Additionally, we identify a single farmer individual, NEO597,
711 carrying U5b, which is a rare example of genetic continuity of a hunter-gatherer associated
712 haplogroup (Figure S3a.1). This contrasts the genetic transition otherwise observed with the
713 arrival of the early farmers. We find a similar overall pattern within the genetic variation of
714 U4, which is mainly influenced by Eastern European hunter-gatherers. Furthermore, we
715 identify individuals with steppe-ancestry as well as two Danish hunter-gatherer remains
716 clustering within the same clade. In the clade of haplogroup U2 we find a single Mesolithic
717 Iberian carrying haplogroup U2'3'4'7'8'9, while the rest of the remains in the U2-clade belong
718 to the sub-haplogroup U2e. U2e is carried by Eastern European hunter-gatherers, although
719 we also identify a significant number of remains from the forest steppe clustering in U2e as
720 well. The Danish Neolithic remains of NEO792 are interesting as this individual carried the
721 highest proportion of steppe-ancestry among the Danish individuals. Haplogroup K, a
722 descending haplogroup of U8, includes a combination of farmers and hunter-gatherers.
723 Specifically, we find the haplogroup K1e to be carried by Danish hunter-gatherers, while K1a
724 and K1b are mainly influenced by Neolithic individuals from Scandinavia (Figure S3a. 2). The
725 highest frequency of early farmers is found within the genetic variation of haplogroup H, a
726 descending haplogroup of HV. Both HV and JT are mainly influenced by Western European
727 farmers (Figure S3a.3).

728 Conclusion

729 In overall we find most of the Scandinavian hunter-gatherers clustering within the variation of
730 haplogroup U. Given the high number of human remains represented, we were able to obtain
731 a phylogeny of a relatively high resolution of this particular haplogroup. Our results show
732 evidence of a continuous migration of especially Eastern hunter-gatherers into Scandinavia.
733 Most of the early farmers carried haplogroup H or fell within haplogroup JT, while a few
734 farmers carried haplogroup K. We find individuals with steppe ancestry mainly clustering
735 together under macro haplogroup M, although we also identify a few steppe individuals in U2
736 closely related to a Danish Neolithic individual, NEO792, who also carried a high proportion
737 of steppe-ancestry in the nuclear genome.
738



739

740 **Figure S3a.1. Maximum likelihood tree of haplogroup U5**

741 Maximum likelihood tree displaying the phylogenetic relationship between the human
 742 remains carrying haplogroup U5 and more specifically the two subclades U5a and U5b.
 743 Labels include information on sample ID, group ID, haplogroup, and age of the respective
 744 remains. Symbols indicate the specific fine scale IBD cluster listed in the legend.

745

746

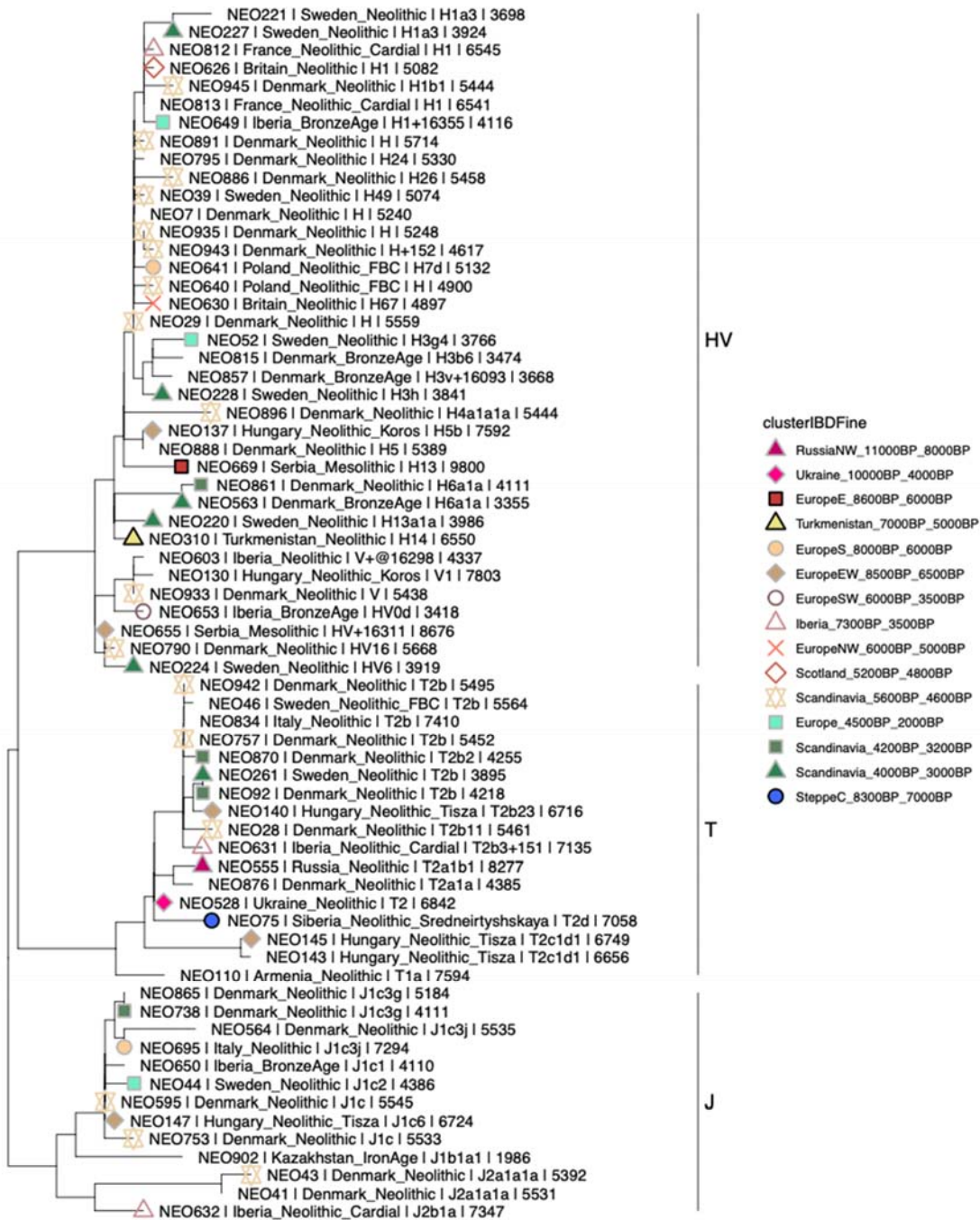
747

748

749 0

750 **Figure S3a.2. Maximum likelihood tree of haplogroup U excluding U5**

751 Maximum likelihood tree displaying the phylogenetic relationship between the human
752 remains carrying a descending haplogroup of U with the exception of U5. Labels include
753 information on sample ID, group ID, haplogroup, and age of the respective remains. Symbols
754 indicate the specific fine scale IBD cluster listed in the legend.
755



756

757 **Figure S3a.3. Maximum likelihood tree of haplogroup farmer-associated haplogroups**

758 Maximum likelihood tree displaying the phylogenetic relationship between the human
759 remains carrying an HV or JT descending haplogroup. Labels include information on sample
760 ID, group ID, haplogroup, and age of the respective remains. Symbols indicate the specific
761 fine scale IBD cluster listed in the legend.

762
763

764 References

- 765 1. Kazutaka Katoh, Kazuharu Misawa, Kei-ichi Kuma, Takashi Miyata, MAFFT: a novel
766 method for rapid multiple sequence alignment based on fast Fourier transform,
767 *Nucleic Acids Research*, Volume 30, Issue 14, 15 July 2002, Pages 3059–3066,
768 <https://doi.org/10.1093/nar/gkf436>
- 769 2. Alexey M Kozlov, Diego Darriba, Tomáš Flouri, Benoit Morel, Alexandros Stamatakis,
770 RAxML-NG: a fast, scalable and user-friendly tool for maximum likelihood
771 phylogenetic inference, *Bioinformatics*, Volume 35, Issue 21, 1 November 2019,
772 Pages 4453–4455, <https://doi.org/10.1093/bioinformatics/btz305>

773
774
775
776
777

778 3b) Y chromosome / Sex determination

779 Martin Sikora¹

780
781
782
783

¹Lundbeck Foundation GeoGenetics Centre, GLOBE Institute, University of Copenhagen,
Copenhagen, Denmark

784 Methods

785 We calculated the ratio of reads aligning to either of the sex chromosomes (R_Y statistic)¹ to
786 determine genetic sex of the study individuals. Y chromosomes of inferred male individuals
787 were further analysed using phylogenetic placement². We built a reference phylogenetic tree
788 of 1,244 male individuals from the 1000 Genomes project with *RAxML-NG*³, using the
789 general time-reversible model including among-site rate heterogeneity and ascertainment
790 correction (model GTR+G+ASC_LEWIS). For each ancient sample, haploid genotypes given
791 the positions and alleles in the reference panel were called using '*bcftools call*' (options '-C
792 *alleles -ploidy 1 -i*'). The resulting genotypes were converted to fasta format and placed onto

793 the reference tree using *EPA-ng*². Phylogenetic placements were processed and visualised
794 using *gappa*⁴.

795

796 To convert phylogenetic placements into haplogroup calls, we assigned each branch of the
797 reference phylogeny to its representing haplogroup, using SNP annotations from ISOGG
798 (version 15.73). For each ancient sample, haplogroups were then called using the most
799 basal branch accumulating 99% of the placement weights, obtained using the '*accumulate*'
800 command in *gappa*.

801

802 For in-depth phylogenetic analyses of haplogroups I, R1, and Q, we compiled extended
803 reference panels of high-coverage modern individuals belonging to those haplogroups from
804 publicly available sources⁵⁻⁸. To increase resolution for the placement of ancient samples,
805 we also included ancient individuals with Y chromosome coverage $\geq 1.5X$ in the reference
806 panels. For each haplogroup panel, we called haploid genotypes individually per sample
807 using '*bcftools call*', setting genotypes with read depth < 2 or quality score < 30 to missing.
808 Individual VCF files were then merged and filtered to retain only biallelic SNPs polymorphic
809 in the reference panel. For each haplogroup reference panel, we built phylogenetic trees
810 using *RAxML-NG* and performed phylogenetic placement as described above, restricting to
811 target samples with $> 0.1X$ coverage.

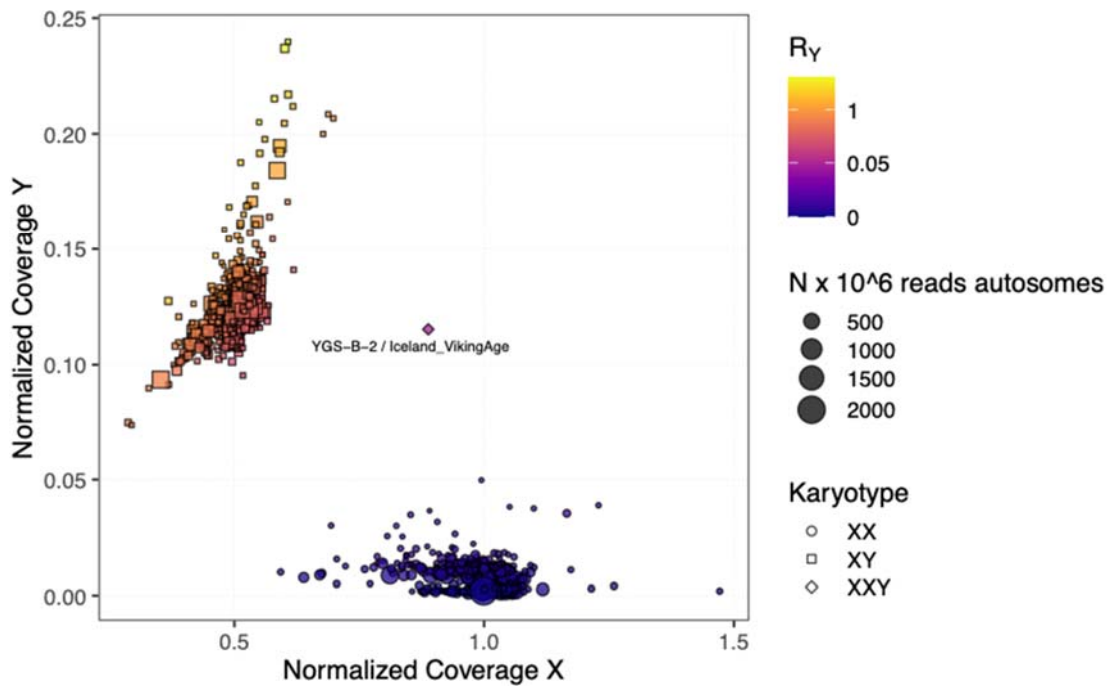
812

813 Results

814 Sex determination

815 We unambiguously determined genetic sex for all 317 study individuals (118 female, 199
816 males; **Supplementary Table VII**). In a plot of normalised sequencing depth across the X and
817 Y chromosomes, the final dataset individuals form two clearly separated clusters
818 corresponding to XX and XY karyotypes (**Fig. S3b.1**). The exception is individual YGS-B-2,
819 an Icelandic Viking Age individual previously found to carry an XXY karyotype⁹.

820



821

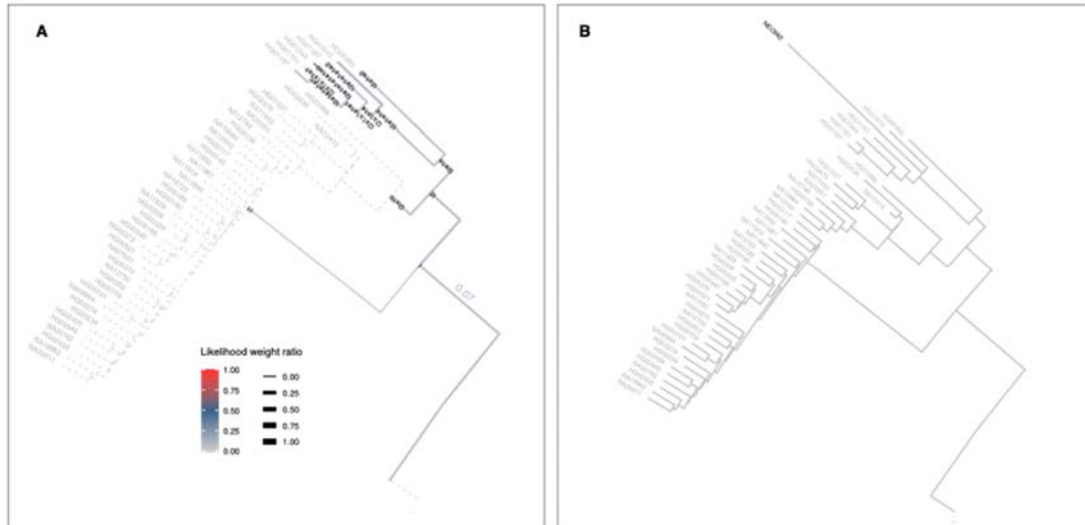
822 **Fig. S3b.1. Sex determination.** Plot shows coverage on X and Y chromosomes normalised
823 by autosomal coverage, for each individual. Symbol color indicates R_Y ration values, and
824 shape the inferred sex chromosome karyotype. Total number of autosomal reads are
825 indicated by symbol size.

826

827 Phylogenetic placement

828 We used phylogenetic placement to analyse Y chromosome diversity in our dataset. For
829 each ancient sample we obtain a distribution of placement weights across the reference
830 phylogeny, hereby incorporating uncertainty in the placement due to low coverage. The
831 placements can be subjected to analyses such as grafting as a pendant edge to the most
832 likely placement, or clustering of multiple samples. As Y chromosome haplogroups are labels
833 for clades of the phylogeny descending from specific ancestral branches, we can convert the
834 placements into haplogroup calls for a specific sample by assigning haplogroup labels to
835 each branch in the reference phylogeny, and finding the most basal branch that accumulates
836 placement weights up to a specified threshold for the sample. We chose a conservative
837 threshold of 0.99 for the weight accumulation; lower thresholds result in more derived
838 haplogroup calls but with potentially higher uncertainty. **Fig. S3b.2** gives an example of this
839 approach for NEO962, a Mesolithic individual from Denmark with low coverage of 0.036X.

840



841
 842 **Fig. S3b.2. Phylogenetic placement.** Plot showing phylogenetic placement weights (A) and
 843 graft tree with most likely placement (B) for individual NEO962 on a subtree representing
 844 reference individuals with haplogroup I. (A) Weights for individual branches are indicated
 845 with edge colour and width, edges without placements are indicated with dashed line. While
 846 the majority of the placement weight mass is distributed among branches of haplogroup I2,
 847 non-zero weights are also found on branches ancestral to I1 (0.02) and I (0.07). The
 848 individual is hence conservatively assigned to haplogroup I. (B) The most likely branching of
 849 NEO962 was found within the subclade I2a1a1a, albeit this placement is associated with
 850 considerable uncertainty.

851

852 Sub-haplogroup analyses of newly reported samples

853

854 Haplogroup I

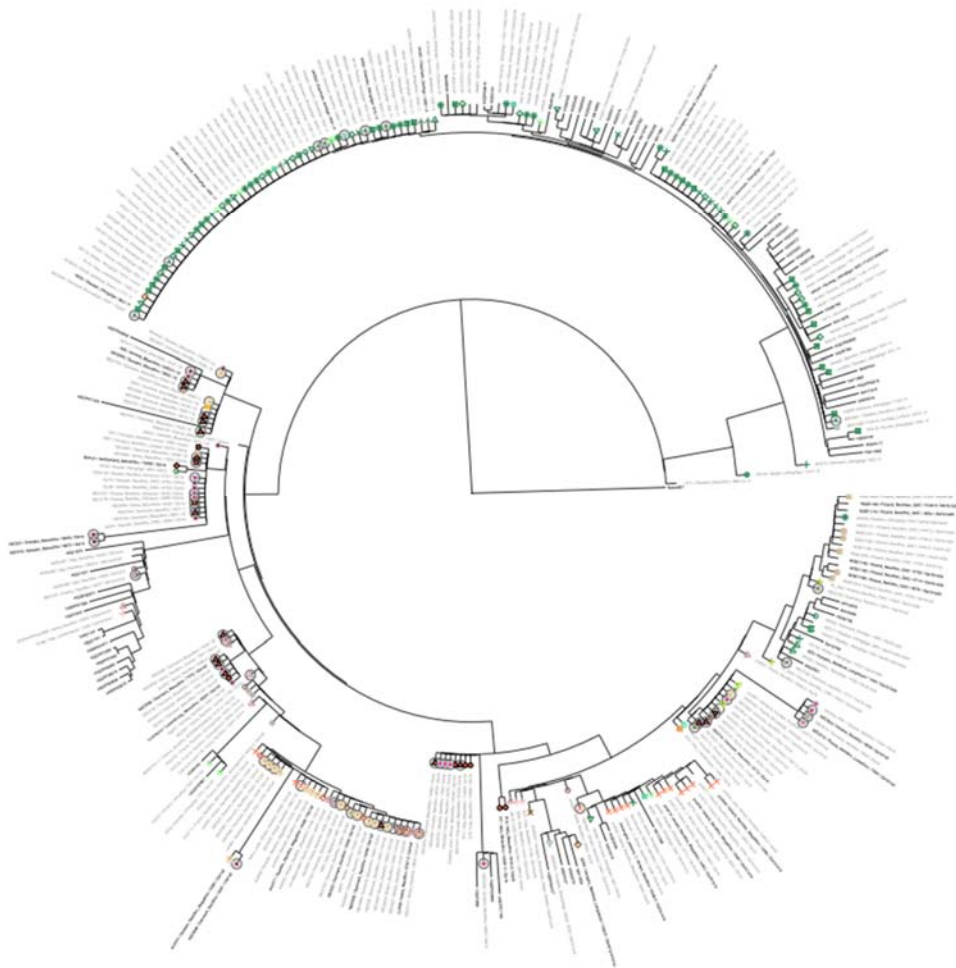
855

856 Haplogroup I2 was common among newly reported samples from western Eurasian hunter-
 857 gatherer contexts, as well as later Neolithic farmers. In particular, the 25 Danish Mesolithic
 858 male individuals were exclusively carriers of haplogroup I2, albeit with considerable diversity
 859 across different sub-haplogroups (Fig. S3b.3). Neolithic farmer individuals from Scandinavia
 860 were predominantly placed within an ancient-only subclade of haplogroup I2a1a2, containing
 861 other individuals from Neolithic farmer contexts across Europe.

862

863 The earliest presence of haplogroup I1, which is the most common haplogroup among
 864 present-day Scandinavians, was found ~4,000BP among late Neolithic and early Bronze Age
 865 Scandinavians newly reported in this study (Fig. S3b.3). A single Swedish Mesolithic

866 individual (sf11) was placed at the base of the I1 clade; however, its low coverage (0.1X)
867 precludes to conclude with certainty whether early I1-related lineages were indeed present
868 among Scandinavian hunter-gatherers.
869



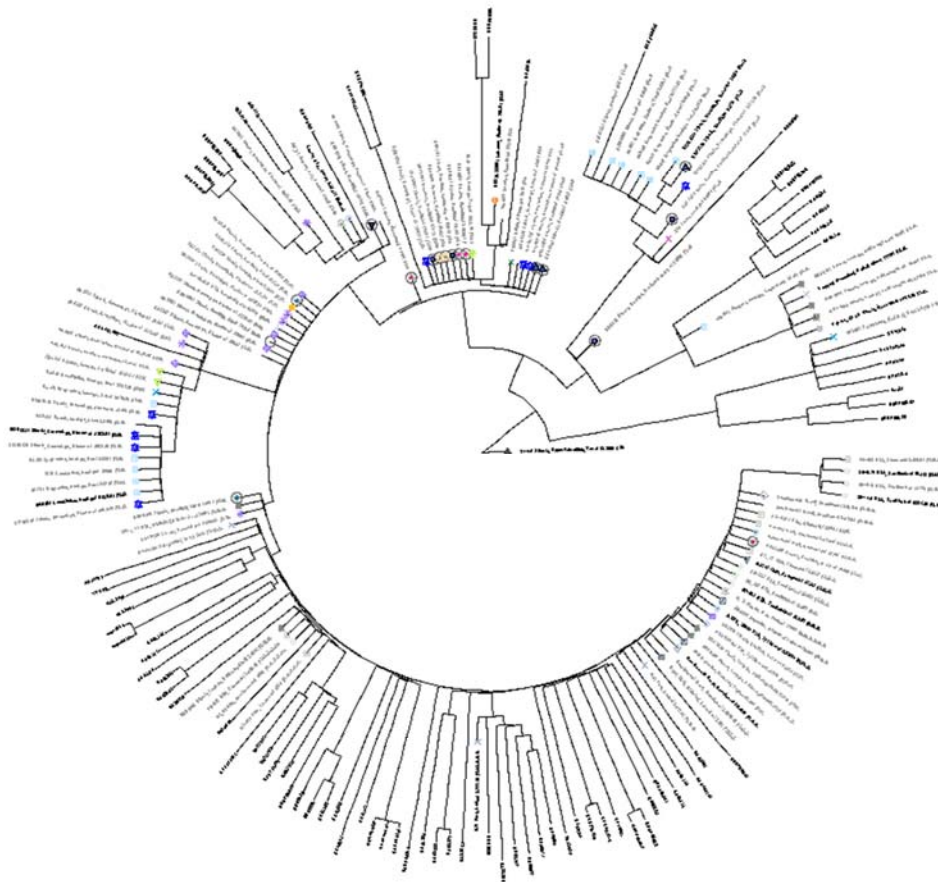
870
871 **Fig. S3b.3. Phylogeny of haplogroup I.** Phylogenetic tree with most likely placements of
872 ancient samples. Samples labelled with black colour were used to infer the reference tree,
873 whereas samples with grey labels were grafted from phylogenetic placement. Terminal
874 branches for ancient samples were shortened to aid visualisation. Symbol colours and
875 shapes indicate genetic clusters from IBD-based clustering. Newly reported individuals are
876 highlighted with circled symbols.

877

878 Haplogroup Q

879

880 Haplogroup Q1 was common among newly reported Neolithic hunter-gatherer individuals
881 from the Siberian Forest steppe and the Lake Baikal region (Fig. S3b. 4). We observed
882 haplogroup Q1b2, rare among ancient West Eurasians, in two Ukrainian hunter-gatherers
883 (NEO501, NEO516) as well as two Danish Neolithic farmer individuals (NEO599, NEO744).



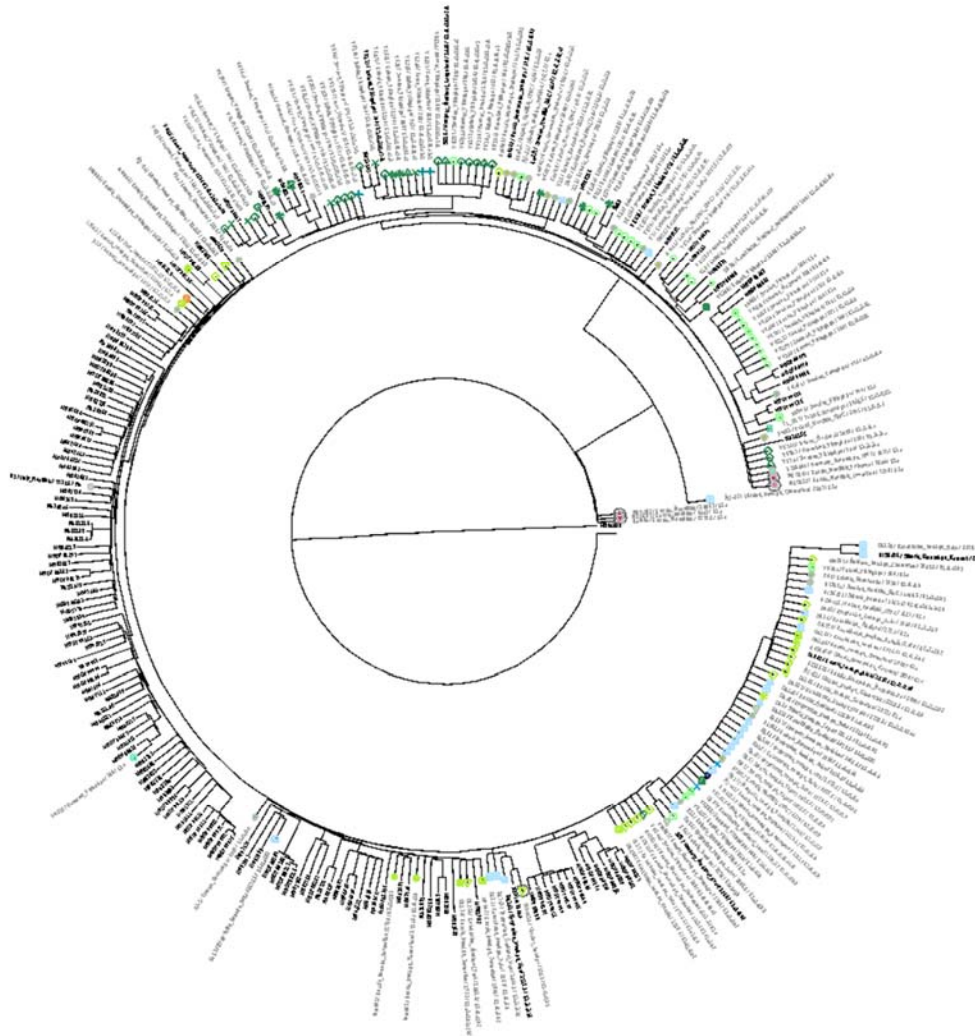
884
885 **Fig. S3b.4.** Phylogeny of haplogroup Q. Phylogenetic tree with most likely placements of
886 ancient samples. Samples labelled with black colour were used to infer the reference tree,
887 whereas samples with grey labels were grafted from phylogenetic placement. Terminal
888 branches for ancient samples were shortened to aid visualisation. Symbol colours and
889 shapes indicate genetic clusters from IBD-based clustering. Newly reported individuals are
890 highlighted with circled symbols.

891

892 Haplogroup R1a

893

894 Haplogroup R1a was found in the newly reported samples mainly among Eastern European
895 hunter-gatherer individuals. Phylogenetic placement suggests that the oldest individuals from
896 Mesolithic and Neolithic Russia represent early diverging lineages (Fig. S3b.5). Notably, a
897 ~7,300-year-old Neolithic individual from the Middle Don region (NEO113) was placed in a
898 basal R1a clade together with early individuals associated with the Corded Ware complex
899 (poz81, RISE446), which would make it the earliest observation of this lineage reported to
900 date.



901
902 **Fig. S3b.5. Phylogeny of haplogroup R1a.** Phylogenetic tree with most likely placements
903 of ancient samples. Samples labelled with black colour were used to infer the reference tree,
904 whereas samples with grey labels were grafted from phylogenetic placement. Terminal
905 branches for ancient samples were shortened to aid visualisation. Symbol colours and
906 shapes indicate genetic clusters from IBD-based clustering. Newly reported individuals are
907 highlighted with circled symbols.
908

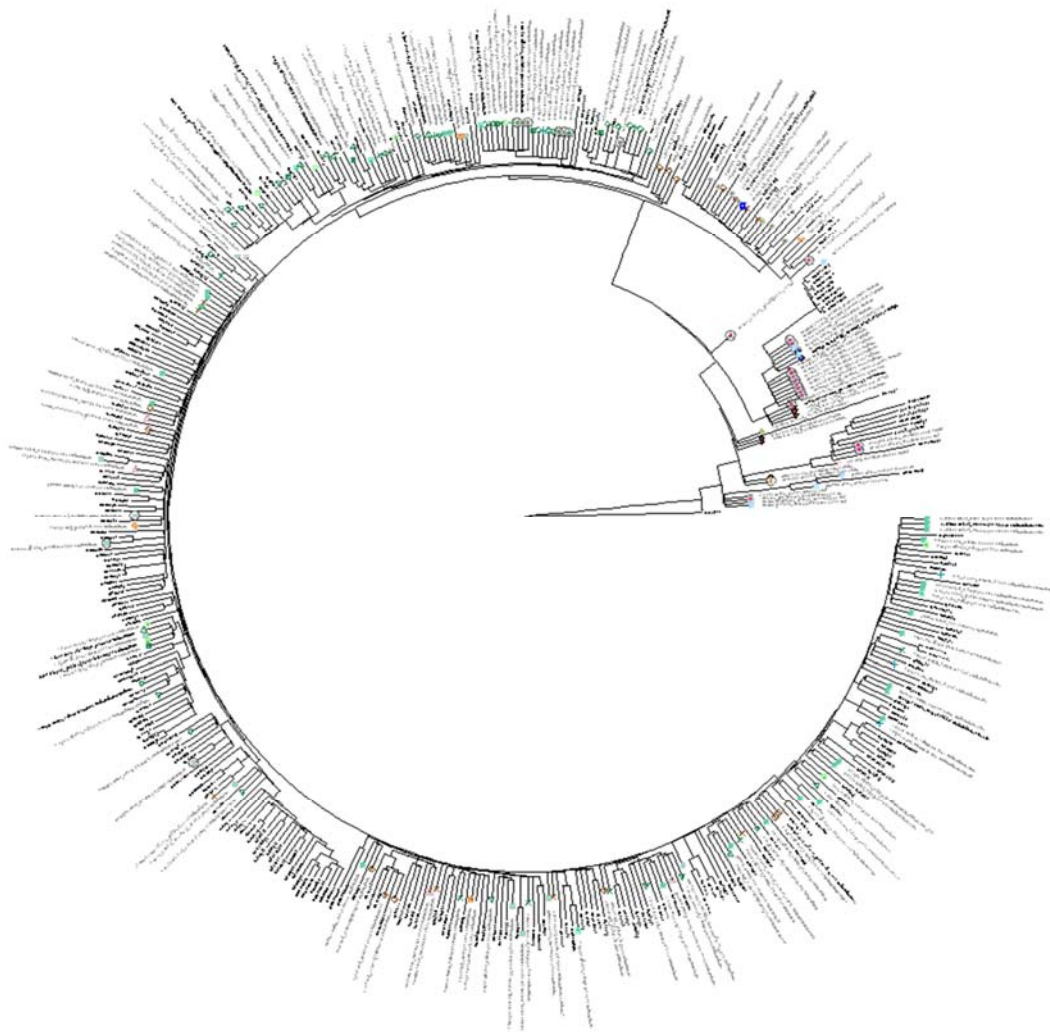
909 Haplogroup R1b

910

911 Newly reported samples belonging to haplogroup R1b were distributed between two distinct
912 groups depending on whether they formed part of the major European subclade R1b1a1b
913 (R1b-M269). Individuals placed outside this subclade were predominantly from Eastern
914 European Mesolithic and Neolithic contexts, and formed part of rare early diverging R1b
915 lineages (Fig. S3b.6). Two Ukrainian individuals belonged to a subclade of R1b1b (R1b-V88)
916 found among present-day Central and North Africans, lending further support^{5,10} to an
917 ancient Eastern European origin for this clade. Haplogroup R1b1a1a (R1b-M73) was
918 frequent among Russian Neolithic individuals.

919

920 Individuals placed within the R1b-M269 clade on the other hand were from Scandinavian
921 Late Neolithic and early Bronze Age contexts (Fig. S3b.6). Interestingly, more fine-scale sub-
922 haplogroup placements of those individuals revealed that Y chromosome lineages
923 distinguished samples from distinct genetic clusters inferred from autosomal IBD sharing
924 (Fig. S3b.6, S3b.7). In particular, individuals associated with the Scandinavian cluster
925 *Scandinavia_4200BP_3200BP* were all placed within the sub-haplogroup R1b1a1b1a1a1
926 (R1b-U106), whereas the two Scandinavian males associated with the Western European
927 cluster *Europe_4500BP_2000BP* were placed within R1b1a1b1a1a2 (R1b-P312) (Fig.
928 S3b.7).



929

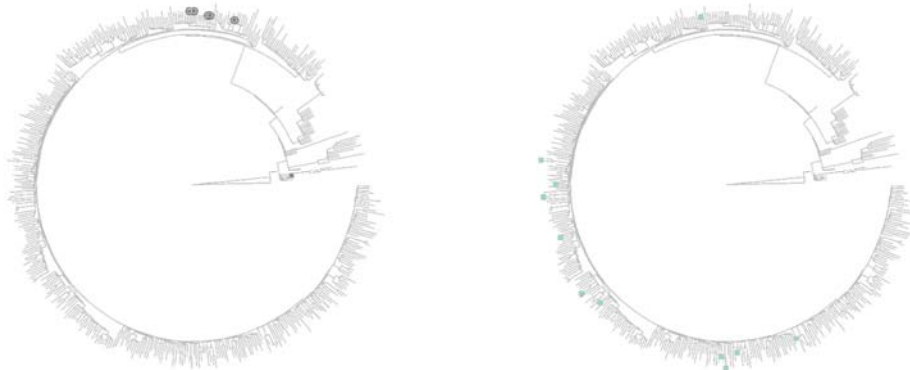
930 **Fig. S3b.6. Phylogeny of haplogroup R1b.** Phylogenetic tree with most likely placements
 931 of ancient samples. Samples labelled with black colour were used to infer the reference tree,
 932 whereas samples with grey labels were grafted from phylogenetic placement. Terminal
 933 branches for ancient samples were shortened to aid visualisation. Symbol colours and
 934 shapes indicate genetic clusters from IBD-based clustering. Newly reported individuals are
 935 highlighted with circled symbols.

936

937

Scandinavia_4200BP_3200BP

Europe_4500BP_2000BP



938

939 **Fig. S3b.7. Phylogeny of haplogroup R1b for genetic clusters.** Phylogenetic trees
940 showing most likely placements of ancient samples from Danish Late Neolithic and Bronze
941 Age genetic clusters Scandinavia_4200BP_3200BP (left) and Europe_4500BP_2000BP
942 (right). Terminal branches for ancient samples were shortened to aid visualisation. Symbol
943 colours and shapes indicate genetic clusters from IBD-based clustering. Newly reported
944 individuals are highlighted with circled symbols.

945 References

946

- 947 1. Skoglund, P., Storå, J., Götherström, A. & Jakobsson, M. Accurate sex identification
948 of ancient human remains using DNA shotgun sequencing. *J. Archaeol. Sci.* **40**, 4477–
949 4482 (2013).
- 950 2. Barbera, P. *et al.* EPA-ng: Massively Parallel Evolutionary Placement of Genetic
951 Sequences. *Syst. Biol.* **68**, 365–369 (2019).
- 952 3. Kozlov, A. M., Darriba, D., Flouri, T., Morel, B. & Stamatakis, A. RAXML-NG: a fast,
953 scalable and user-friendly tool for maximum likelihood phylogenetic inference.
954 *Bioinformatics* **35**, 4453–4455 (2019).
- 955 4. Czech, L., Barbera, P. & Stamatakis, A. Genesis and Gappa: processing, analyzing
956 and visualizing phylogenetic (placement) data. *Bioinformatics* **36**, 3263–3265 (2020).
- 957 5. Haber, M. *et al.* Chad Genetic Diversity Reveals an African History Marked by
958 Multiple Holocene Eurasian Migrations. *Am. J. Hum. Genet.* **0**, (2016).
- 959 6. Mallick, S. *et al.* The Simons Genome Diversity Project: 300 genomes from 142
960 diverse populations. *Nature* **538**, 201–206 (2016).
- 961 7. Bergström, A. *et al.* Insights into human genetic variation and population history from
962 929 diverse genomes. *Science* **367**, (2020).
- 963 8. Byrska-Bishop, M. *et al.* *High coverage whole genome sequencing of the expanded*
964 *1000 Genomes Project cohort including 602 trios.* 2021.02.06.430068
965 <https://www.biorxiv.org/content/10.1101/2021.02.06.430068v1> (2021)
966 doi:10.1101/2021.02.06.430068.
- 967 9. Ebenesersdóttir, S. S. *et al.* Ancient genomes from Iceland reveal the making of a
968 human population. *Science* **360**, 1028–1032 (2018).
- 969 10. Marcus, J. H. *et al.* Genetic history from the Middle Neolithic to present on the
970 Mediterranean island of Sardinia. *Nat. Commun.* **11**, 1–14 (2020).

971

972

3c) Relatedness

973

Martin Sikora¹

974

975 ¹Lundbeck Foundation GeoGenetics Centre, GLOBE Institute, University of Copenhagen,
976 Copenhagen, Denmark

977

978

979 Methods

980 To infer genetic relatedness between the study individuals we used the allele-frequency free
981 inference method introduced by Waples et al.¹. For each pair of individuals, we calculated
982 the three relatedness estimators R0, R1 and KING-robust² using the site-frequency-
983 spectrum (SFS)-based approach. We used the *realSFS*³ method implemented in the
984 *ANGSD*⁴ package to infer the 2D-SFS, selecting the SFS with the highest likelihood across
985 ten replicates. We used a set of 1,191,529 autosomal transversion SNPs with minor allele
986 frequency ≥ 0.05 from the 1000 Genomes Project⁵ for the analysis.

987

988 We used previously established cut-offs² for the KING-robust estimator to assign individual
989 pairs to first-, second- or third-degree relationships. Parent-offspring relationships were
990 distinguished from sibling relationships using R0 and R1 ratios, by requiring that $R0 \leq 0.02$
991 and $0.4 \leq R1 \leq 0.6$ to infer a parent-offspring relative pair. We excluded individual pairs with
992 less than 20,000 sites contributing to the estimators.

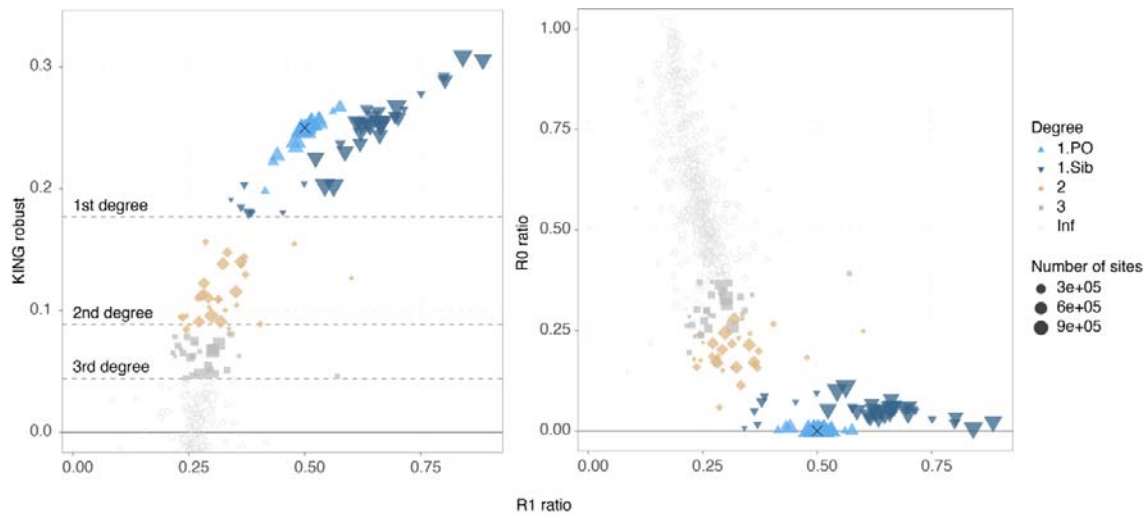
993

994 Results

995 We detected a total of 92 close relative pairs among the 1,664 dataset individuals, including
996 24 parent-offspring pairs, 36 siblings and 30 2nd degree pairs (Fig S3c.1, Supplementary
997 Table VI). We further found evidence of two duplicate / monozygotic twin relationships.
998 Sample NEO70 presented in this study was inferred to be from the same individual as
999 RISE554 previously reported in Allentoft *et al*⁶. Additionally, two male individuals MJ-15 and
1000 MJ-32 reported in Järve *et al*⁷ were also inferred as duplicate/twin pairs. In both cases
1001 genetic sex, mitochondrial as well as Y chromosome haplogroups were all consistent with
1002 their inferred relatedness status.

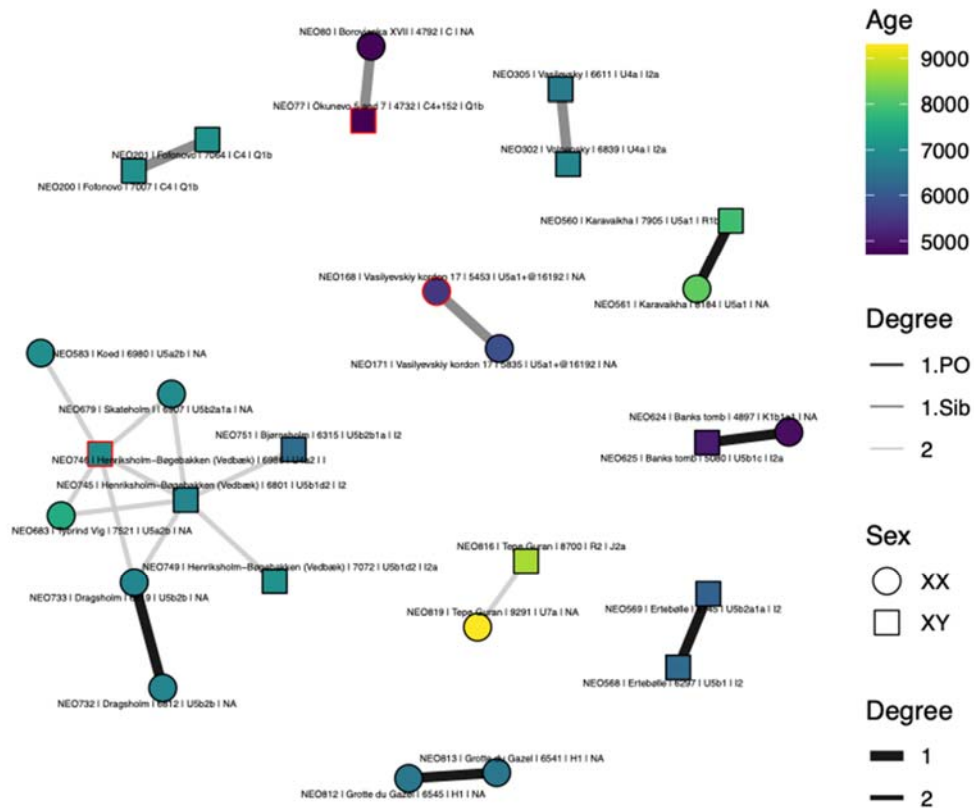
1003

1004



1005
 1006
 1007
 1008
 1009
 1010

Fig. S3c.1. Relatedness inference. Plots showing relatedness estimators R1/KING-robust (left) and R1/R2 (right) for pairs of individuals. Inferred relatedness status for each pair is indicated by plot symbol and colour, with symbol size scaled by the total number of informative sites. Black crosses indicate expected values for parent-offspring relationships.



1011

1012 **Fig. S3c.2. Relatedness among newly reported individuals.** Network showing first and
 1013 second-degree relationships, indicated by edge width and colour. Age of individuals is
 1014 indicated by fill colour, and individuals are labelled with site name, age, mitochondrial and Y
 1015 chromosome haplogroups. Individuals flagged for possible contamination are indicated in
 1016 red.

1017

1018 We identified a total of 10 first- and 12 second-degree relative pairs among the newly
 1019 reported individuals. However, inspection of the relatedness network revealed that the
 1020 majority of 2nd degree connections are between Danish Mesolithic individuals from different
 1021 sites and two individuals from Henriksholm (NEO745, NEO746), one of which was flagged
 1022 as contaminated (**Fig. S3c.2**). As excess heterozygosity due to contamination can lead to
 1023 artificially increased relatedness estimates, we excluded any pair involving those two
 1024 individuals, as well as two other pairs involving contaminated individuals from the final list of
 1025 close relatives (**Table S3c.1**). Finally, three individuals reported here were inferred to be
 1026 either the same individual (2) or close relatives (1) of samples previously published using
 1027 targeted SNP capture (**Table S3c.2**).

1028

1029

Individual 1	Individual 2	Site	Country	Degree	Notes
NEO568	NEO569	Ertebølle	Denmark	1.PO	NEO568 ("Ertebølle man") is the father of infant NEO569
NEO732	NEO733	Dragsholm	Denmark	1.PO	Mother-daughter relationship, direction unknown
NEO624	NEO625	Banks tomb	UK	1.PO	NEO625 is the father of juvenile NEO624
NEO813	NEO812	Grotte du Gazel	France	1.PO	NEO813 is the mother of infant NEO812
NEO560	NEO561	Karavaikha	Russia	1.PO	NEO561 likely the mother of NEO560 (age and MT haplogroup)
NEO201	NEO200	Fofonovo	Russia	1.Sib	
NEO302	NEO305	Volnensky / Vasilevsky	Ukraine	1.Sib	
NEO816	NEO819	Tepe Guran	Iran	2	

1030

1031 **Table S3c.1. Close relatives among newly reported individuals.**

1032

1033

Study Individual	Related individual	Publication	Relationship
------------------	--------------------	-------------	--------------

NEO73	I1960	Narasimhan et al 2019 Science	same individual
NEO669	I5407	Mathieson et al 2018 Nature	same individual
NEO60	BOO005	Lamnidis et al 2018 Nature Communications	first degree, infant NEO60 likely daughter of BOO005

1034

1035 **Table S3c.2. Study individuals with related published individuals.**

1036

1037 References

1038

1039 1. Waples, R. K., Albrechtsen, A. & Moltke, I. Allele frequency-free inference of close
1040 familial relationships from genotypes or low-depth sequencing data. *Mol. Ecol.* **28**, 35–48
1041 (2019).

1042 2. Manichaikul, A. *et al.* Robust relationship inference in genome-wide association
1043 studies. *Bioinformatics* **26**, 2867–2873 (2010).

1044 3. Nielsen, R., Korneliussen, T., Albrechtsen, A., Li, Y. & Wang, J. SNP Calling,
1045 Genotype Calling, and Sample Allele Frequency Estimation from New-Generation
1046 Sequencing Data. *PLoS ONE* **7**, e37558 (2012).

1047 4. Korneliussen, T. S., Albrechtsen, A. & Nielsen, R. ANGSD: Analysis of Next
1048 Generation Sequencing Data. *BMC Bioinformatics* **15**, 356 (2014).

1049 5. The 1000 Genomes Project Consortium,. A global reference for human genetic
1050 variation. *Nature* **526**, 68–74 (2015).

1051 6. Allentoft, M. E. *et al.* Population genomics of Bronze Age Eurasia. *Nature* **522**, 167–+
1052 (2015).

1053 7. Järve, M. *et al.* Shifts in the Genetic Landscape of the Western Eurasian Steppe
1054 Associated with the Beginning and End of the Scythian Dominance. *Curr. Biol.* **29**, 2430-
1055 2441.e10 (2019).

1056

1057 3d) Pop structure general, PCA/Admixture (Martin)

1058

Martin Sikora¹

1059

1060 ¹Lundbeck Foundation GeoGenetics Centre, GLOBE Institute, University of Copenhagen,
1061 Copenhagen, Denmark

1062

1063 Methods

1064

We generated a dataset for population genetic analysis by combining the 317 newly

1065

sequenced individuals with 1,347 previously published ancient genomes with genomic

1066 coverage >0.1X generated using shotgun sequencing (Supplementary Table VII). Imputed
1067 genotype data (Supplementary note S2) for this set of 1,664 ancient genomes was merged
1068 with genotypes of 2,504 modern individuals from the 1,000 Genomes project¹ used as a
1069 reference panel in the imputation. We retained only SNPs passing the 1000 Genomes strict
1070 mask, resulting in a final dataset of 4,168 individuals genotyped at 7,321,965 autosomal
1071 SNPs (“1000G” dataset). In addition to imputed genotypes, we also generated pseudo-
1072 haploid genotypes for each ancient individual by randomly sampling an allele from
1073 sequencing reads covering those SNPs. For population structure analyses in the context of
1074 global genetic diversity, we generated a second dataset by intersecting the ancient genotype
1075 data with SNP array data of 2,180 modern individuals from 213 world-wide populations²⁻⁵
1076 (“HO” dataset).

1077

1078 To facilitate filtering for downstream analyses, we flagged individuals to potentially exclude
1079 based on the following criteria:

1080

- 1081 - Contamination estimate >5% (“*contMT5pct*”, “*contNuc5pct*”, Supplementary note S1)
- 1082 - Autosomal coverage < 0.1X (“*lowcov*”)
- 1083 - Genome-wide average imputation genotype probability < 0.98 (“*lowGpAvg*”)
- 1084 - Individual is the lower quality sample in a close relative pair (“*1d_rel*”, “*2d_rel*”;
1085 Supplementary note S3c)

1086

1087 A total of 1,492 individuals (213 newly reported) passed all filters, which were used in the
1088 majority of downstream analyses unless otherwise noted.

1089

1090 We investigated overall population structure among the dataset individuals using principal
1091 component analyses (PCA) and model-based clustering (ADMIXTURE⁶). We carried out
1092 PCA using different subsets of individuals in the “HO” dataset. For the PCA including only
1093 imputed diploid samples, we used GCTA⁷, excluding SNPs with minor allele frequency
1094 (MAF) < 0.05 in the respective panel. For PCA projecting low coverage or flagged
1095 individuals, we used *smartpca*⁸ with options ‘*lsqproject: YES*’ and ‘*autoshrink: YES*’ on a
1096 fixed set of 400,186 SNPs with MAF ≥ 0.05 in non-African individuals passing all filters.

1097

1098 We ran *ADMIXTURE* on a set of 1,593 ancient individuals from the “1000G” dataset,
1099 excluding individuals flagged as close relatives or a contamination estimate >5%. For the
1100 1,492 individuals passing all filters we used imputed genotypes, the remaining 101 lower
1101 coverage samples were represented by pseudo-haploid genotypes. We restricted the
1102 analysis to transversion SNPs with imputation INFO score ≥ 0.8 and MAF ≥ 0.05. We further

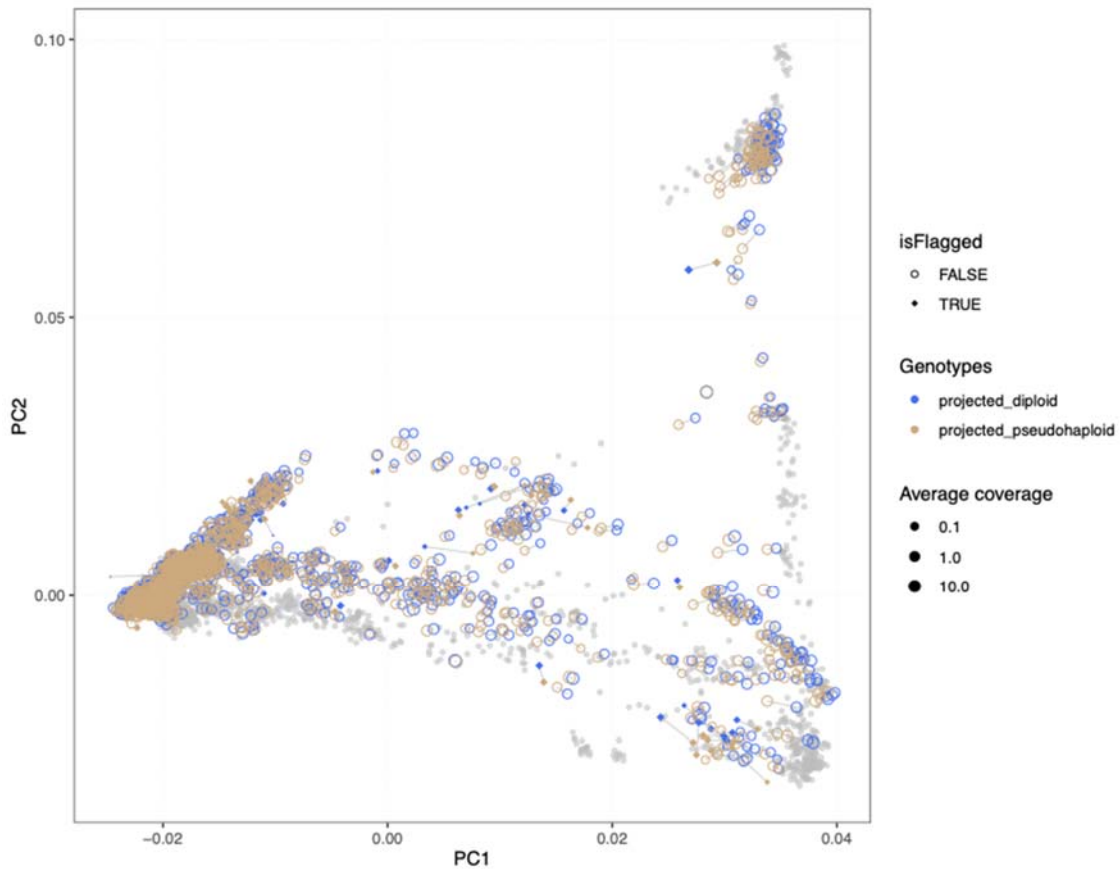
1103 performed linkage disequilibrium (LD) pruning and filtering for missingness using *plink*⁹
1104 (options '*--indep-pairwise 500 50 0.4 --geno 0.8*'), for a final analysis set of 142,550 SNPs.
1105
1106
1107

1108 Results

1109 Comparison of imputed and pseudo-haploid genotypes in PCA space

1110 To determine the consistency of imputed and pseudo-haploid genotypes when used in PCA,
1111 we followed the approach of Antonio *et al.*¹⁰ comparing the coordinates of both sets of
1112 genotypes for each individual when projected onto principal components inferred from
1113 modern individuals (Fig. S3d.1, S3d.3). We did not use the “autoshrink” option of *smartpca*
1114 for this analysis to avoid possible systematic differences in the projection correction between
1115 the two sets of genotypes. Projected PCA positions for samples passing all filters were
1116 consistent between imputed and pseudo-haploid genotypes, with no evidence for systematic
1117 shifts (Fig. S3d.2, S3d.4) and only a very subtle relationship of PCA distance between
1118 genotypes with genomic coverage (Fig. S3d.3, S3d.6). More substantial shifts were only
1119 observed with low coverage (<0.1X) flagged samples.

1120
1121
1122

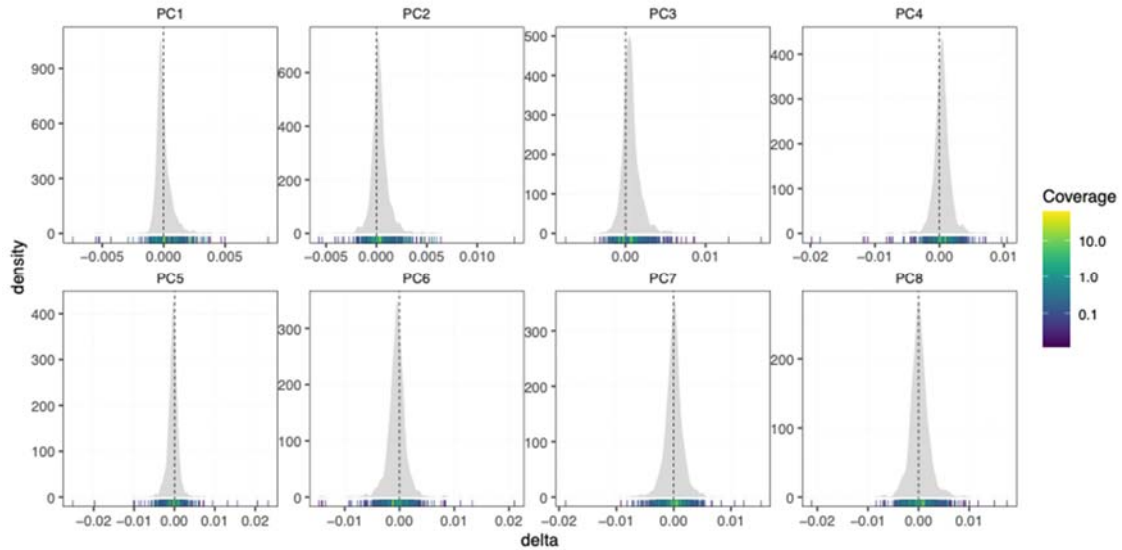


1123

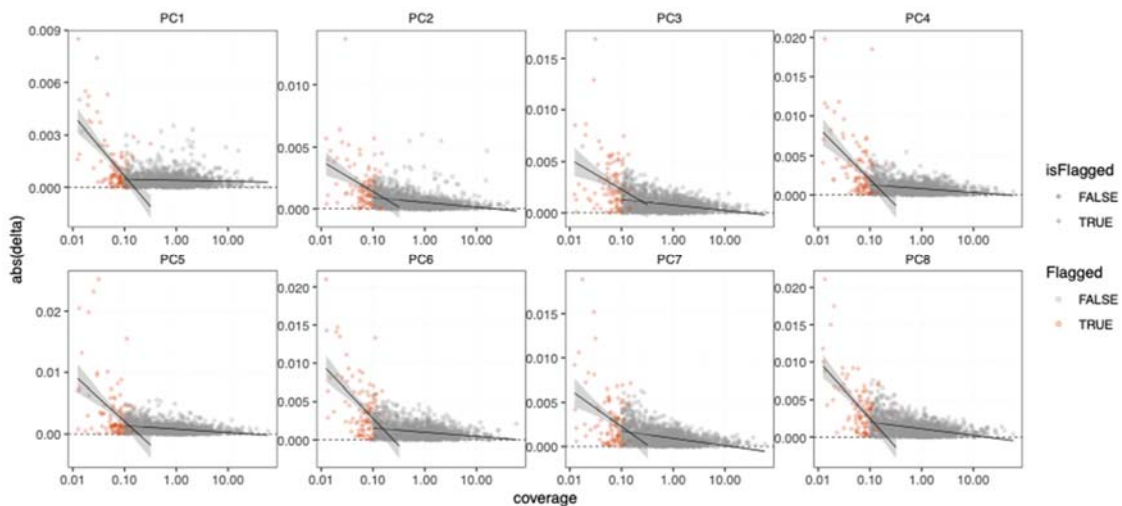
1124 **Fig. S3d.1. PCA projection of imputed and pseudo-haploid genotypes.** Colored symbols

1125 show the position of imputed diploid (blue) and pseudo-haploid (beige) genotypes for each

1126 ancient individual, projected onto principal components inferred from modern individuals from
1127 Eurasia, Oceania, and the Americas. Genotype pairs from the same individual are connected
1128 by grey lines.
1129

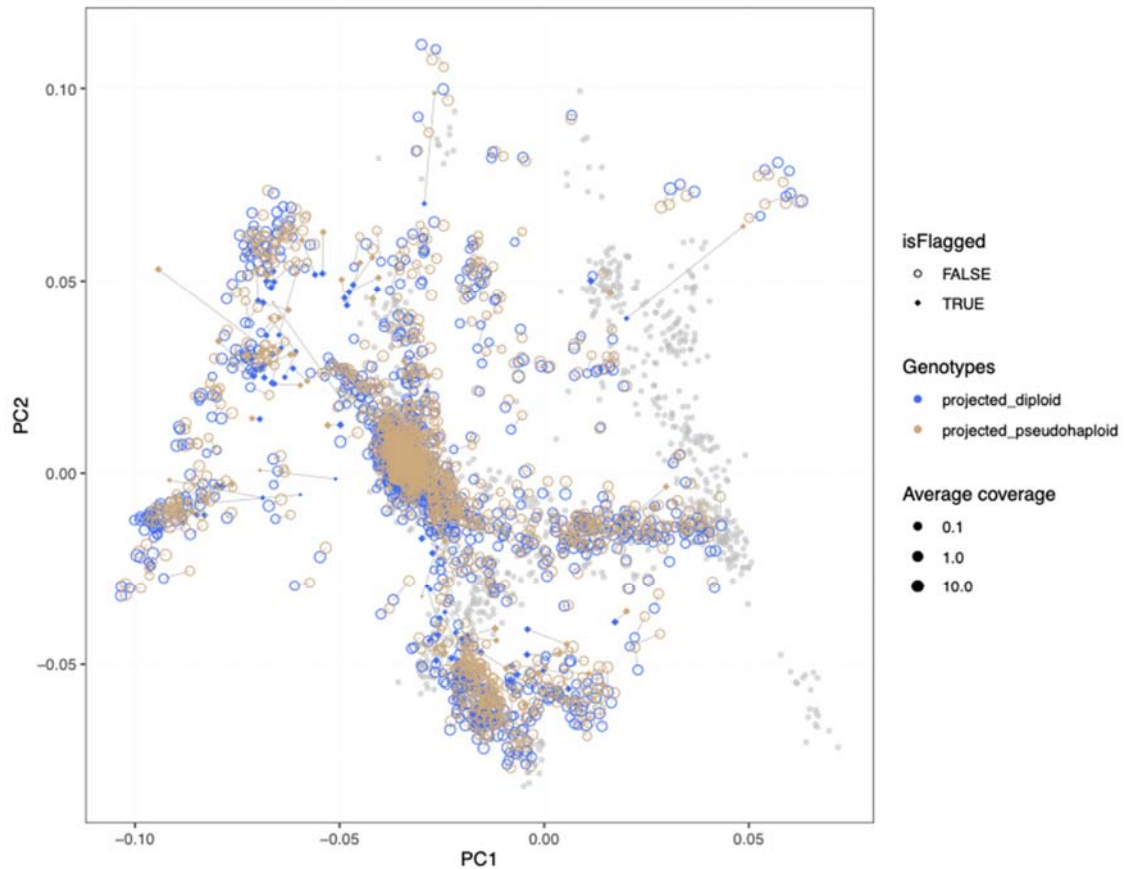


1130
1131 **Fig. S3d.2.** Distribution of differences between genotypes in PC space. Density plots
1132 show differences along individual PCs between imputed and pseudo-haploid genotypes for
1133 each individual in Fig. S3d.1, as a function of their average read depth. Marginal rug plots
1134 show individual observations, colored by the average read depth of the respective individual.
1135

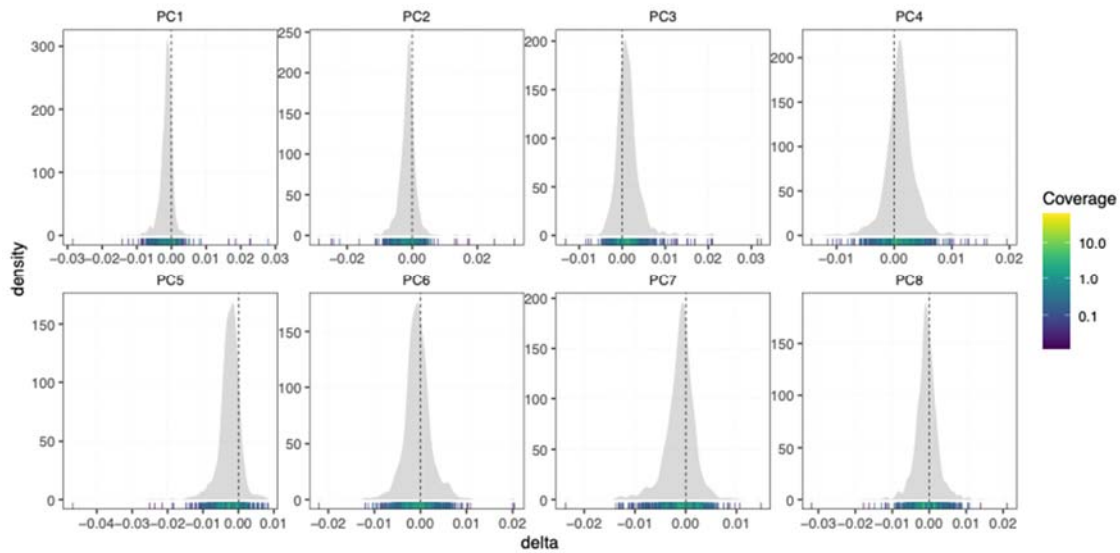


1136
1137 **Fig. S3d.3.** Relationship of read depth and PCA position. Plot shows absolute value of
1138 differences along individual PCs between imputed and pseudo-haploid genotypes for each
1139 individual in Fig. S3d.1, as a function of their average read depth. Individuals flagged for low

1140 coverage or low GP average are indicated with red symbols. Linear regression lines for
1141 flagged and unflagged individuals are shown with black lines.



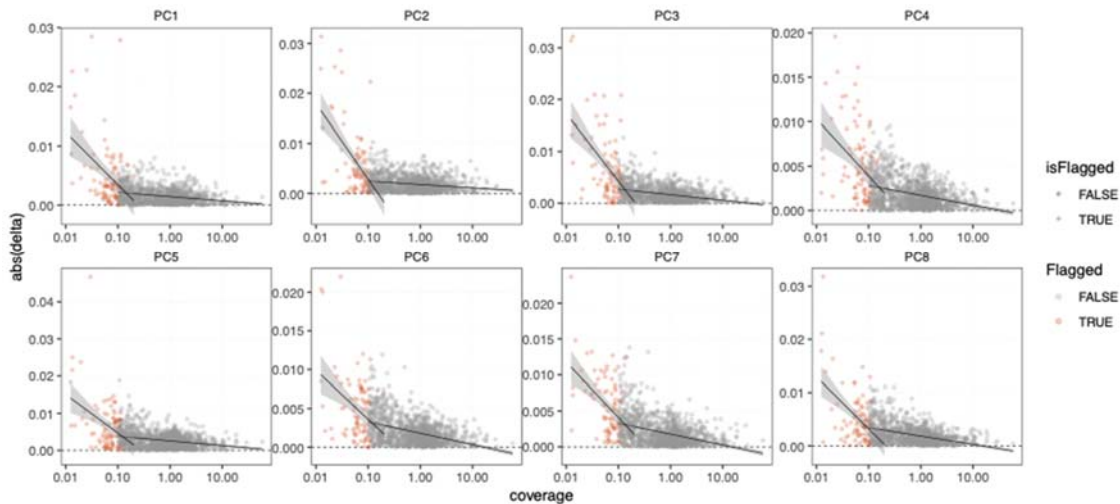
1142 **Fig. S3d.4. PCA projection of imputed and pseudo-haploid genotypes.** Colored symbols
1143 show the position of imputed diploid (blue) and pseudo-haploid (beige) genotypes for each
1144 ancient individual, projected onto principal components inferred from modern individuals from
1145 Western Eurasia. Genotype pairs from the same individual are connected by grey lines.
1146
1147



1148

1149 **Fig. S3d.5. Distribution of differences between genotypes in PC space.** Density plots
 1150 show differences along individual PCs between imputed and pseudo-haploid genotypes for
 1151 each individual in Fig. S3d.4, as a function of their average read depth. Marginal rug plots
 1152 show individual observations, coloured by average read depth of the respective individual.

1153



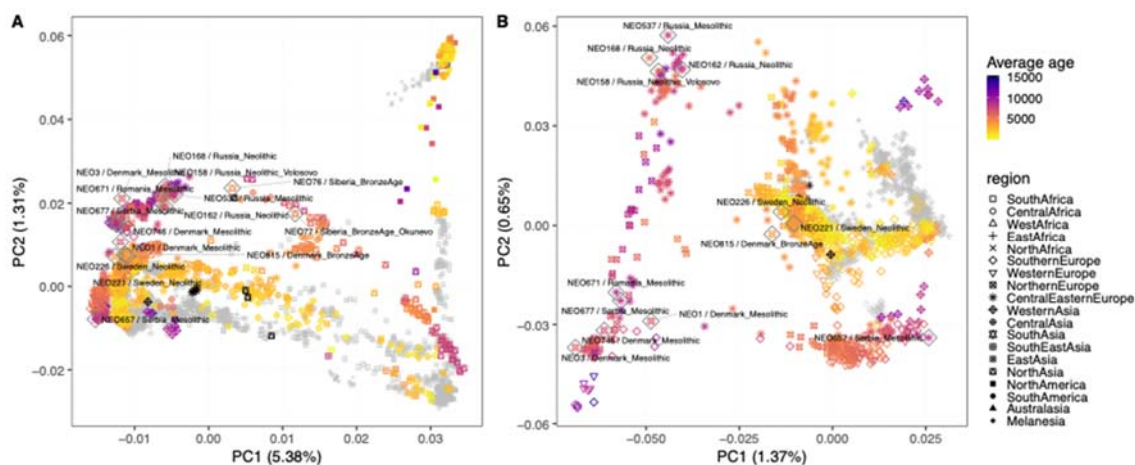
1154

1155 **Fig. S3d.6. Relationship of read depth and PCA position.** Plot shows absolute difference
 1156 along individual PCs between imputed and pseudo-haploid genotypes for each individual in
 1157 Fig. S3d.4, as a function of their average read depth. Individuals flagged for low coverage or
 1158 low GP average are indicated with red symbols. Linear regression lines for flagged and
 1159 unflagged individuals are shown with black lines.

1160

1161 PCA position of samples flagged as contaminated

1162 To investigate the effect of elevated contamination estimates on the position of individuals
1163 flagged as possibly contaminated, we projected them onto the principal components inferred
1164 from modern and ancient individuals passing all filters. We found that the majority of those
1165 individuals projected consistently with ancient samples of related age and regional contexts
1166 (Fig. S3d.7). An exception to this is seen in the Mesolithic Danish individual NEO1, which
1167 shows a clear shift towards present-day Europeans along PC1 and PC2. Overall, our results
1168 suggest that inferences about broad patterns of deep Eurasian population structure are likely
1169 not affected in the majority of the flagged individuals. We nevertheless opted for a
1170 conservative approach and excluded those individuals from in-depth analyses further
1171 downstream.
1172



1173 **Fig. S3d.7. PCA positions of individuals flagged as contaminated.** Flagged individuals
1174 are labelled and outlined with grey diamonds. Principal components were inferred using
1175 ancient and modern individuals from (A) Eurasia, Oceania, and the Americas or (B) Western
1176 Eurasia. Plot symbols indicate geographic region, coloured by age of the respective
1177 individual. Present-day individuals are indicated in grey.

1179

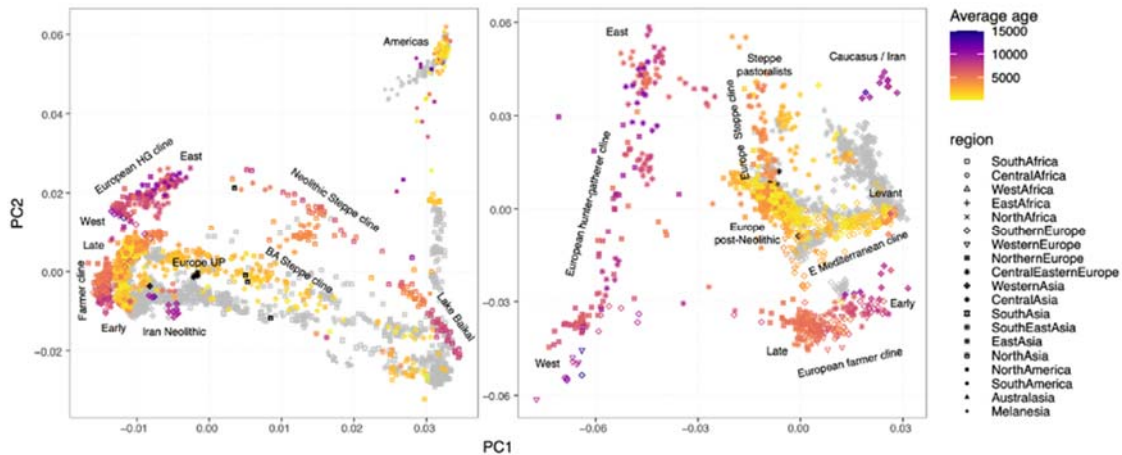
1180

1181 Genetic ancestry of newly reported samples

1182 Overview

1183

1184



1185

1186 **Fig. S3d.8. Overview of genetic structure.** PCA of ancient and modern individuals from
1187 Eurasia, Oceania, and the Americas (left), or Western Eurasia (right). Plot symbols indicate
1188 geographic regions, colored by the age of the respective individual. Present-day individuals
1189 are indicated in grey. Terms for spatiotemporal ancestry clusters and clines are indicated.

1190 Genetic structure in a PCA using 3,316 individuals from regions outside Africa is dominated
1191 by continental-scale differentiation among western Eurasia (defined here as west of the
1192 Urals), east Asia and the Americas (Fig S3d.8). Two west-east clines of ancient individuals
1193 connect western and eastern Eurasia: A “Neolithic Steppe cline” between hunter-gatherers of
1194 the West Siberian Forest Steppe and Lake Baikal; as well as a later “BA Steppe cline” linking
1195 Western Steppe pastoralists with the Altai mountain region (Fig S3d.8).

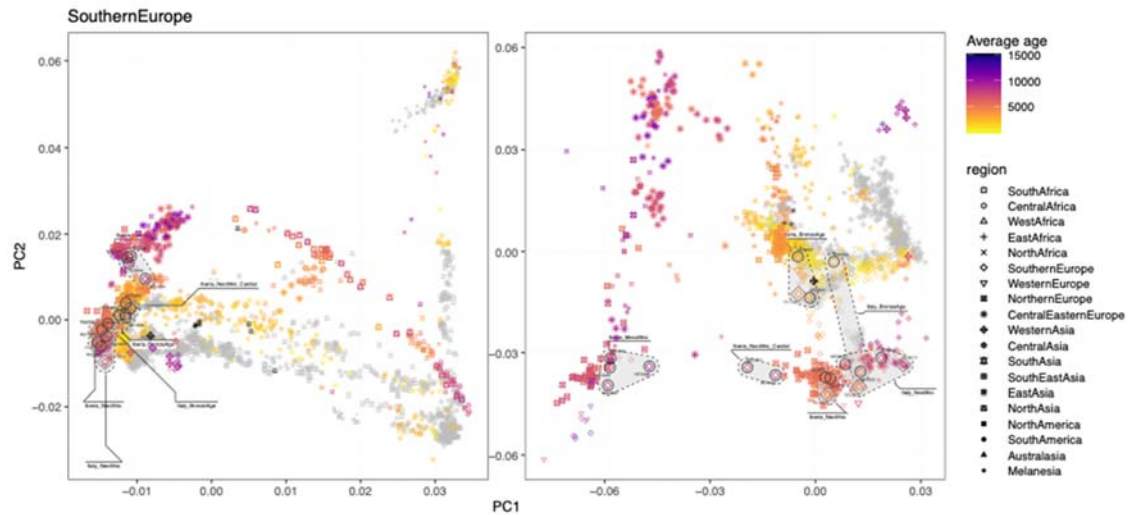
1196 Focusing the PCA on 2,126 modern and ancient individuals from Western Eurasia, the
1197 extremes of the PCA space are defined by clines and clusters related to previously described
1198 “deep” ancestry sources, including: A “European hunter-gatherer cline” between western and
1199 eastern European Mesolithic individuals; A “European farmer cline” ranging from early
1200 Neolithic individuals from Anatolia and Southern Europe to mid- and late Neolithic European
1201 individuals; and hunter-gatherers and early farmers from Iran and the Caucasus. European
1202 individuals from the late Neolithic and early Bronze Age onwards form an extended
1203 “European post-Neolithic” cluster in the centre of the PCA, differentiated along either a
1204 “European Steppe cline” between Steppe pastoralists and late European farmers, or an
1205 “Eastern Mediterranean cline” anchored in the east by Anatolian and Levantine Bronze Age
1206 individuals (Fig S3d.8). The newly reported genomes from western Eurasia cluster across

1207 the entire range of the PCA, resulting in increased fine-scale resolution along the major
1208 ancestry clines, particularly the European hunter-gatherer and farmer clines. The following
1209 sections provide regional descriptions for the patterns of ancestry observed in the newly
1210 reported samples.

1211

1212 Southern Europe

1213



1214

1215 **Fig. S3d.9. Newly reported individuals from Southern Europe.** PCA positions of newly
1216 reported individuals are highlighted with black circles (imputed) or grey diamonds (pseudo-
1217 haploid, projected). Individuals from the same spatiotemporal group are connected with
1218 shaded hulls.

1219

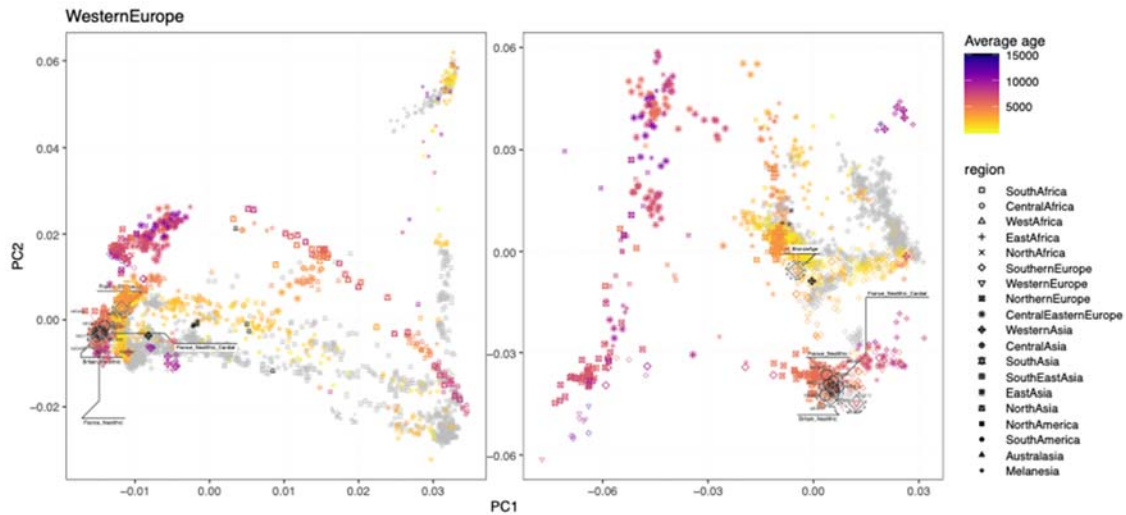
1220 We report 18 new individuals from Italy (6) and the Iberian Peninsula (12), distributed across
1221 European hunter-gatherer (HG), farmer, and post-Neolithic ancestry clusters (Fig S3d.9).

1222 Four Iberian Mesolithic individuals cluster with other Southern European Mesolithic
1223 individuals at the “western” end of the European HG cline. Among the individuals falling
1224 within the European farmer cline, two early Neolithic individuals from Portugal (NEO631,
1225 NOE632; ~ 7,300 BP) are shifted towards the European HG cline suggestive of increased
1226 HG ancestry. The four most recent individuals (from ~4,100 BP) form part of the extended
1227 European post-Neolithic cluster.

1228

1229 Western Europe

1230



1231

1232 **Fig. S3d.10. Newly reported individuals from Western Europe.** PCA positions of newly
 1233 reported individuals are highlighted with black circles (imputed) or grey diamonds (pseudo-
 1234 haploid, projected). Individuals from the same spatiotemporal group are connected with
 1235 shaded hulls.

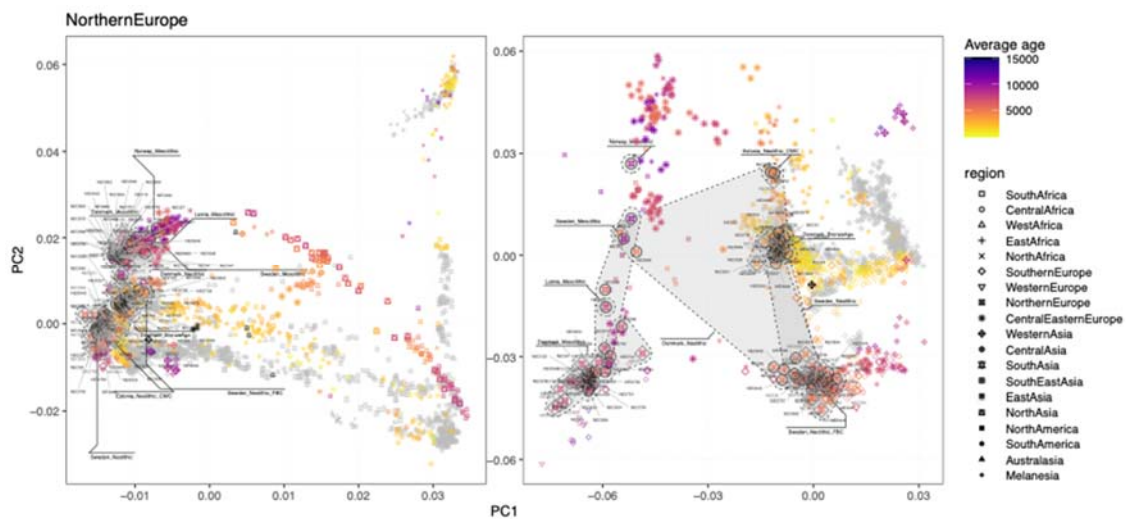
1236

1237 We report 12 new individuals from France (5) and the UK (7), from the early Neolithic to the
 1238 Bronze Age. All 11 Neolithic individuals fall within the European farmer cline, whereas a
 1239 single Bronze Age individual from Grotte Mandrin (NEO120, ~3,400BP) clustered with post-
 1240 Neolithic Europeans (Fig. S3d.10).

1241

1242 Northern Europe

1243



1244

1245 **Fig. S3d.11. Newly reported individuals from Northern Europe.** PCA positions of newly
 1246 reported individuals are highlighted with black circles (imputed) or grey diamonds (pseudo-

1247 haploid, projected). Individuals from the same spatiotemporal group are connected with
1248 shaded hulls.

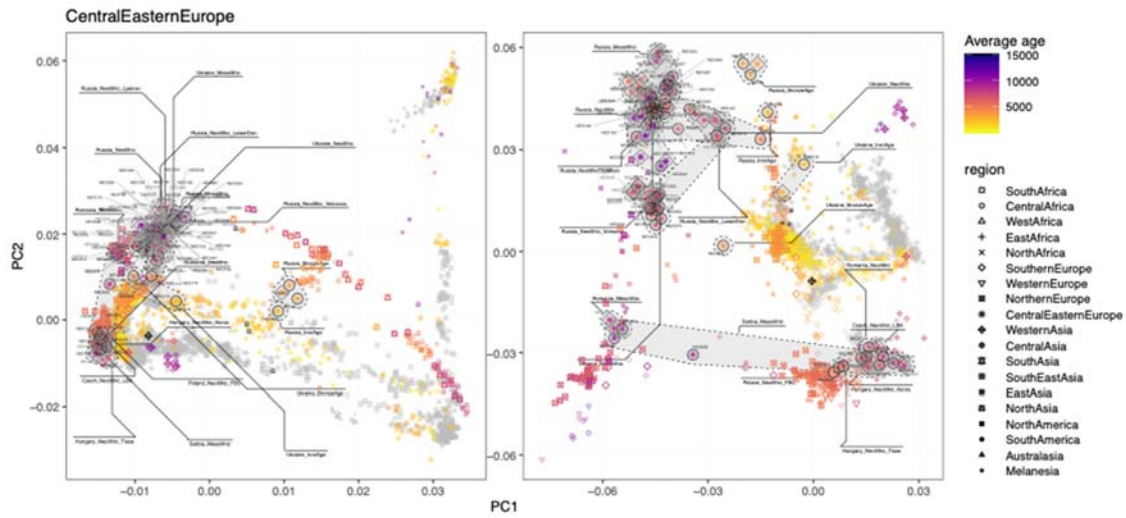
1249

1250 We report 124 new individuals from Denmark (100), Sweden (21), Norway (1) and the Baltic
1251 (2), spanning a period from ~10,500BP to 3,100BP. This transect includes 46 Mesolithic
1252 individuals, all of which cluster within the European HG cline (Fig. S3d.10). The 40 HG
1253 individuals from Denmark fall towards the “western” end of the cline, whereas the other
1254 Scandinavian and Baltic individuals occupy varied positions shifted towards the “eastern”
1255 end of the cline. Neolithic Scandinavian individuals generally fall towards the “late” end of the
1256 European farmer cline, with later Neolithic individuals also found among the extended post-
1257 Neolithic Europe cluster. Three late Neolithic (~4,500BP) individuals from Denmark
1258 (NEO876, NEO792) and Estonia (NEO306) are shifted further up along PC2 towards the
1259 Steppe pastoralist cluster, suggesting higher amounts of Steppe-related ancestry (Fig.
1260 S3d.11).

1261

1262 Central and Eastern Europe

1263



1264

1265 **Fig. S3d.12. Newly reported individuals from Southern Europe.** PCA positions of newly
1266 reported individuals are highlighted with black circles (imputed) or grey diamonds (pseudo-
1267 haploid, projected). Individuals from the same spatiotemporal group are connected with
1268 shaded hulls.

1269

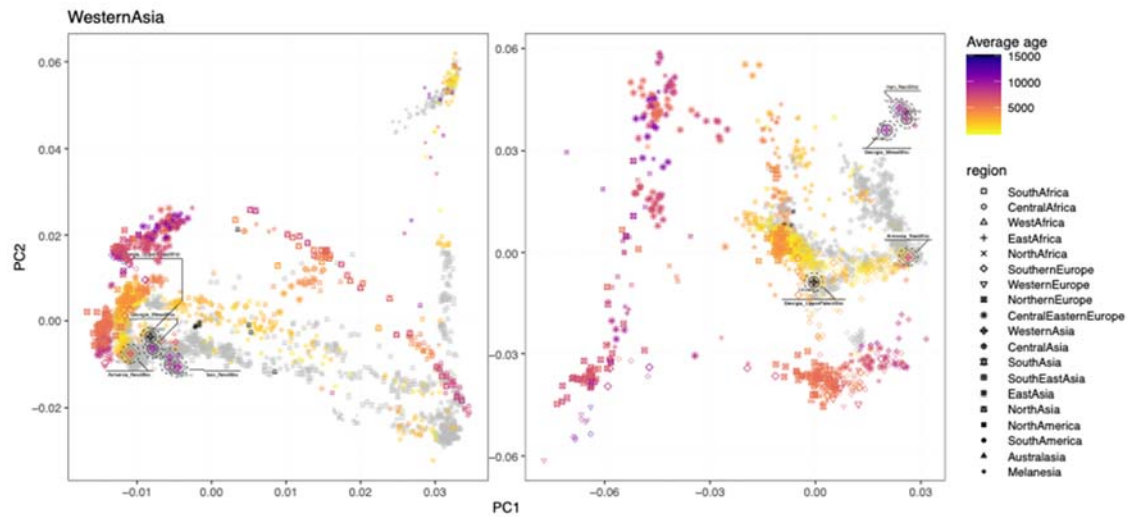
1270 We report 112 new individuals from Central and Eastern Europe, falling into two distinct
1271 groups. The 92 individuals from Russia (57) and Ukraine (35) predominantly occupy a broad
1272 area between the centre and “eastern” end of the European HG cline, roughly corresponding
1273 to a geographic cline from the south (Ukraine) to the north (Russia) (Fig. S3d.11). Among the
1274 southern Russian samples, six individuals from Golubaya Krinitsa in the Middle Don region
1275 are shifted on a cline along PC1 towards Iranian and Caucasus Mesolithic and Neolithic at
1276 the other extreme, falling close to later Steppe pastoralists from the region. Three Bronze
1277 Age individuals from Northwestern Russia (Bol'shoy Oleni Ostrov; NEO60, NEO61, NEO62)
1278 are positioned between the Neolithic and BA Steppe clines, centrally between the West and
1279 East Eurasian poles in the extended PCA of all non-Africans (Fig. S3d.12).

1280

1281 The remaining 20 samples from Central and Southeastern Europe include Mesolithic
1282 individuals at the “western” end of the European HG cline, as well as early Neolithic
1283 individuals on the farmer cline. An early Neolithic individual from Iron Gates, Serbia
1284 (NEO658) is found intermediate between the HG and farmer clines, suggestive of recent
1285 farmer/HG admixture (Fig. S3d.12).

1286

1287 Western Asia



1289

1290

Fig. S3d.13. Newly reported individuals from Western Asia. PCA positions of newly

1291

reported individuals are highlighted with black circles (imputed) or grey diamonds (pseudo-

1292

haploid, projected). Individuals from the same spatiotemporal group are connected with

1293

shaded hulls.

1294

1295

We report 6 new individuals from Iran (3) and the South Caucasus region (3). The oldest

1296

sample in the dataset, a ~25,000-year-old individual from Georgia is positioned intermediate

1297

between Upper Paleolithic Europeans and early Neolithic farmers in both the west Eurasian

1298

and extended non-African PCA (**Fig. S3d.13**).

1299

Three Iranian Neolithic individuals (~9,200 BP) as well as one Mesolithic Georgian individual (NEO281; ~9,700 BP) fall with other previously

1300

published samples of similar provenance, defining one of the extremes of PC1/PC2 space.

1301

One Neolithic individual from Armenia (NEO110, ~7,600 BP) is found at the “eastern”

1302

extreme of an eastern Mediterranean cline between ancient Levantine individuals and

1303

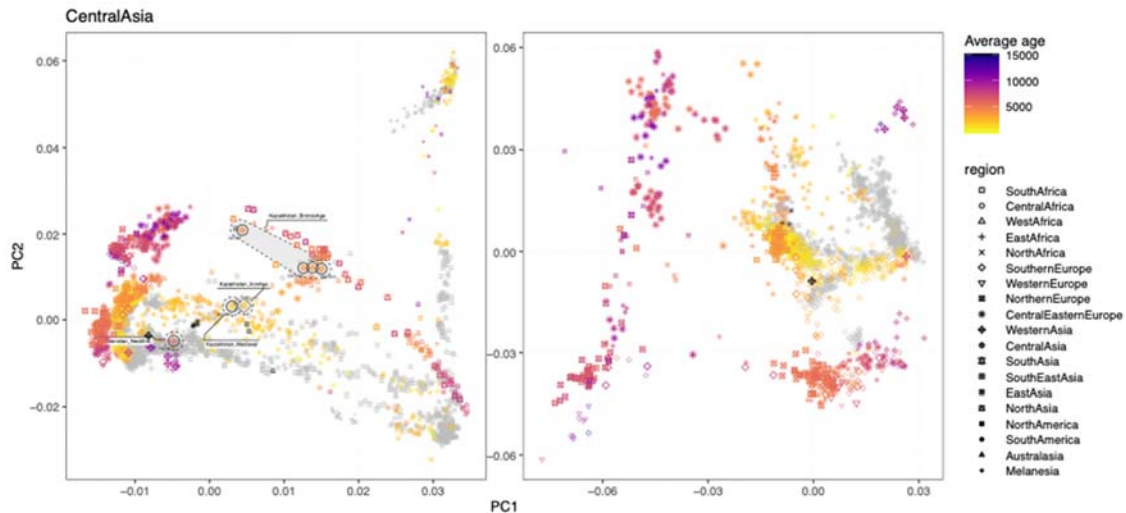
Southern post-Neolithic Europeans.

1304

1305

Central Asia

1306



1307

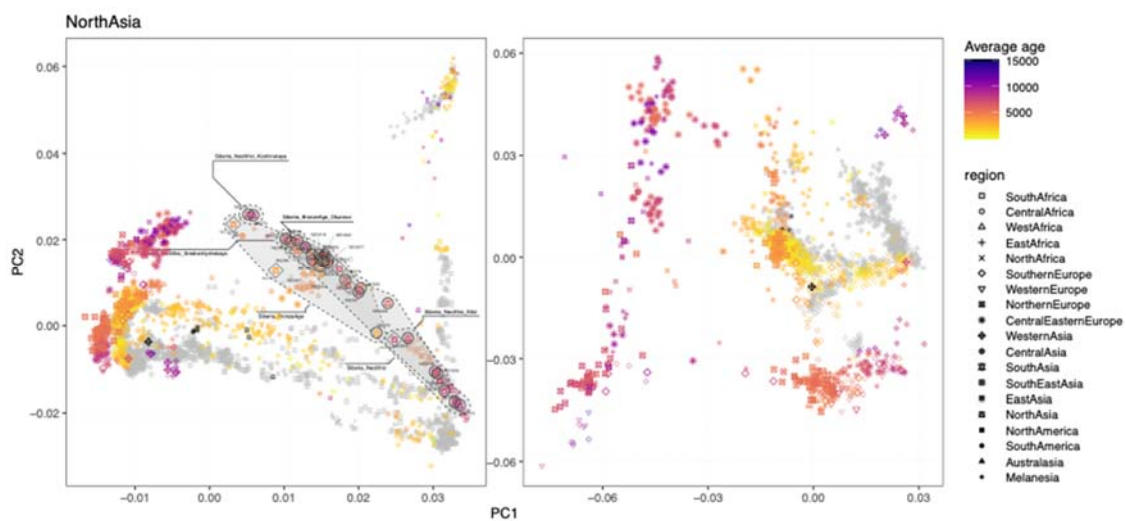
1308 **Fig. S3d.14. Newly reported individuals from Central Asia.** PCA positions of newly
 1309 reported individuals are highlighted with black circles (imputed) or grey diamonds (pseudo-
 1310 haploid, projected). Individuals from the same spatiotemporal group are connected with
 1311 shaded hulls.

1312 We report 7 new individuals from Kazakhstan (6) and Turkmenistan (1). The Neolithic
 1313 individual from Turkmenistan (~6,500 BP) clusters close to Neolithic Iranians. The individuals
 1314 from Kazakhstan are more recent (~4,500BP – 2,000BP), with the older individuals forming
 1315 part of the Neolithic Steppe cline, and the younger individuals along the BA Steppe cline
 1316 (Fig. S3d.14).

1317

1318 North Asia

1319



1320

1321 **Fig. S3d.15. Newly reported individuals from North Asia.** PCA positions of newly
 1322 reported individuals are highlighted with black circles (imputed) or grey diamonds (pseudo-

1323 haploid, projected). Individuals from the same spatiotemporal group are connected with
1324 shaded hulls.

1325

1326 We report 38 new individuals from Western Siberia and Lake Baikal, spanning a period from
1327 ~8,300BP to 2,800 BP. The individuals fall along the entire range of the Neolithic Steppe
1328 cline, spanning from early Forest Steppe hunter-gatherers at the “western” end (NEO72,
1329 NEO73) to Lake Baikal hunter-gatherers at the “eastern” end (Fig. S3d.15)

1330 References

- 1331 1. The 1000 Genomes Project Consortium,. A global reference for human genetic
1332 variation. *Nature* **526**, 68–74 (2015).
- 1333 2. Patterson, N. *et al.* Ancient Admixture in Human History. *Genetics* **192**, 1065–1093
1334 (2012).
- 1335 3. Lazaridis, I. *et al.* Ancient human genomes suggest three ancestral populations for
1336 present-day Europeans. *Nature* **513**, 409–413 (2014).
- 1337 4. Lazaridis, I. *et al.* Genomic insights into the origin of farming in the ancient Near East.
1338 *Nature* **536**, 419–424 (2016).
- 1339 5. Pickrell, J. K. *et al.* The genetic prehistory of southern Africa. *Nat. Commun.* **3**, 1143
1340 (2012).
- 1341 6. Shringarpure, S. S., Bustamante, C. D., Lange, K. & Alexander, D. H. Efficient
1342 analysis of large datasets and sex bias with ADMIXTURE. *BMC Bioinformatics* **17**, 218
1343 (2016).
- 1344 7. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: A Tool for Genome-
1345 wide Complex Trait Analysis. *Am. J. Hum. Genet.* **88**, 76–82 (2011).
- 1346 8. Patterson, N., Price, A. L. & Reich, D. Population Structure and Eigenanalysis. *PLoS*
1347 *Genet* **2**, e190 (2006).
- 1348 9. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and
1349 richer datasets. *GigaScience* **4**, 7 (2015).
- 1350 10. Antonio, M. L. *et al.* Ancient Rome: A genetic crossroads of Europe and the
1351 Mediterranean. *Science* **366**, 708–714 (2019).

1352

1353

1354

1355
1356
1357
1358
1359
1360
1361
1362

3e) Inferring the spatiotemporal spread of population movements in the past 13 millennia

Fernando Racimo¹

¹Lundbeck Foundation GeoGenetics Centre, GLOBE Institute, University of Copenhagen, Copenhagen, Denmark

1363 Introduction

1364 We aimed to infer the geographic and temporal spread of major population movements in the
1365 past 13 millennia of Western Eurasian history. We used a method developed in ¹, which uses
1366 spatiotemporal ordinary kriging on latent ancestry proportion estimates from ancient and
1367 present-day genomes. This way, we obtained detailed spatiotemporal maps reflecting the
1368 dynamics of the spread of ancestry during the transition from the Mesolithic to the Neolithic,
1369 Bronze Age, Iron Age and more recent periods, finally resulting in the complex ancestry
1370 make-up of present-day populations in the region.

1371 Methods

1372 We obtained ancestry proportions estimated using Admixture² with K=9 latent ancestry
1373 clusters ([Supplementary Note S3d](#)) on a sequence dataset including both whole-genome
1374 shotgun-sequenced genomes and genomic sequences obtained via SNP capture, after
1375 imputation ([Supplementary Note S2](#)). We performed spatiotemporal kriging³ of these
1376 proportions over the last 12,900 years, in intervals of 300 years, with a 5,000-point spatial
1377 grid spanning Western and Central Eurasia. We used the R package *gstat* to fit a
1378 spatiotemporal variogram via a metric covariance model, and perform ordinary kriging⁴. We
1379 focused on the ancestry clusters for which we could fit variogram models that were not static
1380 over time.

1381 Results

1382 We were able to fit spatiotemporal variogram functions to six of the nine ancestries. We label
1383 these as WHG, EHG, IRN, LVN, SIB and EAS. The first four are roughly maximised in
1384 Mesolithic western European hunter-gatherers, Mesolithic eastern European hunter-
1385 gatherers, Iranian Neolithic populations, Levant Neolithic populations and ancient Siberian
1386 populations, respectively (see ⁵ for a model positing the first four of these populations as the
1387 major sources of ancestry in present-day Europeans). We depict the spatiotemporal spread

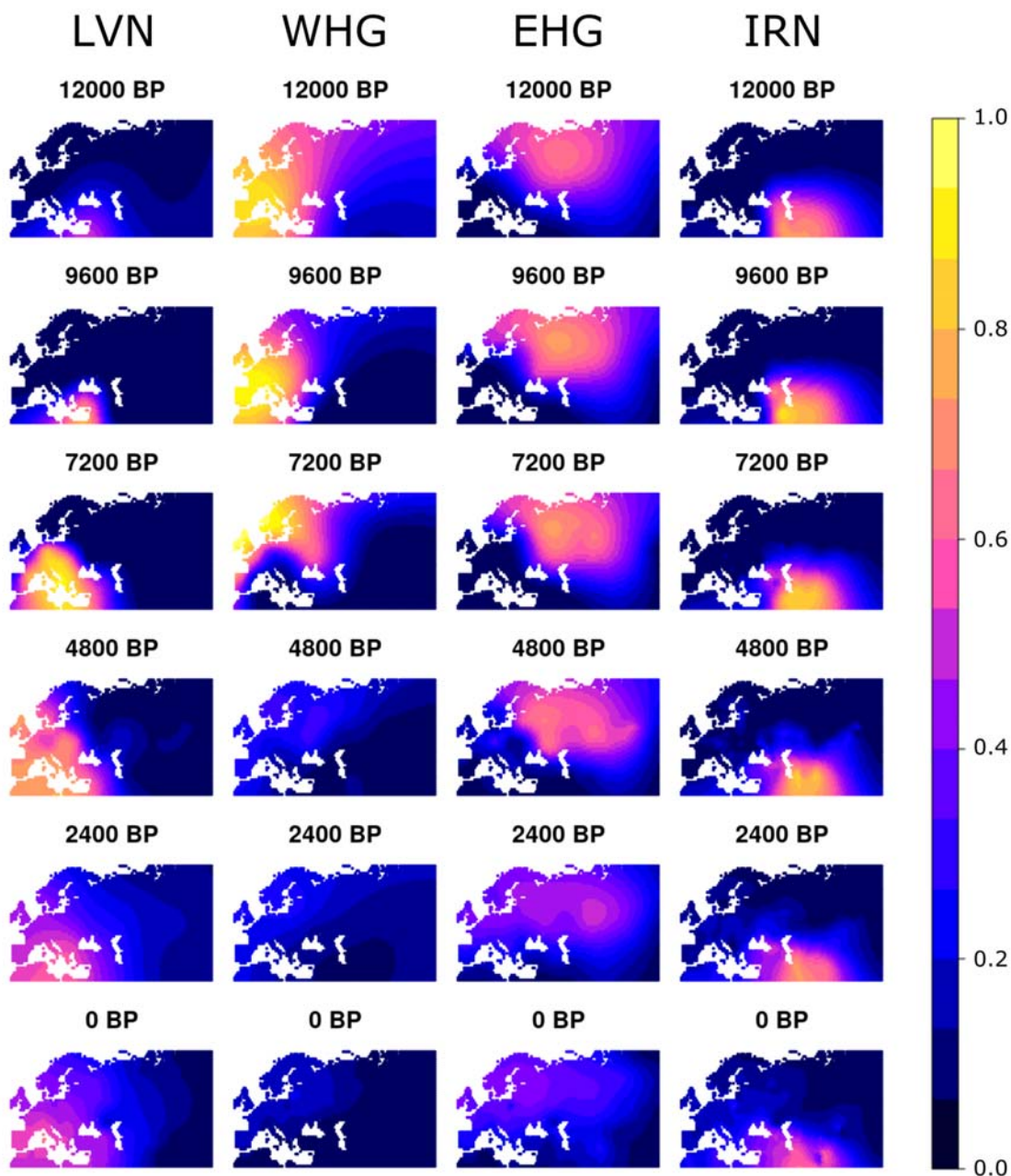
1388 of the first four of these ancestries in Fig. S3e.1 and Supplementary Animations 1-4. The fifth
1389 ancestry (SIB) occurs at much lower rates in western Eurasia, and rises in frequency in
1390 northeastern Europe during the Iron Age (Fig. S3e.2, Supplementary Animation 5)⁶⁻⁸. A sixth
1391 ancestry (EAS) has affinities to East Asians, and expands into the Caucasus in recent times
1392 (Fig. S3e.2, Supplementary Animation 6).

1393 These spatiotemporal maps evince interesting patterns of ancestry change across the
1394 landscape. For example, the advancement of Neolithic Levant (LVN) ancestry appears
1395 staggered: we observe different periods of advancement followed by stasis. In addition to the
1396 Bronze Age movement of EHG ancestry, there is a southern incursion of IRN ancestry via
1397 South Europe⁹⁻¹¹. This is particularly obvious in Bronze Age Greek and Iron Age Roman
1398 samples, and may be due to contacts with Anatolia and Northern Africa (where this ancestry
1399 is also present). Additionally, we observe small incursions of very late SIB ancestry into
1400 Eastern Europe (Fig. S3e.2 Supplementary Animation 5). This signal is driven by the
1401 presence of SIB ancestry in Iron Age Cimmerian nomads¹² and in a medieval Serbian¹³, and
1402 could perhaps be linked to the introduction of languages from the Finno-Ugric family into the
1403 Hungarian Plain. An incursion of this ancestry into Western Eurasia can also be seen in
1404 Medieval Ottoman Anatolians¹⁴.

1405 We can focus on local timelines of kriged ancestry changes in different points of the map
1406 (Fig. S3e.3). Here, we observe that the timing and duration of the rise in LVN ancestry was
1407 different in different points in Europe (Fig. S3e.3). We also observe that, in certain regions of
1408 Europe, the rise in IRN and EHG ancestry are largely decoupled from each other (see e.g.
1409 "Rome" in Fig. S3e.2)^{9-11,15}.

1410

1411

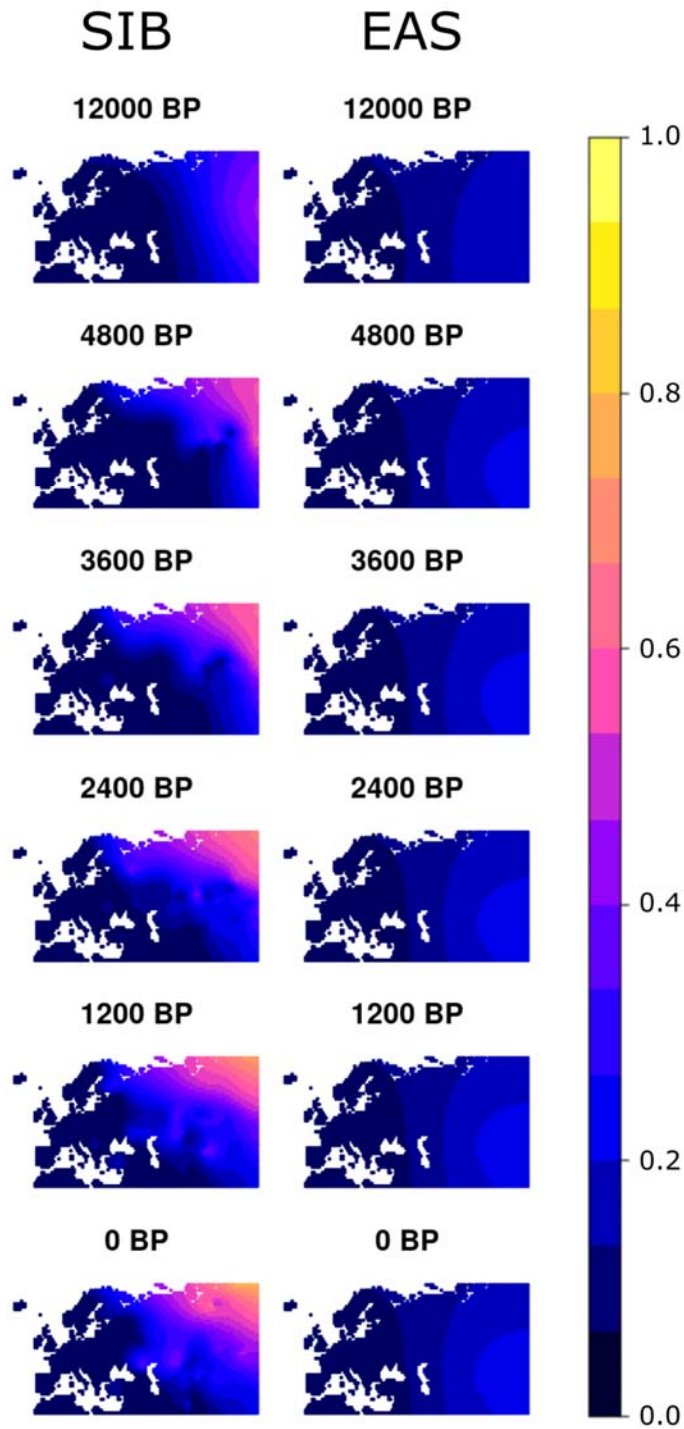


1413

1414 **Figure S3e.1.** Spatiotemporal kriging of four major ancestry clusters over the last 12,000
 1415 years of human history. LVN = ancestry maximised in Anatolian farmer populations. WHG =
 1416 ancestry maximised in western European hunter-gatherers. EHG = ancestry maximised in
 1417 eastern European hunter-gatherers. IRN = ancestry maximised in Iranian Neolithic
 1418 individuals and Caucasus hunter-gatherers.

1419

1420



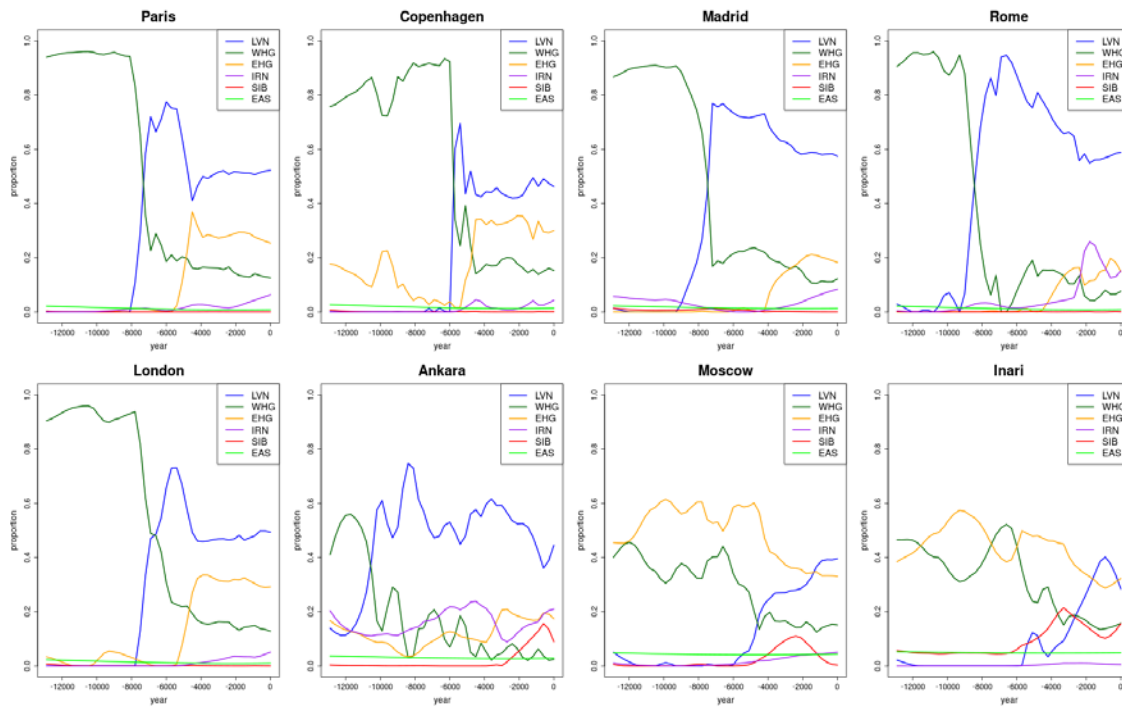
1421

1422 **Figure S3e.2.** Spatiotemporal kriging of two additional ancestry clusters with later incursions
 1423 into Western Eurasia over the last 12,000 years of human history, particularly focusing on
 1424 the last 5,000 years. SIB = ancestry maximised in ancient Siberian individuals. EAS =
 1425 ancestry maximised in East Asian individuals.

1426

1427

1428

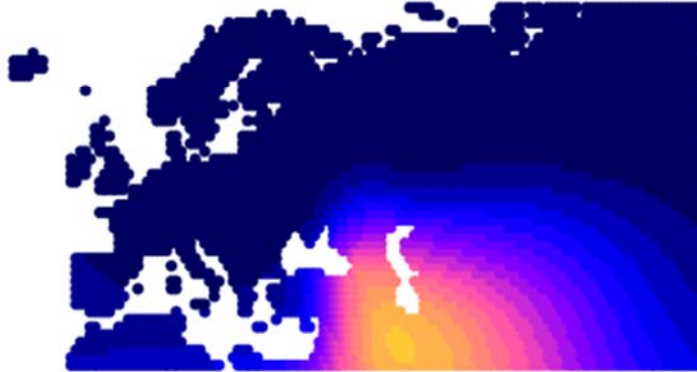


1429

1430 **Figure S3e.3.** Local timelines in what are now 8 urban centres across western Eurasia,
1431 reflecting local differences in the tempo and mode of ancestry changes over time. LVN =
1432 ancestry maximised in Anatolian farmer populations. WHG = ancestry maximised in western
1433 European hunter-gatherers. EHG = ancestry maximised in eastern European hunter-
1434 gatherers. IRN = ancestry maximised in Iranian Neolithic individuals / Caucasus hunter-
1435 gatherers. SIB = ancestry maximised in ancient Siberian individuals. EAS = ancestry
1436 maximised in East Asian individuals.

1437 **Animation S3e.1. IRN**

-12900

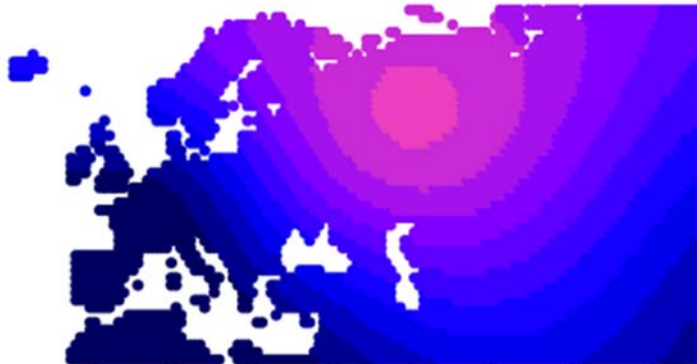


1438

1439

1440 **Animation S3e.2. EHG**

-12900



1441

1442

1443 **Animation S3e.3. LVN**

-12900



1444

1445

1446 **Animation S3e.4. WHG**

-12900



1447

1448 **References**

- 1449 1. Racimo, F. *et al.* The spatiotemporal spread of human migrations during the
1450 European Holocene. *Proc. Natl. Acad. Sci. U. S. A.* (2020)
1451 doi:10.1073/pnas.1920051117.

- 1452 2. Alexander, D. H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry
1453 in unrelated individuals. *Genome Res.* **19**, 1655–1664 (2009).
- 1454 3. Cressie, N. & Wikle, C. K. *Statistics for Spatio-Temporal Data*. (John Wiley & Sons,
1455 2015).
- 1456 4. Gräler, B., Pebesma, E. & Heuvelink, G. Spatio-Temporal Interpolation using gstat.
1457 *The R Journal* vol. 8 204 (2016).
- 1458 5. Lazaridis, I. *et al.* Genomic insights into the origin of farming in the ancient Near East.
1459 *Nature* **536**, 419–424 (2016).
- 1460 6. Lamnidis, T. C. *et al.* Ancient Fennoscandian genomes reveal origin and spread of
1461 Siberian ancestry in Europe. *Nat. Commun.* **9**, 5018 (2018).
- 1462 7. Tambets, K. *et al.* Genes reveal traces of common recent demographic history for
1463 most of the Uralic-speaking populations. *Genome Biology* vol. 19 (2018).
- 1464 8. Saag, L. *et al.* The Arrival of Siberian Ancestry Connecting the Eastern Baltic to
1465 Uralic Speakers further East. *Curr. Biol.* **29**, 1701–1711.e16 (2019).
- 1466 9. Antonio, M. L. *et al.* Ancient Rome: A genetic crossroads of Europe and the
1467 Mediterranean. *Science* **366**, 708–714 (2019).
- 1468 10. Mathieson, I. *et al.* The genomic history of southeastern Europe. *Nature* **555**,
1469 197–203 (2018).
- 1470 11. Marcus, J. H. *et al.* Genetic history from the Middle Neolithic to present on the
1471 Mediterranean island of Sardinia. *Nat. Commun.* **11**, 939 (2020).
- 1472 12. Krzewińska, M. *et al.* Ancient genomes suggest the eastern Pontic-Caspian
1473 steppe as the source of western Iron Age nomads. *Sci Adv* **4**, 4457 (2018).
- 1474 13. Veeramah, K. R. *et al.* Population genomic analysis of elongated skulls reveals extensive female-biased immigration in Early Medieval Bavaria.
1475 *Proc. Natl. Acad. Sci. U. S. A.* **115**, 3494–3499 (2018).
- 1476 1476 14. Damgaard, P. de B. *et al.* The first horse herders and the impact of early
1477 Bronze Age steppe expansions into Asia. *Science* vol. 360 7711 (2018).
- 1478 1478 15. Fernandes, D. M. *et al.* The spread of steppe and Iranian-related ancestry in
1479

1480 the islands of the western Mediterranean. *Nat Ecol Evol* **4**, 334–345 (2020).

1481

1482

1483 3f) HBD/ IBD sharing/ROH/clustering

1484 Martin Sikora¹

1485 ¹Lundbeck Foundation GeoGenetics Centre, GLOBE Institute, University of Copenhagen,
1486 Copenhagen, Denmark

1487

1488 Methods

1489 We used *IBDseq*¹ to detect genomic segments shared identical-by-descent (IBD) between all
1490 individuals in the “1000G” dataset, restricting to transversion SNPs with imputation INFO
1491 score ≥ 0.8 and MAF ≥ 0.01 . We filtered the resulting IBD segments for LOD score ≥ 3 and a
1492 minimum length of 2 centimorgans (cM), and further removed regions of excess IBD
1493 following the approach of Browning and Browning². First, we used the *GenomicRanges*³
1494 package in R to calculate the total number of IBD segments overlapping each position along
1495 the genome, and calculated their 3% trimmed mean and standard deviation (SD). We then
1496 called regions of excess IBD if they were > 3 trimmed SD from the trimmed mean. We split
1497 IBD segments overlapping the excess IBD regions, and removed any segments with length $<$
1498 2 cM after splitting. For analyses of runs of homozygosity (ROH) we used a shorter length
1499 cutoff of 1cM.

1500

1501 We carried out genetic clustering of the ancient individuals using hierarchical community
1502 detection on a network of pairwise identity-by-descent (IBD)-sharing similarities⁴. To facilitate
1503 detection of clusters at a finer scale, we ran *IBDseq* on a dataset restricting to ancient
1504 samples only, and applied more lenient filters of imputation INFO score > 0.5 , and minimum
1505 IBD segment length of 1 cM. We constructed a weighted network of the individuals using the
1506 *igraph*⁵ package in R, with the fraction of the genome shared IBD between pairs of
1507 individuals as weights. We then performed iterative community detection on this network
1508 using the Leiden algorithm⁶ implemented in the *leidenAlg* R package (*v1.01*,
1509 <https://github.com/kharchenkolab/leidenAlg>). We used a resolution parameter of $r=0.5$ as the
1510 starting value for each level of community detection. If more than one community was
1511 detected, we split the network into the respective communities, and repeated the community
1512 detection step. If no communities were detected, we incremented the resolution parameter in
1513 steps of 0.5 until a maximum value of $r=3$. The initial clustering was completed when no
1514 more communities were detected at the highest resolution parameter, across all

1515 subcommunities. To convert the resulting hierarchy into a final clustering, we simplified the
1516 initial clustering by collapsing nodes into single clusters based on observed spatiotemporal
1517 annotations of the samples.

1518

1519 To estimate ancestry proportions from patterns of pairwise IBD sharing, we used an
1520 approach akin to “chromosome painting”⁷. We first inferred an IBD-based “painting profile”
1521 for each target individual, by summing up the total amount of IBD shared with each “donor”
1522 group (using population labels for modern donors or IBD-based genetic clusters for ancient
1523 donors), and normalising them to the interval [0, 1]. We used a leave-one-out approach as in⁸
1524 to account for the fact that recipient individuals cannot be included as donors from their own
1525 group. We then used these painting profiles in supervised modelling of target individuals as
1526 mixtures from different sets of putative source groups^{8,9}, using non-negative least squares
1527 implemented in the R package *limSolve*¹⁰.

1528

1529 To investigate ancestry compositions across the full set of ancient individuals, we used three
1530 sets of source groups reflecting different temporal depths:

1531

- 1532 - “deep”, a set of groups representing highly differentiated deep ancestry sources
- 1533 - “postNeol”, using diverse Neolithic and earlier source groups
- 1534 - “postBA”, using Late Neolithic and Bronze Age source groups

1535

1536 We also performed two analyses of more restricted spatiotemporal scope:

1537

- 1538 - “hgEur”, modelling European hunter-gatherers as mixtures of early European hunter-
1539 gatherers and selected outgroups
- 1540 - “fEur”, modelling later European farmers as mixtures of earlier farmers and hunter-
1541 gatherers
- 1542 - “postNeolScand”, modelling Scandinavian Late Neolithic and early Bronze Age
1543 individuals as mixtures of other European early Bronze Age groups.

1544

1545 Results

1546 IBD-based hierarchical graph clustering

1547 We performed hierarchical graph clustering on the 1,492 ancient individuals passing all
1548 filters, which were assigned into a final curated set of 122 genetic clusters (Fig. S3f.1). The
1549 obtained clusters captured both broad and finer-scale genetic structure, corresponding to
1550 shared ancestry within particular spatiotemporal ranges and/or archaeological contexts (Fig.

1551 **S3f.2**). We named these cluster using a “geographic-temporal” nomenclature¹¹ (e.g.
1552 “Denmark_10500BP_6000BP”), in concert with more traditional names for groups of multiple
1553 clusters with shared archaeological or subsistence contexts (e.g. “Farmer_Europe_early”)
1554 where applicable.

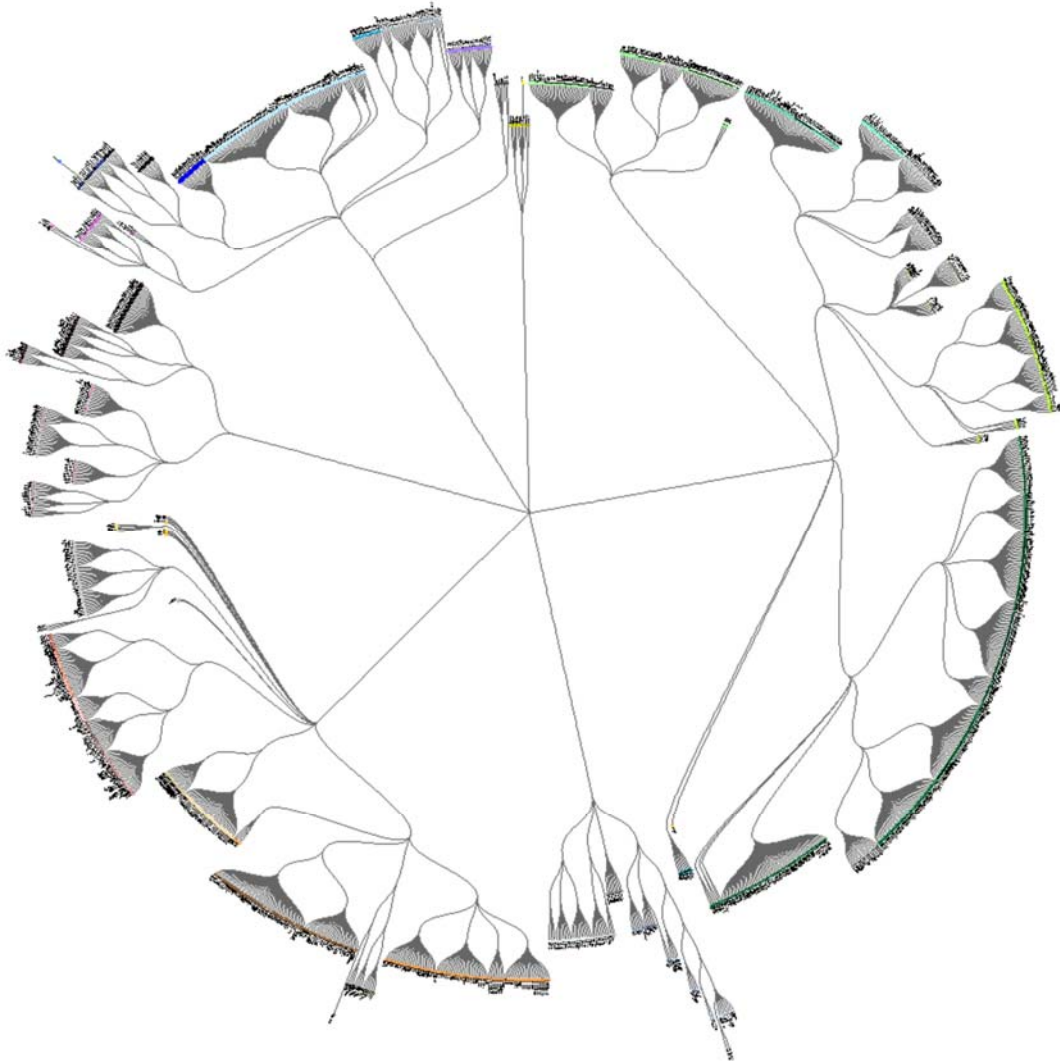
1555

1556 At the highest level of the clustering hierarchy, the individuals were partitioned into six global
1557 clusters representing broad continent-wide genetic structure.

1558

- 1559 - Africa_8000BP_400BP
- 1560 - Europe_15000BP_4000BP
- 1561 - EuropeWCAAsia_25000BP_300BP
- 1562 - Eurasia_5000BP_200BP
- 1563 - Asia_45000BP_200BP
- 1564 - Americas_12000BP_100BP

1565 The following sections provide more detailed descriptions of relevant sub-clusters within the
1566 four global clusters from Eurasia.



1567

1568 **Fig. S3f.1. Hierarchical graph clustering.** Tree diagram showing final curated hierarchical

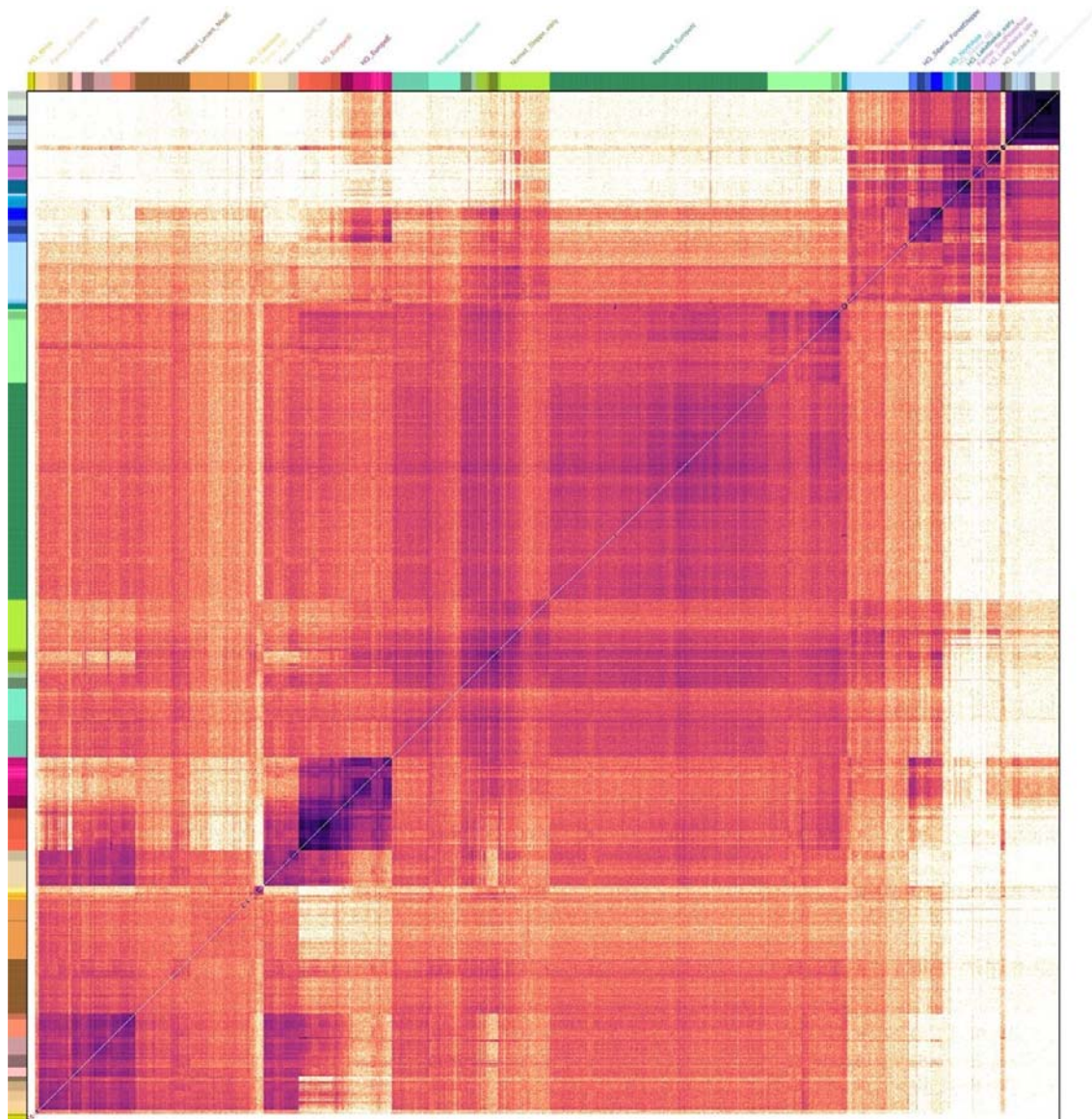
1569 clustering relationship among the 1,492 ancient individuals passing all filters. Genetic

1570 clusters are differentiated using plot symbol colours and shapes.

1571

1572

1573



1574

1575 **Fig. S3f.2. IBD sharing similarities.** Heatmap of pairwise IBD-sharing similarities between
 1576 the 1,492 ancient individuals passing all filters, sorted according to clustering hierarchy.
 1577 Colored bars indicate cluster membership of individuals. Selected cluster group labels are
 1578 shown in the top colour bar.

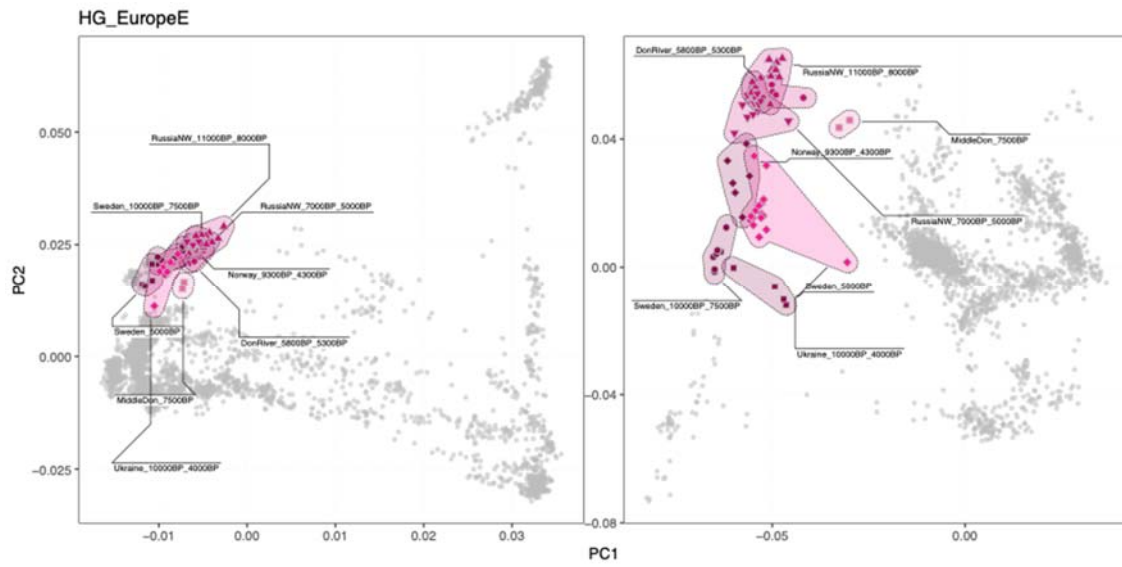
1579

1580 Europe 15000BP 4000BP

1581

1582 This global cluster includes individuals from western Eurasian Mesolithic and Neolithic
 1583 contexts with hunter-gatherer ancestry. The individual genetic clusters are partitioned into
 1584 two cluster groups corresponding to “Eastern hunter-gatherers” and “Western hunter-
 1585 gatherers” as previously used in the literature (Fig. S3f.3-6):

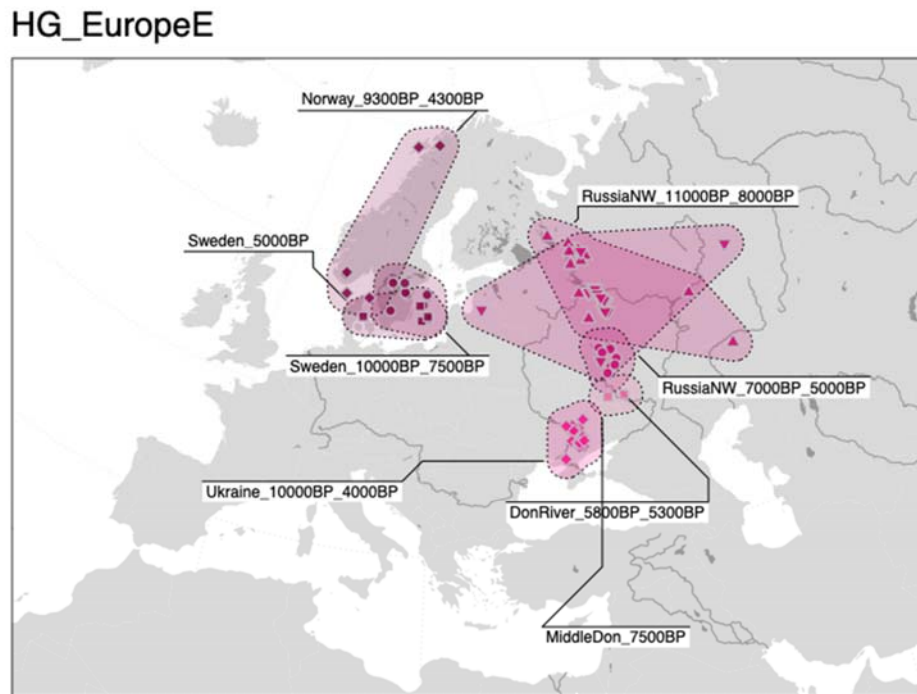
1586



1587

1588 **Fig. S3f.3** PCA for cluster group *HG_EuropeE*. PCA positions of individuals within specific
1589 clusters are highlighted with colored symbols, and connected with shaded hulls (from PCAs
1590 shown in Fig. S3d.7).

1591



1592

1593 **Fig. S3f.4** Geographic distribution of individuals in cluster group *HG_EuropeE*.

1594 Geographic locations of individuals within specific clusters are highlighted with colored

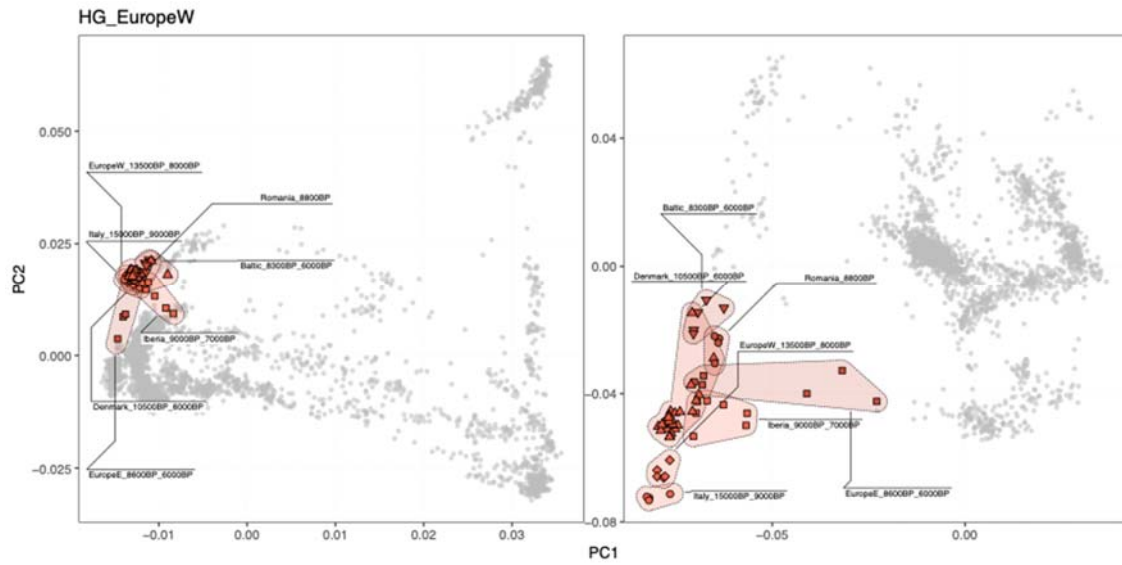
1595 symbols, and connected with shaded hulls (from PCAs shown in Fig. S3d.7).

1596

1597

1598

1599



1600

1601

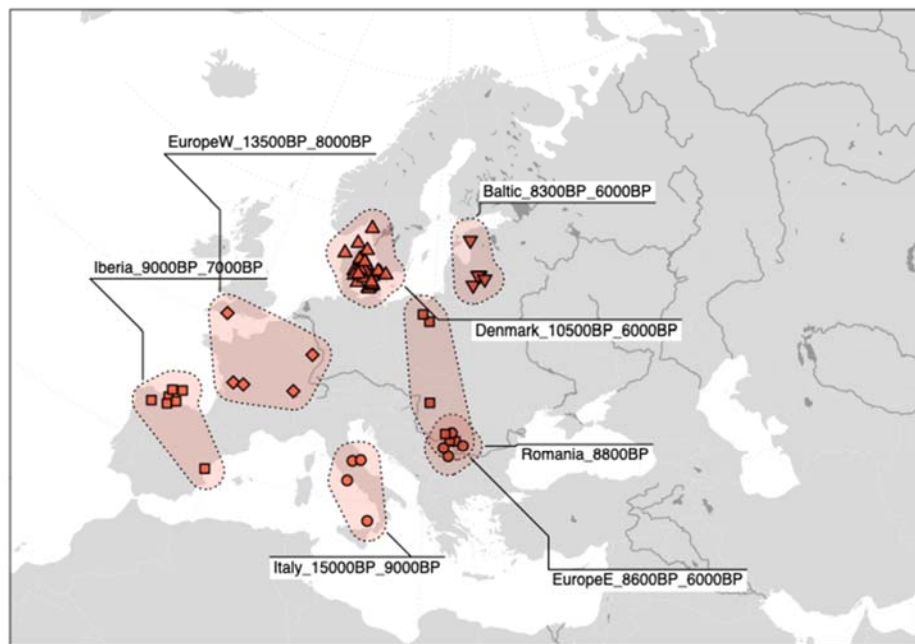
1602

1603

1604

Fig. S3f.5 PCA for cluster group *HG_EuropeW*. PCA positions of individuals within specific clusters are highlighted with colored symbols, and connected with shaded hulls (from PCAs shown in Fig. S3d.7).

HG_EuropeW



1605

1606 **Fig. S3f.6 Geographic distribution of individuals in cluster group *HG_EuropeW*.**

1607 Geographic locations of individuals within specific clusters are highlighted with colored

1608 symbols, and connected with shaded hulls (from PCAs shown in Fig. S3d.7).

1609

1610 EuropeWCAAsia_25000BP_300BP

1611

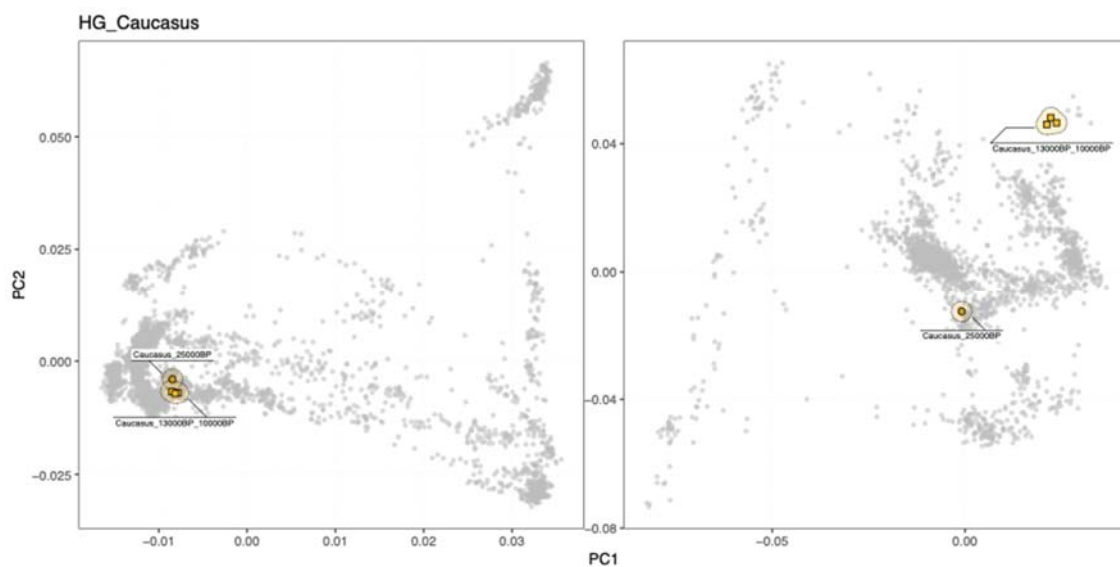
1612 This global cluster includes individuals from western Eurasia with ancestry related to

1613 Mesolithic and Neolithic Near Eastern groups. The individual genetic clusters are partitioned

1614 into a total of eight cluster groups, including all clusters of “Neolithic farmers” and “Caucasus

1615 hunter-gatherers” previously used in the literature (Fig. S3f.7-20).

1616



1617

1618 **Fig. S3f.7 PCA for cluster group *HG_Caucasus*.** PCA positions of individuals within

1619 specific clusters are highlighted with colored symbols, and connected with shaded hulls (from

1620 PCAs shown in Fig. S3d.7).

1621

HG_Caucasus



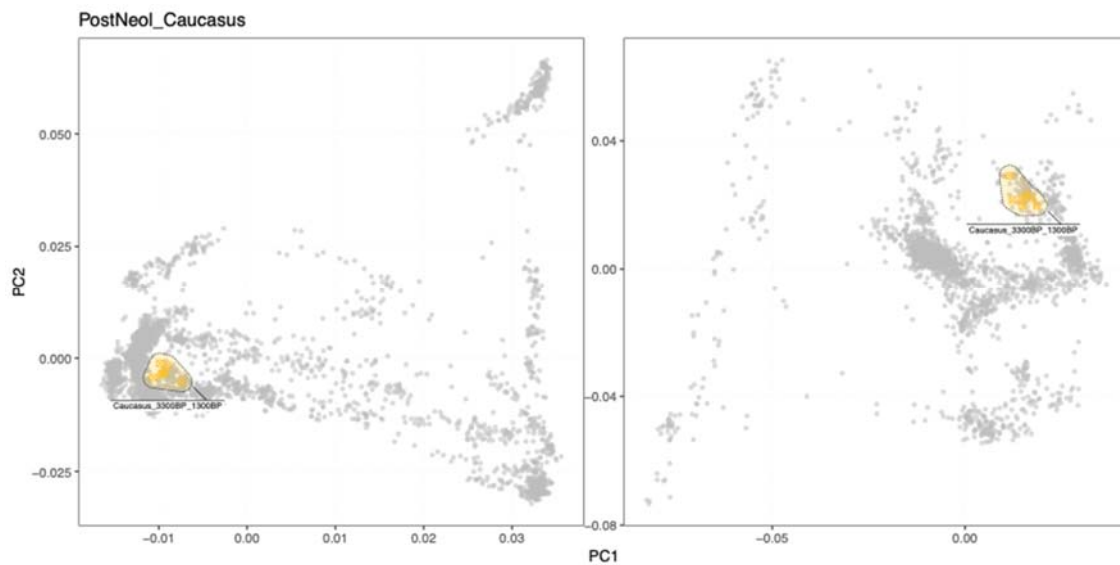
1622

1623 **Fig. S3f.8** Geographic distribution of individuals in cluster group *HG_Caucasus*.

1624 Geographic locations of individuals within specific clusters are highlighted with colored symbols

1625 and connected with shaded hulls (from PCAs shown in Fig. S3d.7).

1626



1627

1628 **Fig. S3f.9** PCA for cluster group *PostNeol_Caucasus*. PCA positions of individuals within

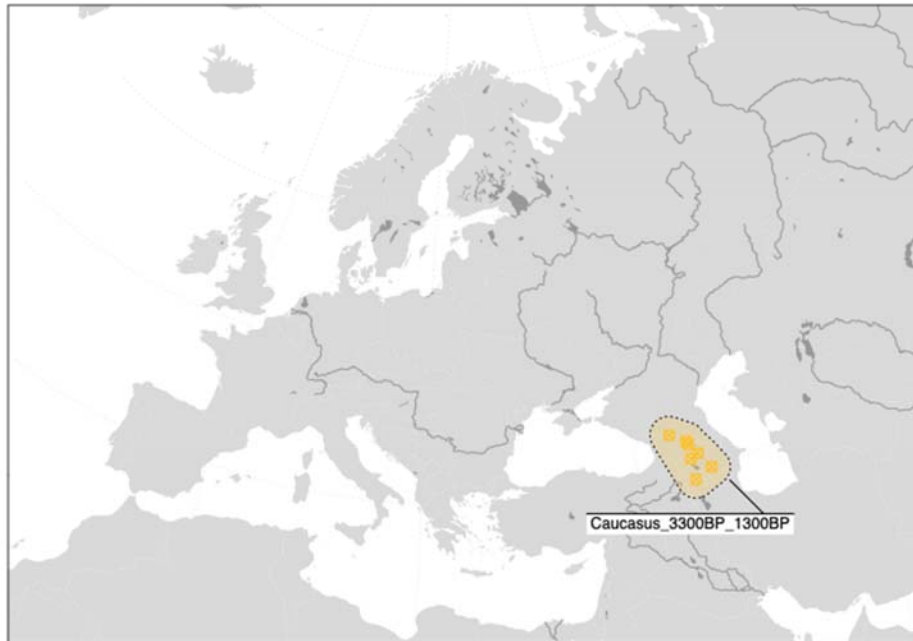
1629 specific clusters are highlighted with colored symbols, and connected with shaded hulls (from

1630 PCAs shown in Fig. S3d.7).

1631

1632

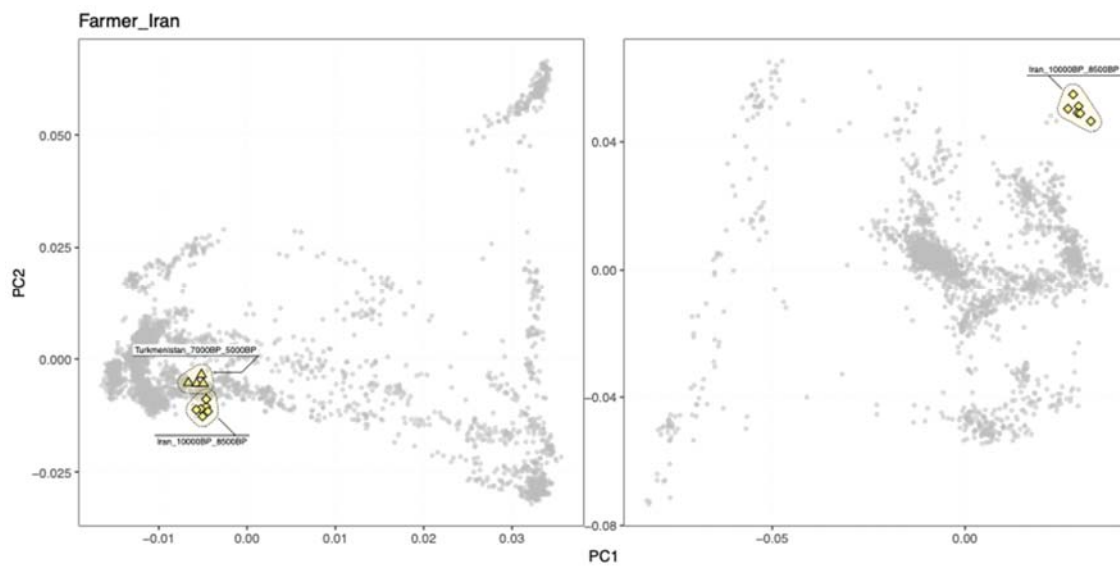
PostNeol_Caucasus



1633

1634 **Fig. S3f.10 Geographic distribution of individuals in cluster group**

1635 *PostNeol_Caucasus*. Geographic locations of individuals within specific clusters are
1636 highlighted with colored symbols, and connected with shaded hulls (from PCAs shown in Fig.
1637 **S3d.7**).



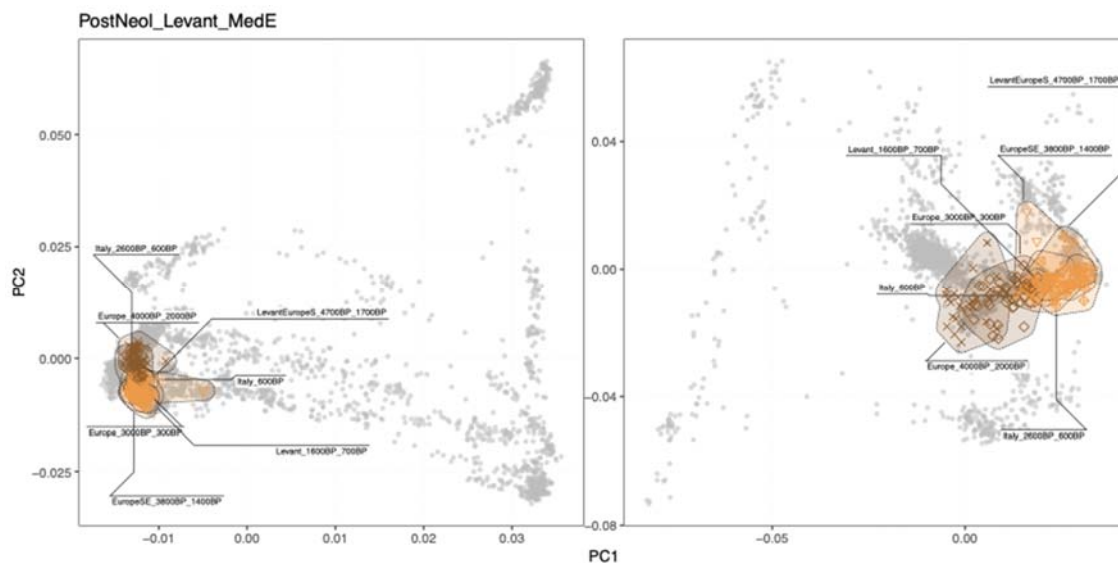
1638

1639

1640 **Fig. S3f.11** PCA for cluster group *Farmer_Iran*. PCA positions of individuals within specific
 1641 clusters are highlighted with colored symbols, and connected with shaded hulls (from PCAs
 1642 shown in Fig. S3d.7).
 1643
 1644

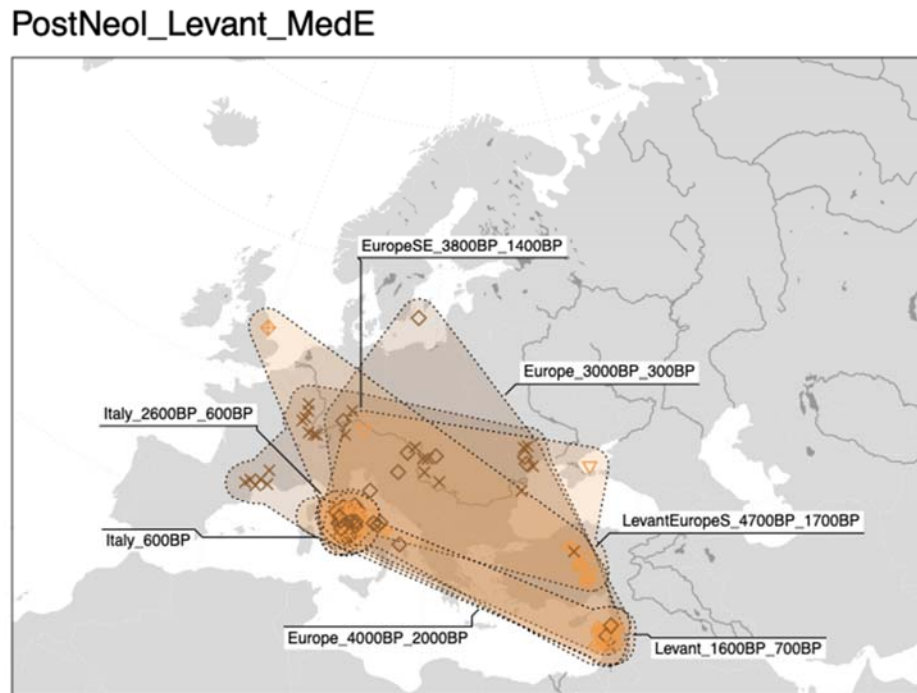


1645
 1646 **Fig. S3f.12** Geographic distribution of individuals in cluster group *Farmer_Iran*.
 1647 Geographic locations of individuals within specific clusters are highlighted with colored
 1648 symbols, and connected with shaded hulls (from PCAs shown in Fig. S3d.7).

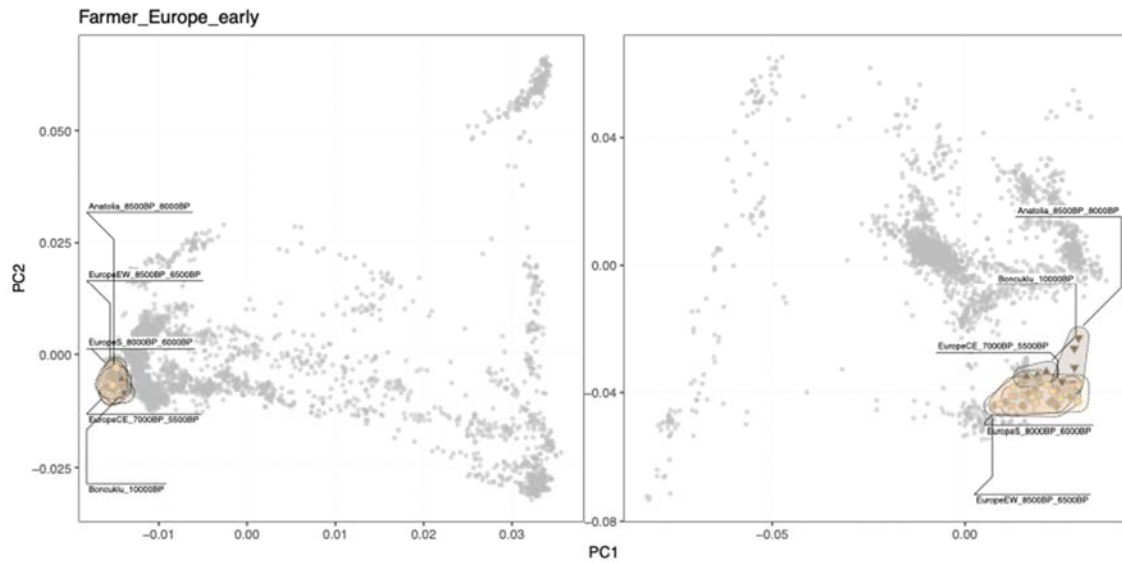


1649

1650 **Fig. S3f.13** PCA for cluster group *PostNeol_Levant_MedE*. PCA positions of individuals
1651 within specific clusters are highlighted with colored symbols, and connected with shaded
1652 hulls (from PCAs shown in Fig. S3d.7).
1653
1654



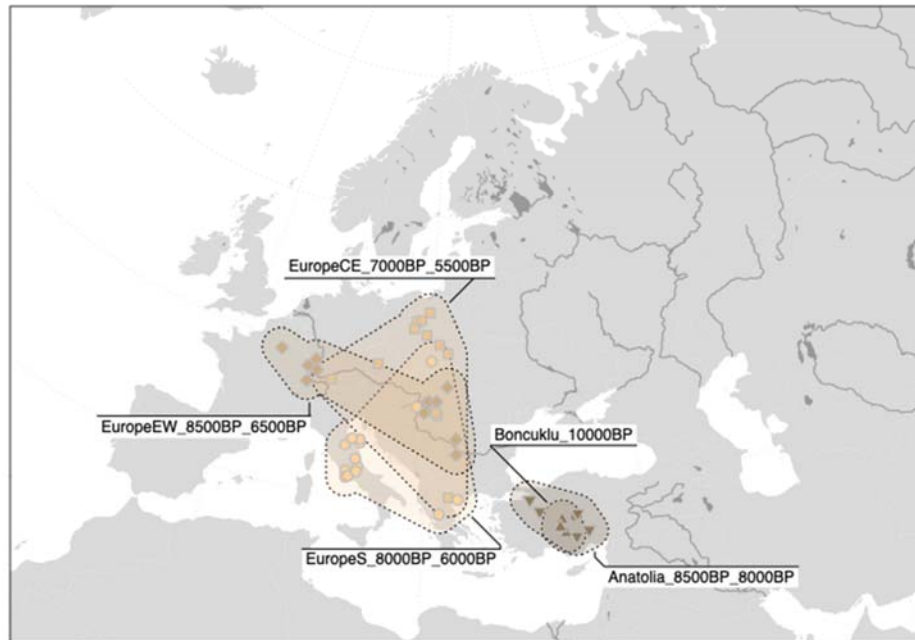
1655
1656 **Fig. S3f.14** Geographic distribution of individuals in cluster group
1657 *PostNeol_Levant_MedE*. Geographic locations of individuals within specific clusters are
1658 highlighted with colored symbols, and connected with shaded hulls (from PCAs shown in Fig.
1659 S3d.7).



1660
 1661
 1662
 1663
 1664
 1665

Fig. S3f.15 PCA for cluster group *Farmer_Europe_early*. PCA positions of individuals within specific clusters are highlighted with colored symbols, and connected with shaded hulls (from PCAs shown in Fig. S3d.7).

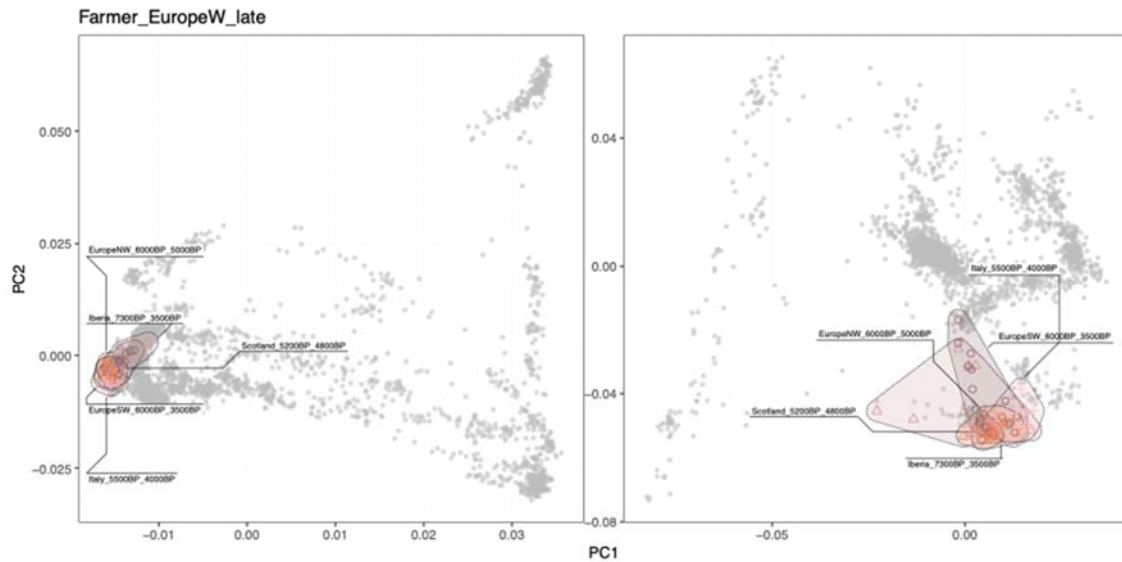
Farmer_Europe_early



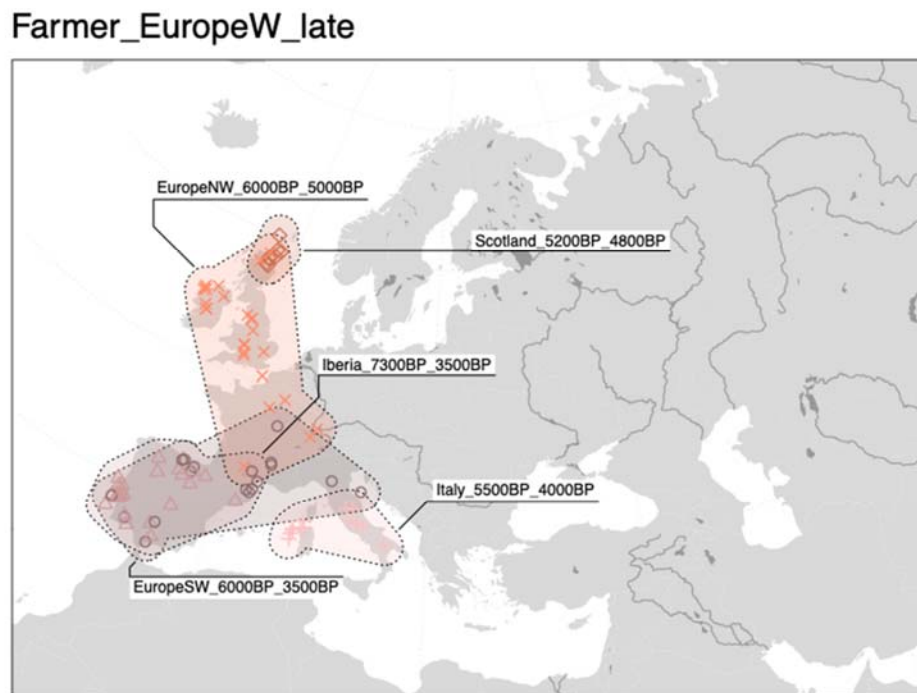
1666
 1667
 1668

Fig. S3f.16 Geographic distribution of individuals in cluster group *Farmer_Europe_early*. Geographic locations of individuals within specific clusters are

1669 highlighted with colored symbols, and connected with shaded hulls (from PCAs shown in Fig.
1670 S3d.7).

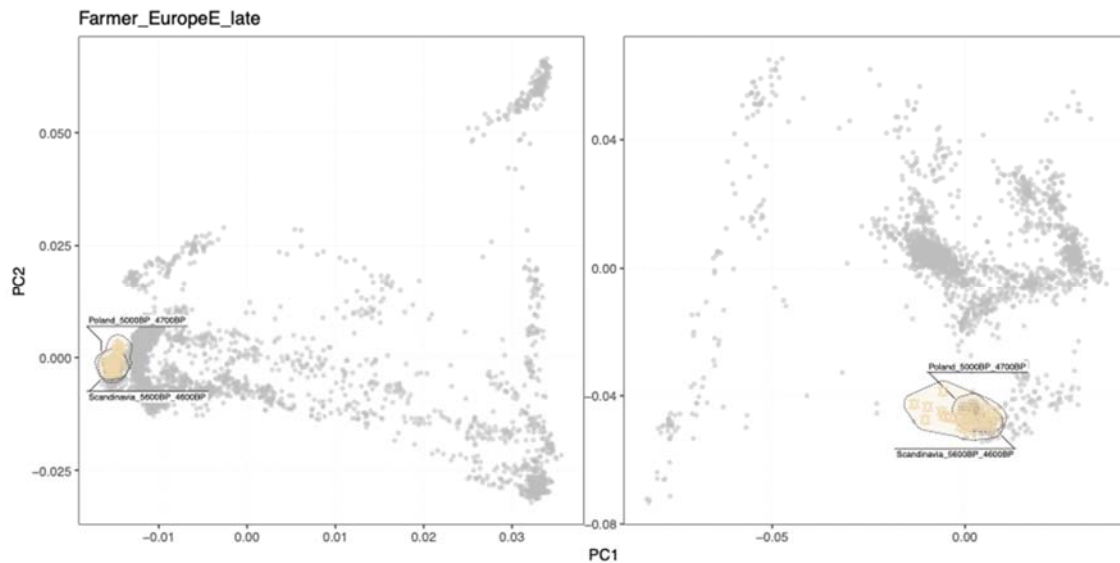


1671
1672 **Fig. S3f.17** PCA for cluster group *Farmer_EuropeW_late*. PCA positions of individuals
1673 within specific clusters are highlighted with colored symbols, and connected with shaded
1674 hulls (from PCAs shown in Fig. S3d.7).
1675
1676



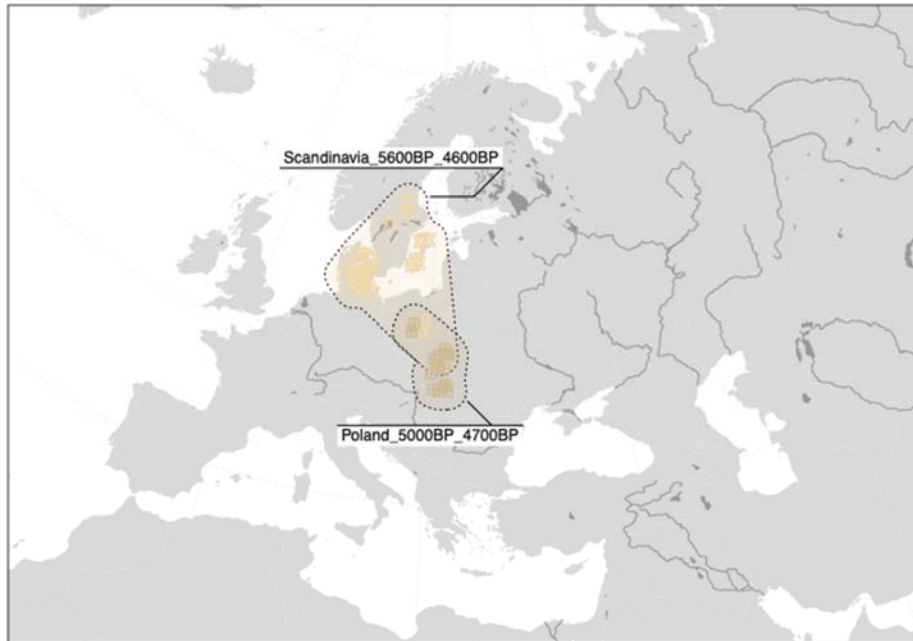
1677

1678 **Fig. S3f.18** Geographic distribution of individuals in cluster group
1679 *Farmer_EuropeW_late*. Geographic locations of individuals within specific clusters are
1680 highlighted with colored symbols, and connected with shaded hulls (from PCAs shown in Fig.
1681 S3d.7).
1682



1683
1684 **Fig. S3f.19** PCA for cluster group *Farmer_EuropeE_late*. PCA positions of individuals
1685 within specific clusters are highlighted with colored symbols, and connected with shaded
1686 hulls (from PCAs shown in Fig. S3d.7).

Farmer_EuropeE_late



1688

1689 **Fig. S3f.20 Geographic distribution of individuals in cluster group**

1690 **Farmer_EuropeE_late.** Geographic locations of individuals within specific clusters are
 1691 highlighted with colored symbols, and connected with shaded hulls (from PCAs shown in Fig.
 1692 S3d.7).

1693

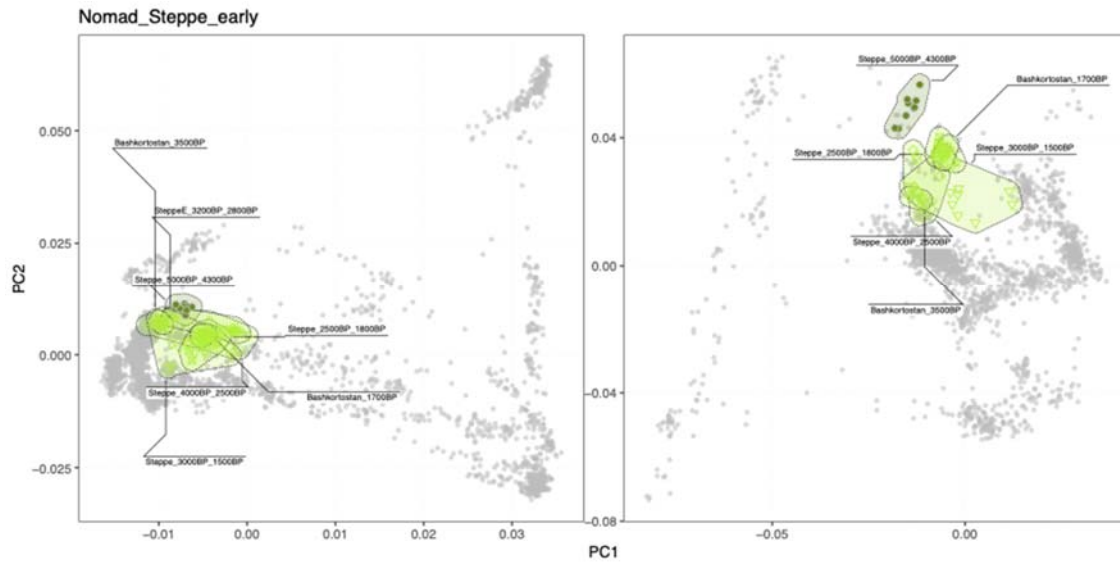
1694 Eurasia 5000BP 200BP

1695 This global cluster includes individuals from western Eurasia from the Bronze Age onwards.

1696 The individual genetic clusters are partitioned into a total of five cluster groups, including all

1697 clusters of individuals with “Steppe ancestry” related to Bronze Age Steppe pastoralists

1698 previously used in the literature (Fig. S3f.21-28).



1700

1701

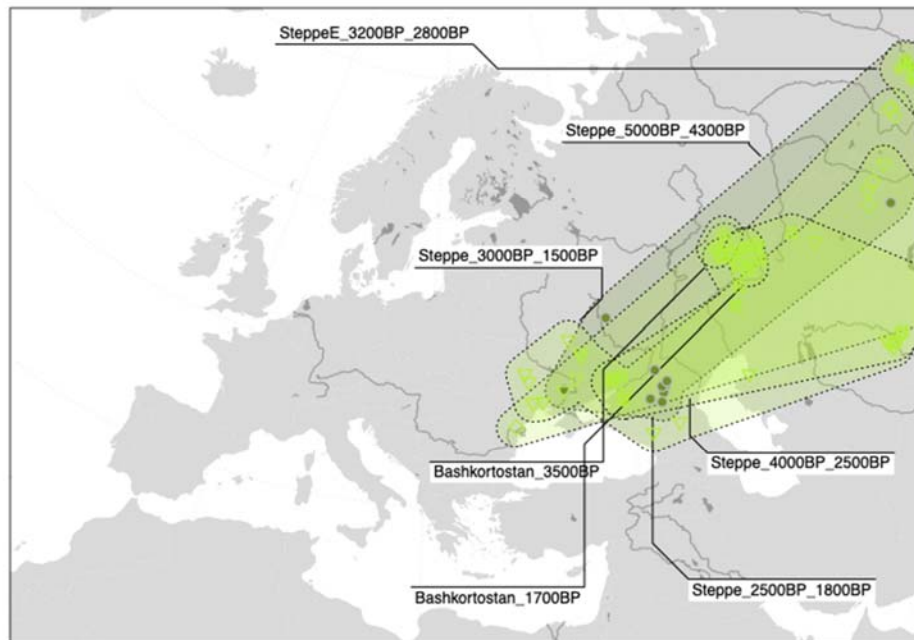
1702

1703

1704

Fig. S3f.21 PCA for cluster group *Nomad_Steppe_early*. PCA positions of individuals within specific clusters are highlighted with colored symbols, and connected with shaded hulls (from PCAs shown in Fig. S3d.7).

Nomad_Steppe_early



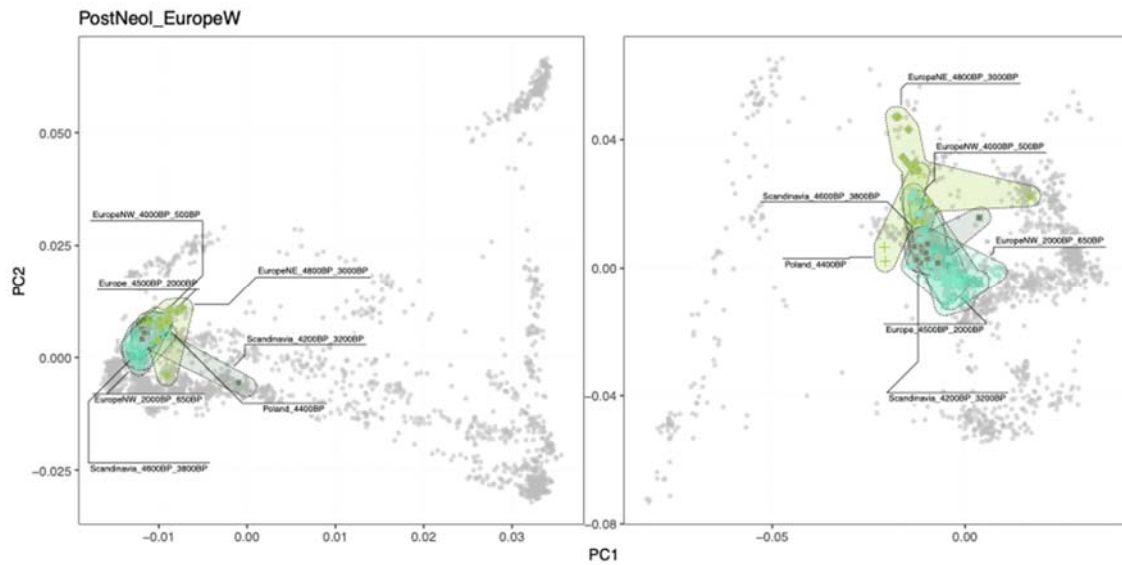
1705

1706

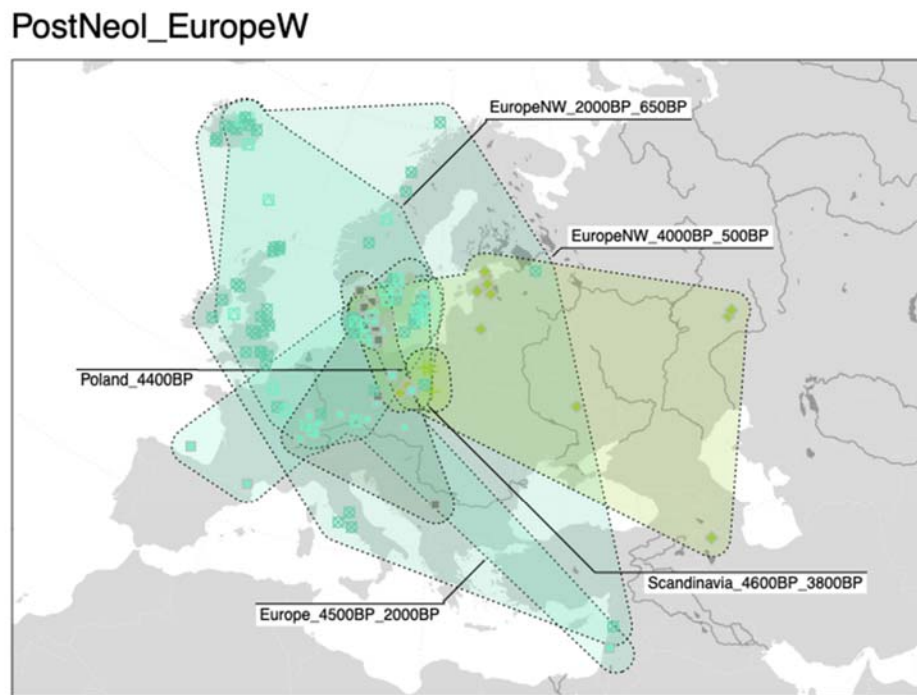
1707

Fig. S3f.22 Geographic distribution of individuals in cluster group *Nomad_Steppe_early*. Geographic locations of individuals within specific clusters are

1708 highlighted with colored symbols, and connected with shaded hulls (from PCAs shown in Fig.
1709 S3d.7).
1710



1711
1712 **Fig. S3f.23** PCA for cluster group *PostNeol_EuropeW*. PCA positions of individuals within
1713 specific clusters are highlighted with colored symbols, and connected with shaded hulls (from
1714 PCAs shown in Fig. S3d.7).
1715



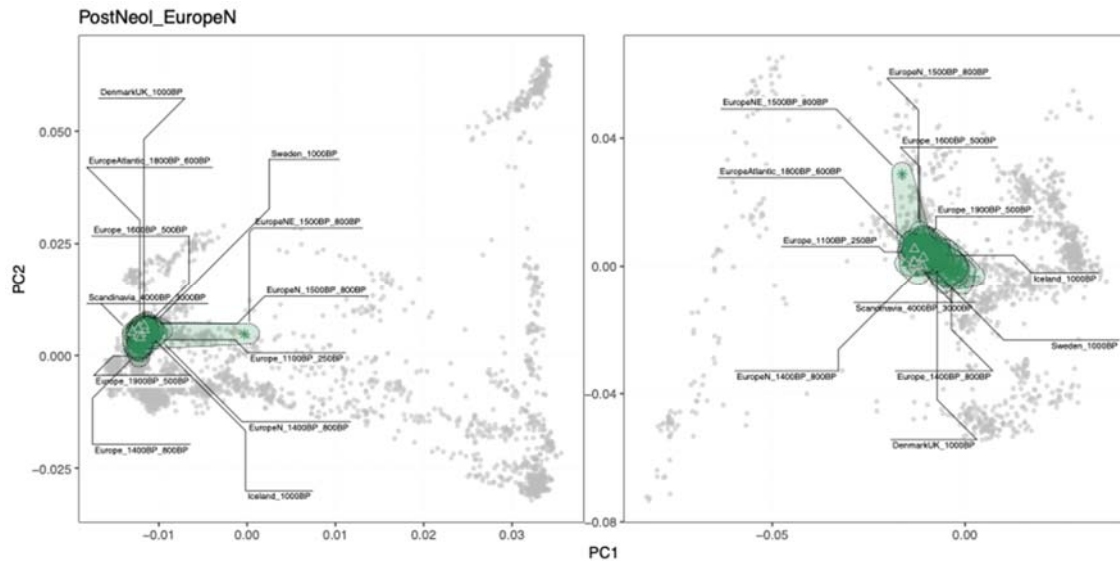
1716

1717 **Fig. S3f.24** Geographic distribution of individuals in cluster group *PostNeol_EuropeW*.

1718 Geographic locations of individuals within specific clusters are highlighted with colored

1719 symbols, and connected with shaded hulls (from PCAs shown in Fig. S3d.7).

1720



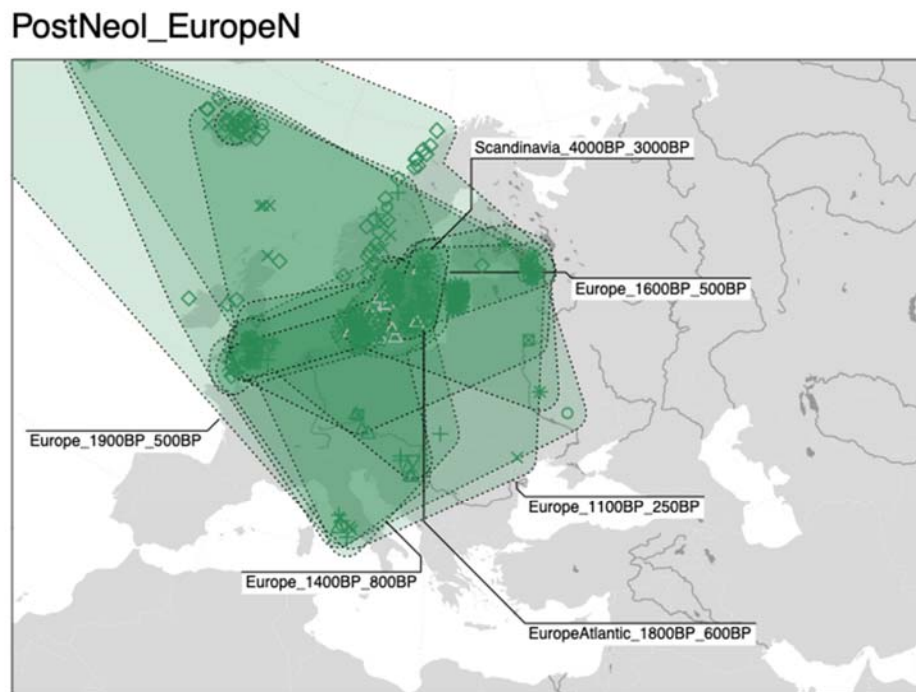
1721

1722 **Fig. S3f.25** PCA for cluster group *PostNeol_EuropeN*. PCA positions of individuals within

1723 specific clusters are highlighted with colored symbols, and connected with shaded hulls (from

1724 PCAs shown in Fig. S3d.7).

1725



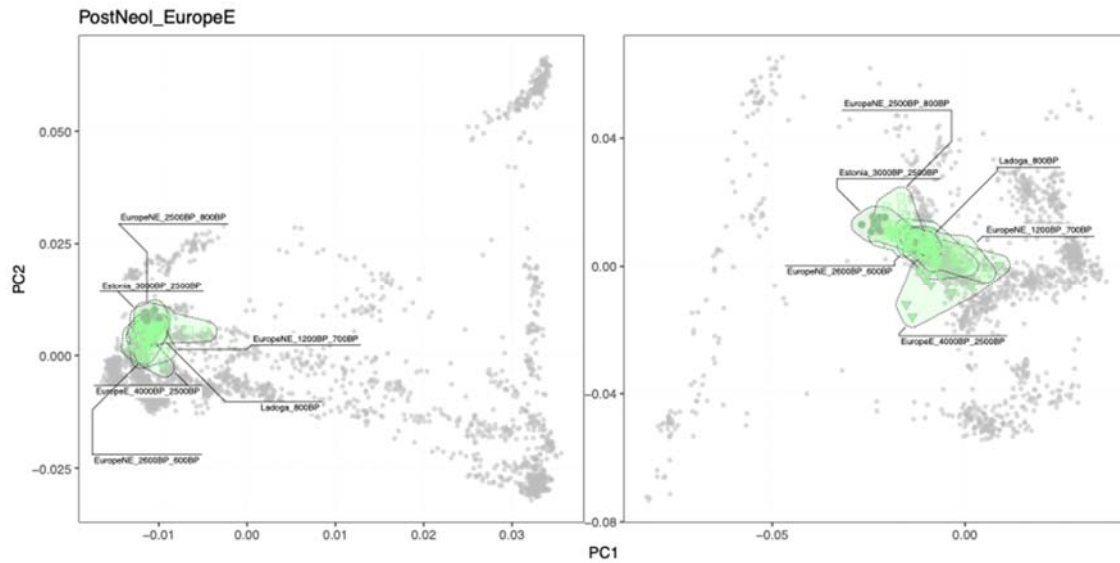
1726

1727 **Fig. S3f.26** Geographic distribution of individuals in cluster group *PostNeol_EuropeN*.

1728 Geographic locations of individuals within specific clusters are highlighted with colored

1729 symbols, and connected with shaded hulls (from PCAs shown in Fig. S3d.7).

1730



1731

1732 **Fig. S3f.27** PCA for cluster group *PostNeol_EuropeE*. PCA positions of individuals within

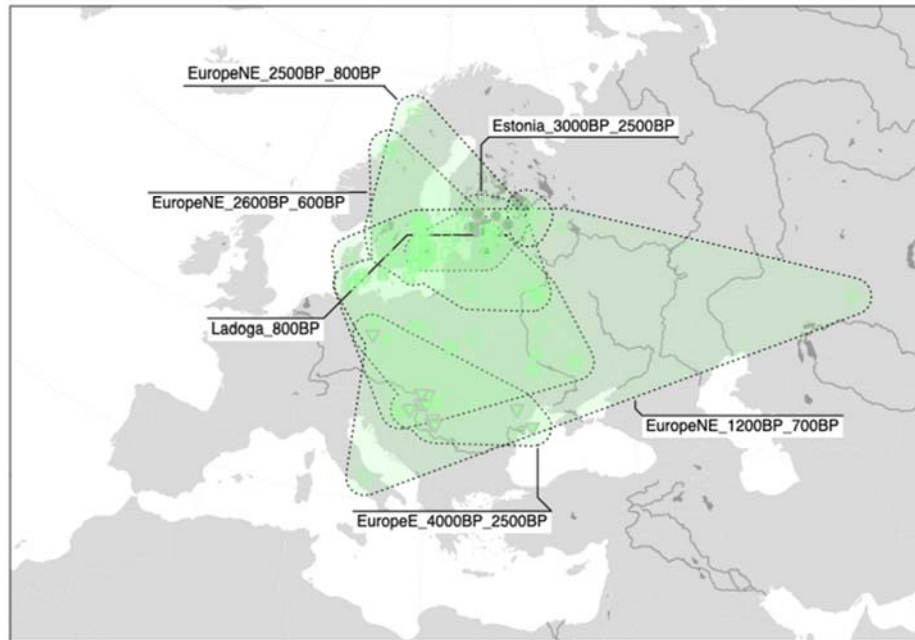
1733 specific clusters are highlighted with colored symbols, and connected with shaded hulls (from

1734 PCAs shown in Fig. S3d.7).

1735

1736

PostNeol_EuropeE



1737

1738 **Fig. S3f.28** Geographic distribution of individuals in cluster group *PostNeol_EuropeE*.

1739 Geographic locations of individuals within specific clusters are highlighted with colored

1740 symbols, and connected with shaded hulls (from PCAs shown in Fig. S3d.7).

1741

1742

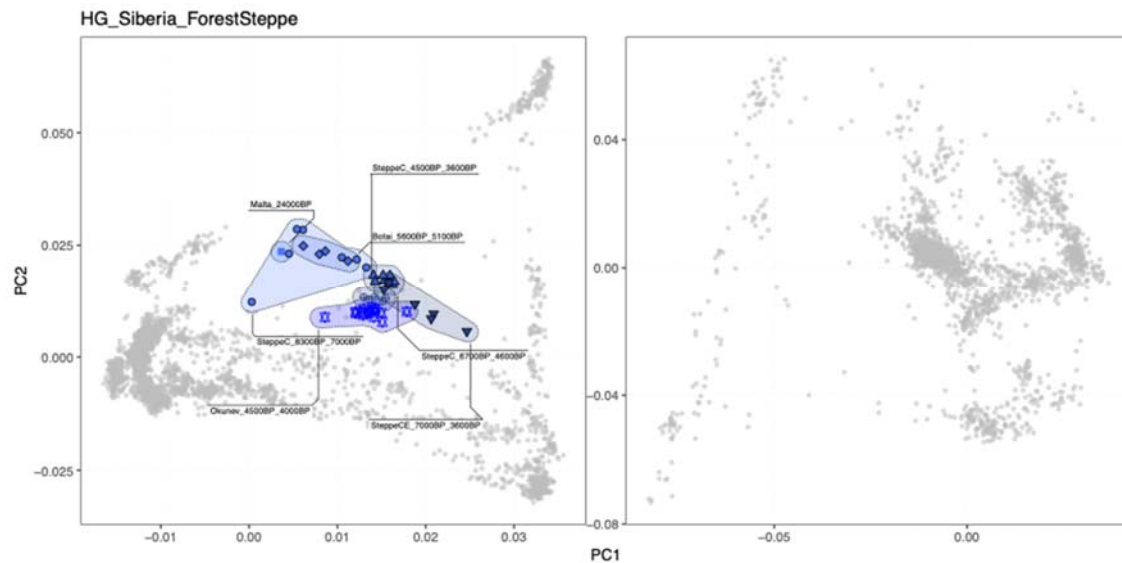
1743 Asia_45000BP_200BP

1744

1745 This global cluster includes diverse sets of clusters of individuals from the Southeast-, East-
1746 and North Asia, broadly characterised by “east Eurasian” ancestry. The individual genetic
1747 clusters are partitioned into a total of 12 cluster groups (Fig. S3f.29-44).

1748

1749

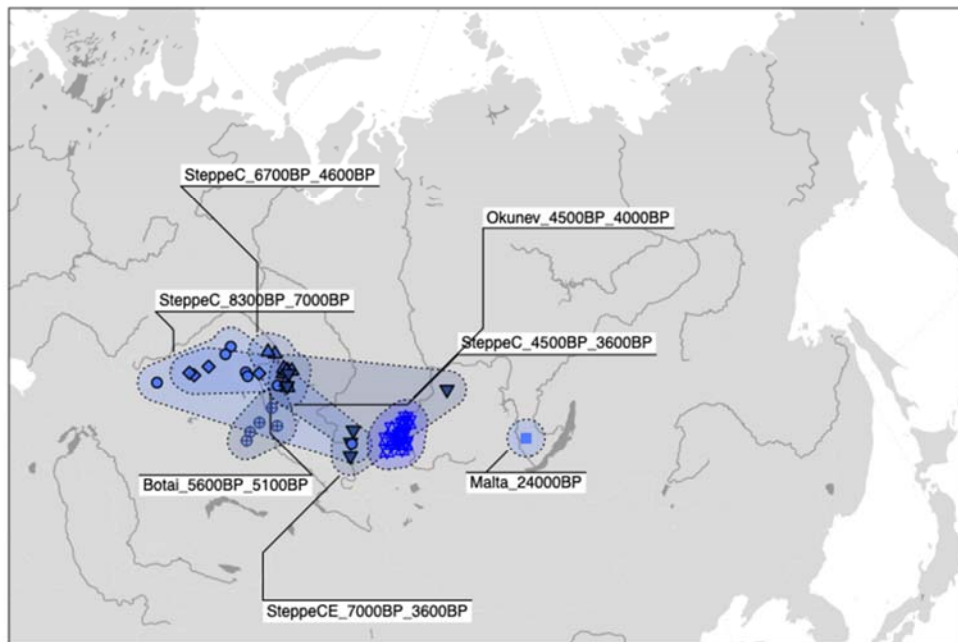


1750

1751 **Fig. S3f.29** PCA for cluster group *HG_Siberia_ForestSteppe*. PCA positions of
1752 individuals within specific clusters are highlighted with colored symbols, and connected with
1753 shaded hulls (from PCAs shown in Fig. S3d.7).

1754

HG_Siberia_ForestSteppe



1755

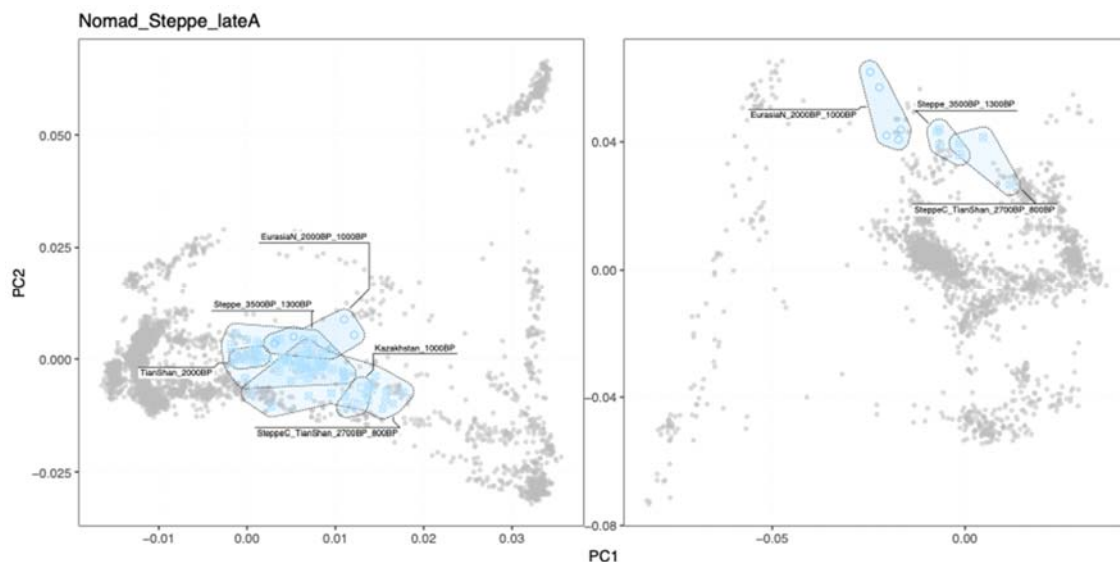
1756 **Fig. S3f.30** Geographic distribution of individuals in cluster group

1757 *HG_Siberia_ForestSteppe*. Geographic locations of individuals within specific clusters are

1758 highlighted with colored symbols, and connected with shaded hulls (from PCAs shown in Fig.

1759 S3d.7).

1760



1761

1762 **Fig. S3f.31** PCA for cluster group *Nomad_Steppe_lateA*. PCA positions of individuals

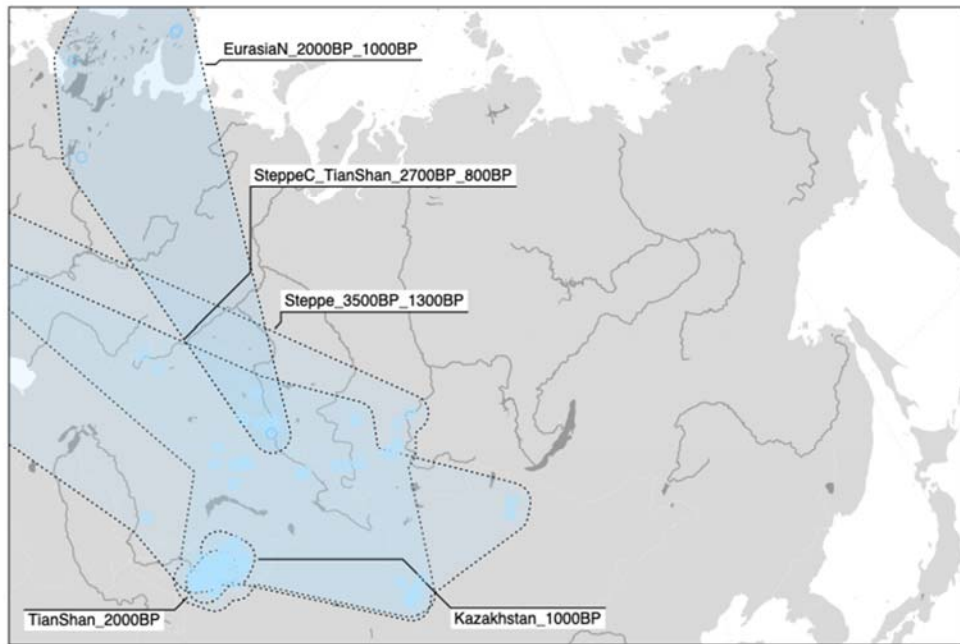
1763 within specific clusters are highlighted with colored symbols, and connected with shaded

1764 hulls (from PCAs shown in Fig. S3d.7).

1765

1766

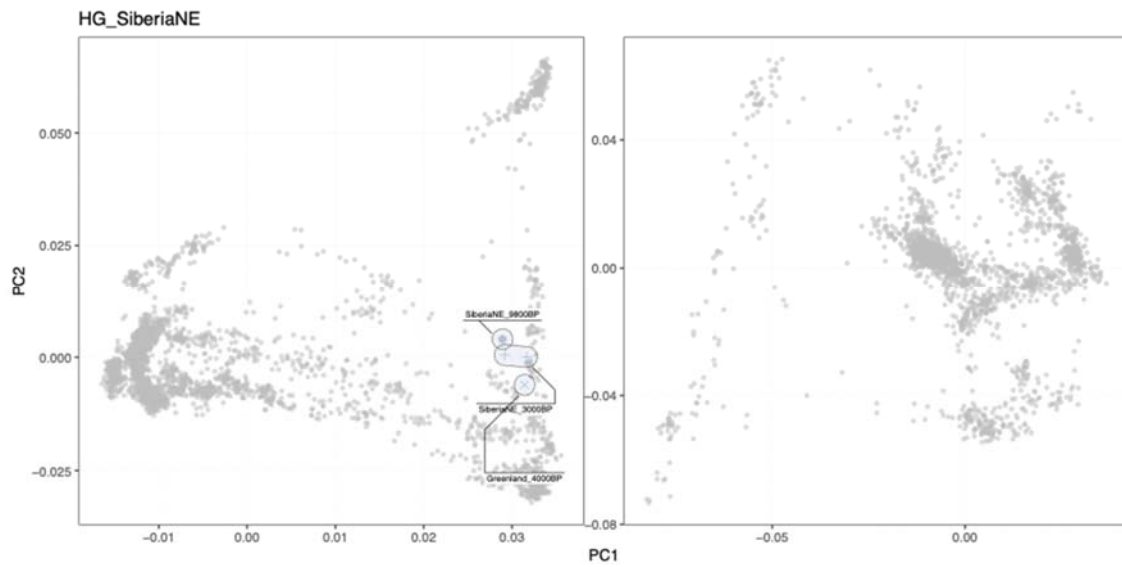
Nomad_Steppe_lateA



1767

1768 **Fig. S3f.32** Geographic distribution of individuals in cluster group

1769 *Nomad_Steppe_lateA*. Geographic locations of individuals within specific clusters are
1770 highlighted with colored symbols, and connected with shaded hulls (from PCAs shown in Fig.
1771 **S3d.7**).



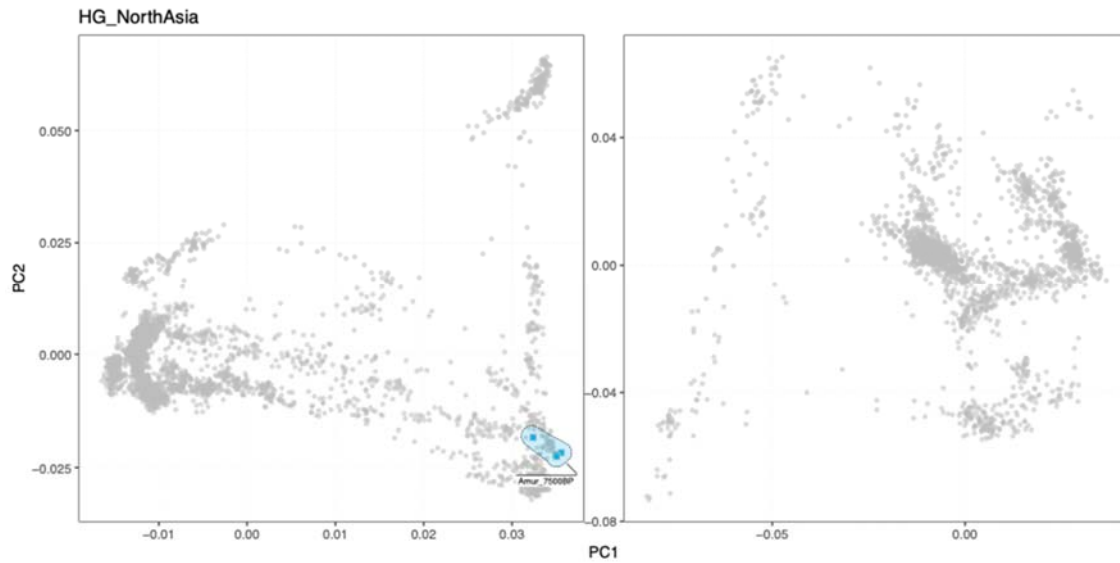
1772

1773 **Fig. S3f.33** PCA for cluster group *HG_SiberiaNE*. PCA positions of individuals within
1774 specific clusters are highlighted with colored symbols, and connected with shaded hulls (from
1775 PCAs shown in Fig. S3d.7).
1776
1777

HG_SiberiaNE



1778
1779 **Fig. S3f.34** Geographic distribution of individuals in cluster group *HG_SiberiaNE*.
1780 Geographic locations of individuals within specific clusters are highlighted with colored
1781 symbols, and connected with shaded hulls (from PCAs shown in Fig. S3d.7).
1782



1783

1784 **Fig. S3f.35** PCA for cluster group *HG_NorthAsia*. PCA positions of individuals within
 1785 specific clusters are highlighted with colored symbols, and connected with shaded hulls (from
 1786 PCAs shown in Fig. S3d.7).

1787

1788

HG_NorthAsia

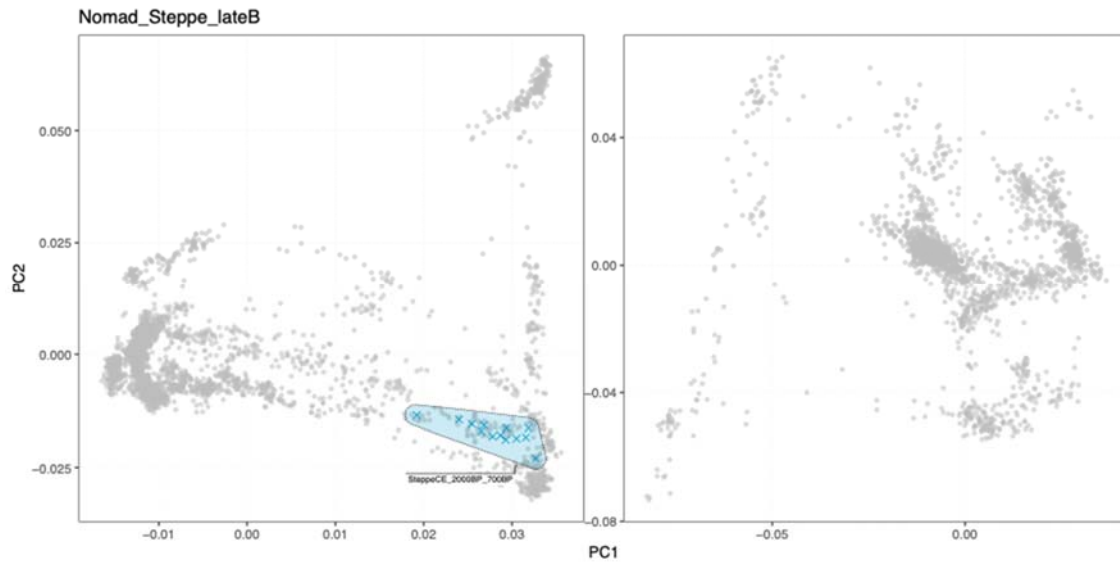


1789

1790 **Fig. S3f.36** Geographic distribution of individuals in cluster group *HG_NorthAsia*.

1791 Geographic locations of individuals within specific clusters are highlighted with colored

1792 symbols, and connected with shaded hulls (from PCAs shown in Fig. S3d.7).



1793

1794

Fig. S3f.37 PCA for cluster group *Nomad_Steppe_lateB*. PCA positions of individuals within specific clusters are highlighted with colored symbols, and connected with shaded hulls (from PCAs shown in Fig. S3d.7).

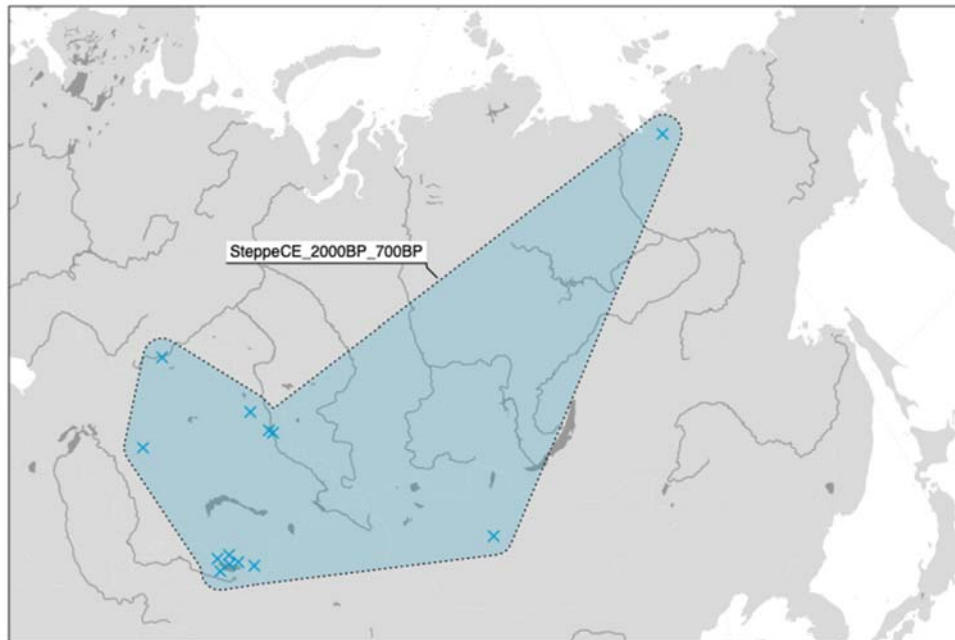
1795

1796

1797

1798

Nomad_Steppe_lateB



1799

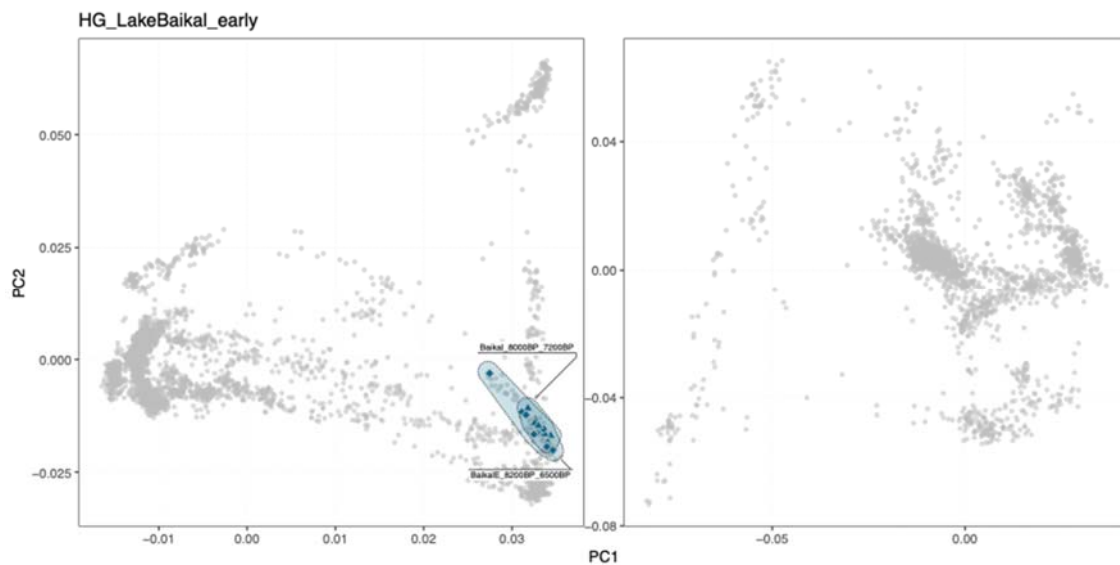
1800

Fig. S3f.38 Geographic distribution of individuals in cluster group

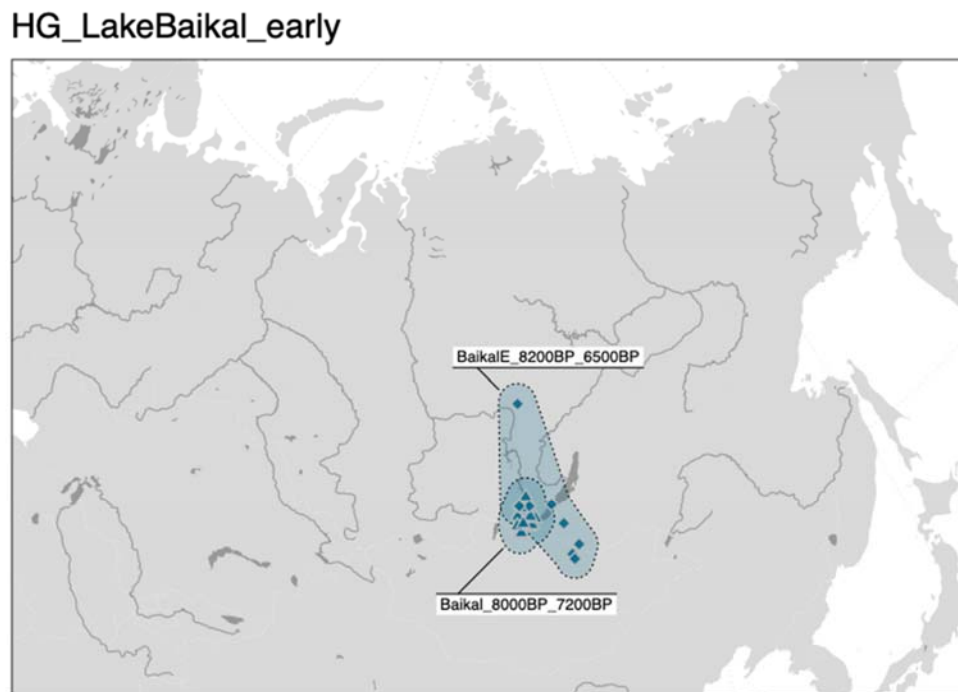
1801

Nomad_Steppe_lateB. Geographic locations of individuals within specific clusters are

1802 highlighted with colored symbols, and connected with shaded hulls (from PCAs shown in Fig.
1803 S3d.7).
1804

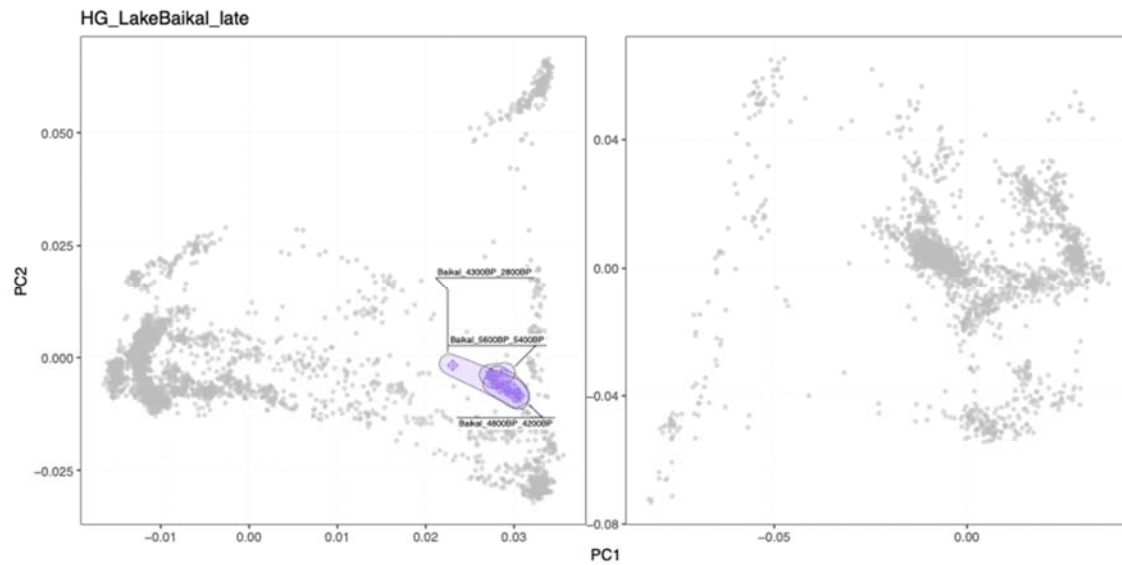


1805
1806 **Fig. S3f.39** PCA for cluster group *HG_LakeBaikal_early*. PCA positions of individuals
1807 within specific clusters are highlighted with colored symbols, and connected with shaded
1808 hulls (from PCAs shown in Fig. S3d.7).
1809
1810



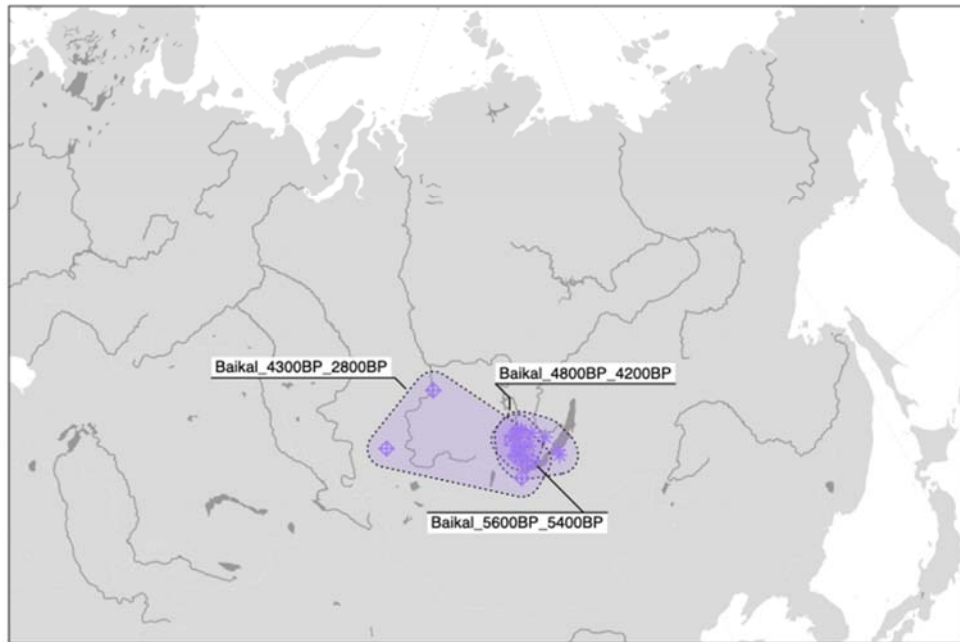
1811

1812 **Fig. S3f.40** Geographic distribution of individuals in cluster group
1813 *HG_LakeBaikal_early*. Geographic locations of individuals within specific clusters are
1814 highlighted with colored symbols, and connected with shaded hulls (from PCAs shown in Fig.
1815 S3d.7).
1816



1817
1818 **Fig. S3f.41** PCA for cluster group *HG_LakeBaikal_late*. PCA positions of individuals
1819 within specific clusters are highlighted with colored symbols, and connected with shaded
1820 hulls (from PCAs shown in Fig. S3d.7).
1821

HG_LakeBaikal_late



1822

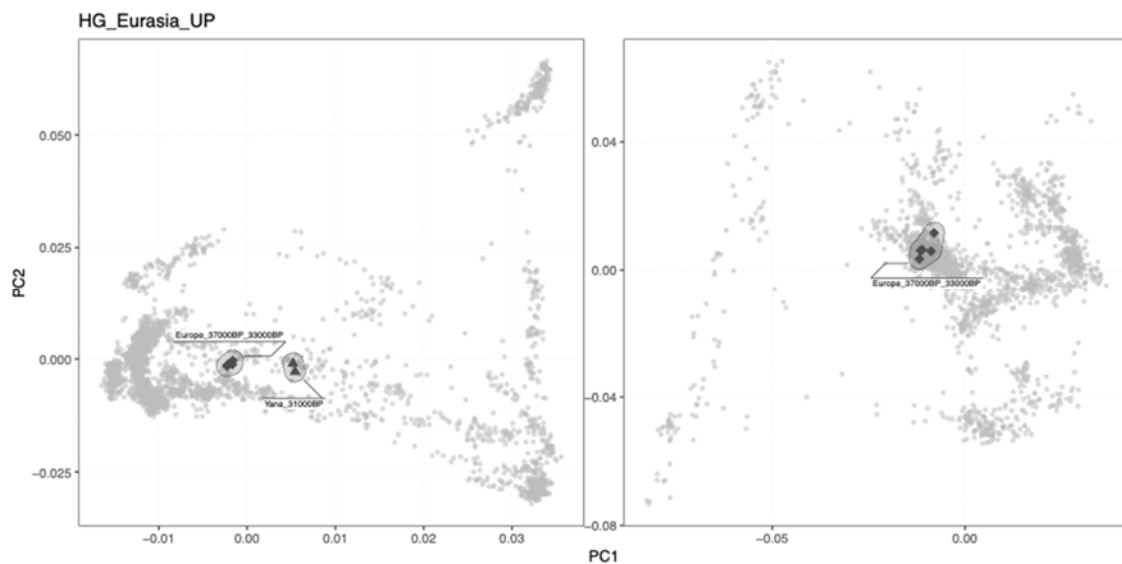
1823 **Fig. S3f.42** Geographic distribution of individuals in cluster group

1824 *HG_LakeBaikal_late*. Geographic locations of individuals within specific clusters are

1825 highlighted with colored symbols, and connected with shaded hulls (from PCAs shown in Fig.

1826 S3d.7).

1827



1828

1829 **Fig. S3f.43** PCA for cluster group *HG_Eurasia_UP*. PCA positions of individuals within

1830 specific clusters are highlighted with colored symbols, and connected with shaded hulls (from

1831 PCAs shown in Fig. S3d.7).

1832

1833

HG_Eurasia_UP



1834

1835 **Fig. S3f.44** Geographic distribution of individuals in cluster group *HG_Eurasia_UP*.

1836 Geographic locations of individuals within specific clusters are highlighted with colored

1837 symbols, and connected with shaded hulls (from PCAs shown in Fig. S3d.7).

1838

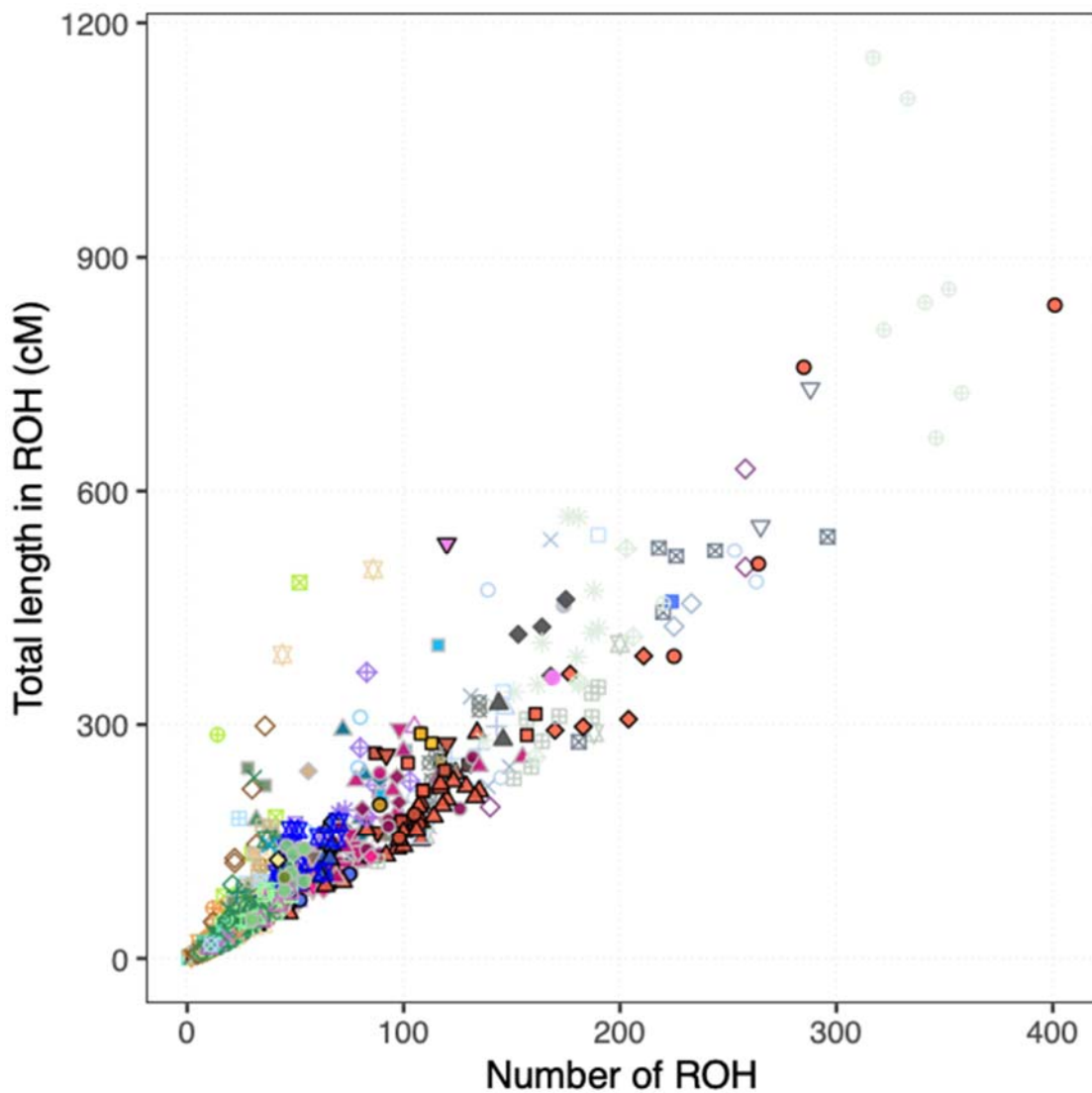
1839

1840 Runs of homozygosity and IBD sharing within clusters

1841 We quantified genetic relatedness within clusters by investigating runs of homozygosity and
1842 pairwise IBD sharing among cluster individuals. The analyses showed broad differences in
1843 patterns of genetic relatedness between different cluster groups across Eurasia, associated
1844 with both spatiotemporal and subsistence contexts of the individuals. The highest amounts of
1845 IBD sharing and ROH were generally found in clusters of individuals from hunter-gatherer
1846 contexts. Individuals of comparable age from farming contexts showed lower sharing,
1847 consistent with overall higher effective population sizes in farming communities compared to
1848 forager groups (Fig. S3f.45-48).

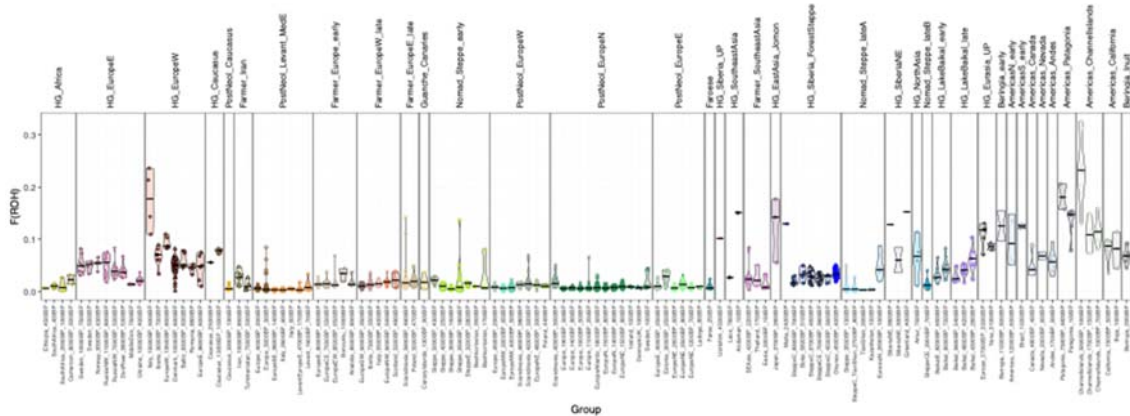
1849

1850

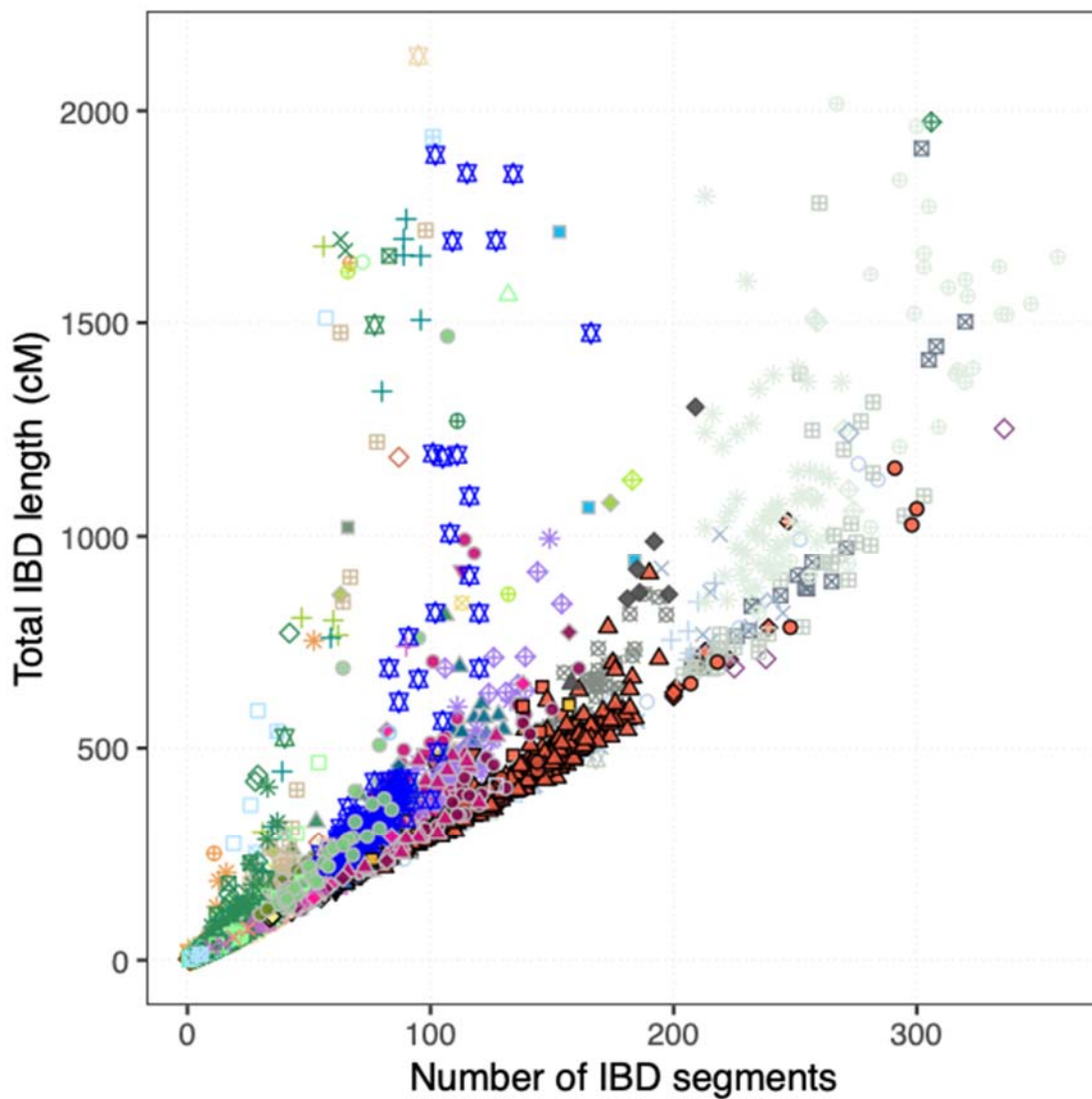


1851

1852 **Fig. S3f.45. Number and total length of ROH segments.** Plot shows the number and total
 1853 length of ROH segments detected in the respective ancient individual. Symbol colour and
 1854 shape indicated genetic cluster membership.
 1855
 1856



1857
 1858 **Fig. S3f.46. Distribution of F(ROH).** Violin plots and symbols showing the distributions of
 1859 F(ROH) within genetic clusters.
 1860



1862

1863

1864

1865

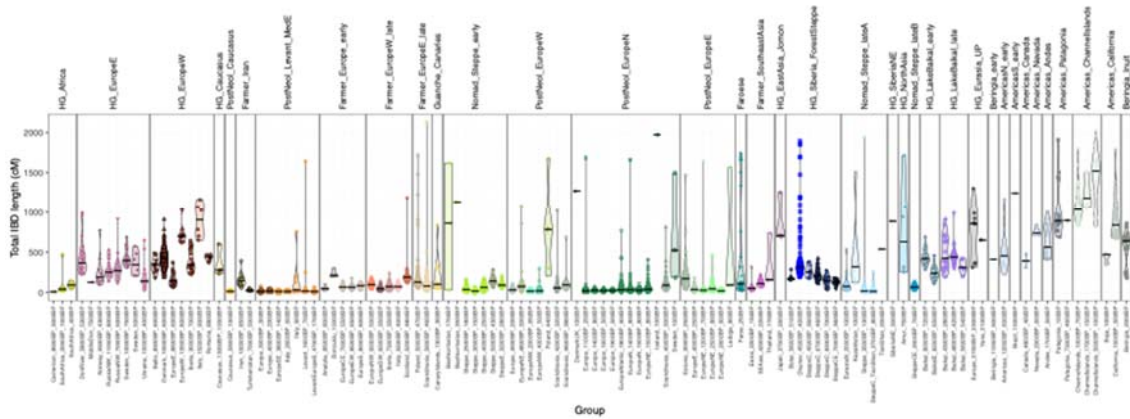
1866

1867

1868

1869

Fig. S3f.47. Number and total length of IBD segments. Plot shows the number and total length of IBD segments detected in the respective pair of individuals. Symbol colour and shape indicated genetic cluster membership.



1870

1871 **Fig. S3f.48.** Violin plots and symbols showing the distributions of total IBD sharing within
 1872 genetic clusters.

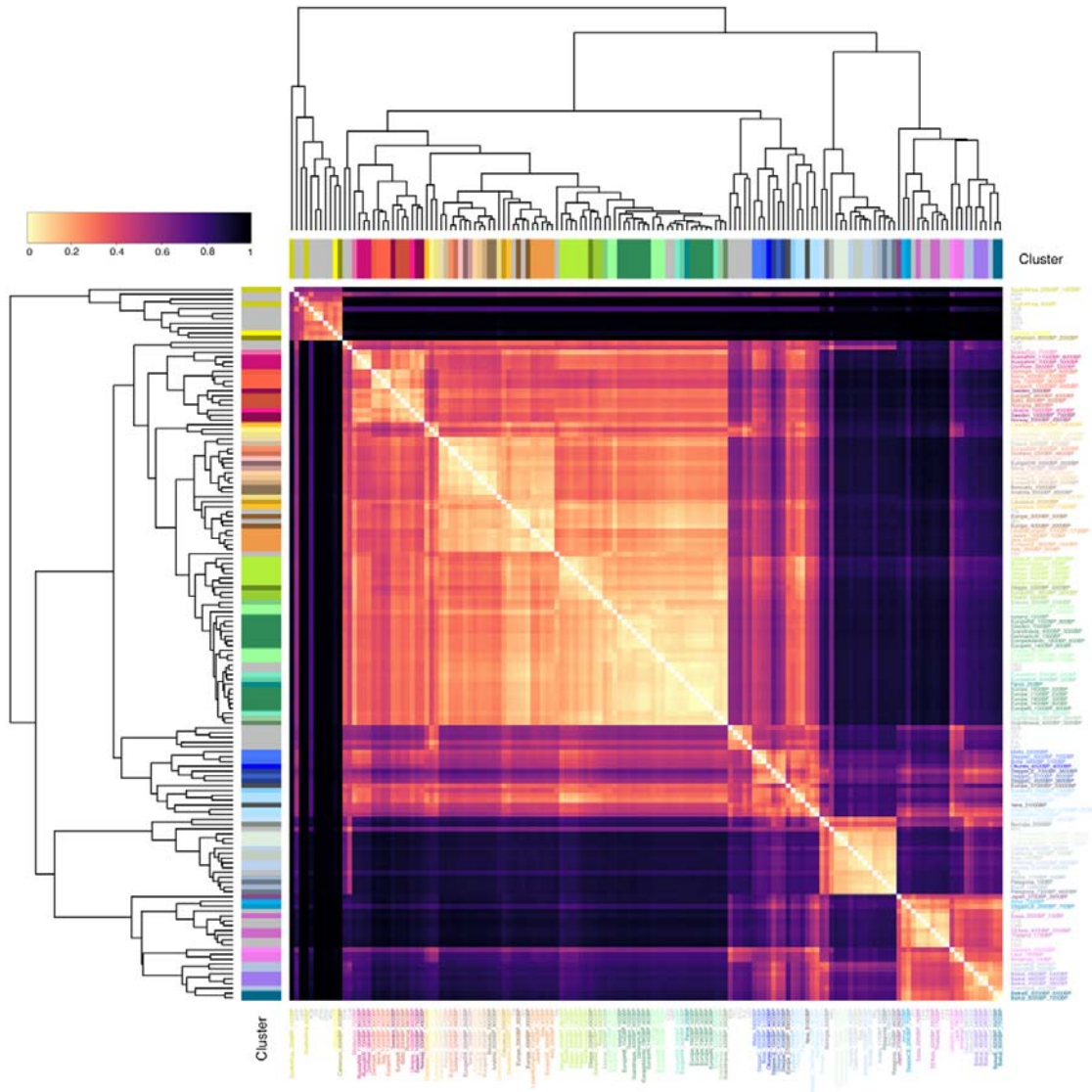
1873

1874

1875 Mixture models

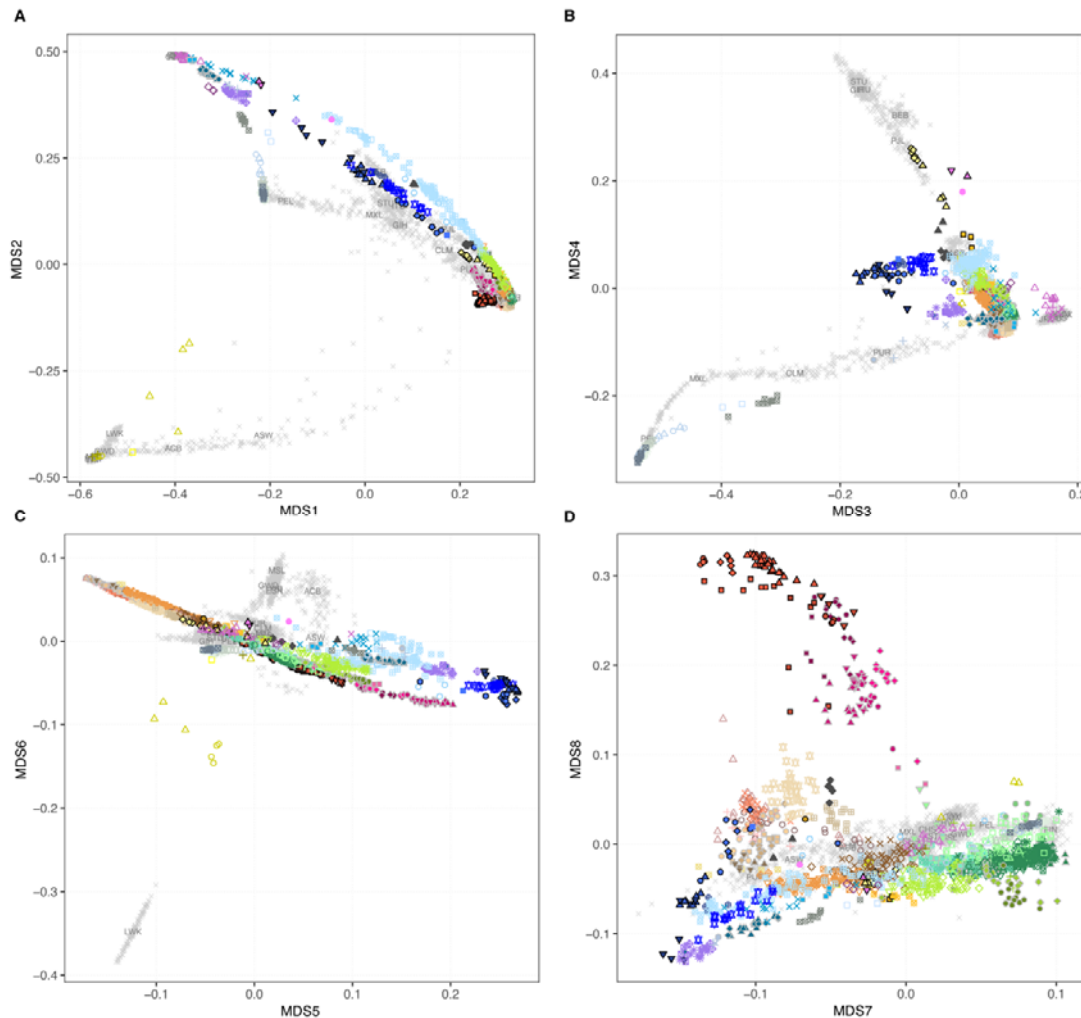
1876 We used IBD “painting profiles” to model sets of target individuals as mixtures of putative
 1877 source groups. To investigate how these IBD profiles are capturing underlying population
 1878 structure, we compared their similarities using the “total-variation- distance“ (TVD)^{12,13}
 1879 measure. We calculated pairwise TVD values for each pair of individual profiles in the
 1880 combined ancient and modern dataset, as well as for average profiles aggregated across all
 1881 individuals within a genetic cluster. Our results show that the painting profiles readily
 1882 distinguish both broad- and fine-scale genetic differentiation among the individuals and
 1883 genetic clusters (**Fig. S3f.49,50**).

1884



1885
 1886
 1887
 1888
 1889
 1890
 1891
 1892

Fig. S3f.49. Cluster IBD painting profile distances. Heatmap showing pairwise distance between genetic cluster IBD painting profiles, measured using TVD. Colored bars indicate cluster membership.



1893

1894

Fig. S3f.49. Genetic structure inferred from IBD painting profiles. (A)-(D) Plots show the first 8 dimensions of a multidimensional scaling (MDS) of individual painting profile TVDs across ancient and modern individuals. Genetic cluster membership for ancient individuals is indicated by symbol colour and shape. Present-day individuals are indicated with grey crosses, with labels indicating population median coordinates.

1899

1900

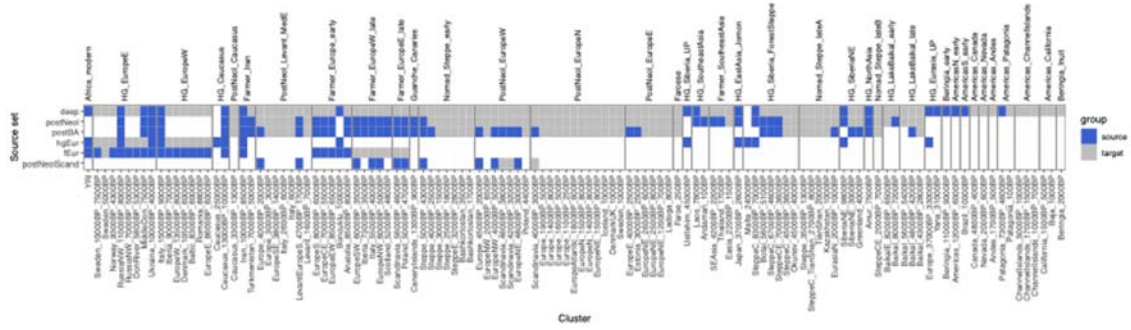
1901

We used these painting profiles in supervised modelling of target individuals as mixtures from different sets of putative source groups (Fig. S3f.51; Supplementary Table VIII). In each source set analysis, we computed source group painting profiles by averaging the profiles of the included individuals within each source group. We then estimated the mixture and proportions of source profiles that best fits the profile observed in the target individuals, using non-negative least squares (Fig. S3f.52, 53; Supplementary Tables IX-XIV)

1907

1908

1909



1910

1911

Fig. S3f.52. Mixture model source and target groups. Matrix showing the source and

1912

target groups used across the five source set analyses.

1913



1916 **Fig. S3f.52. Mixture model results.** Heatmap showing estimated ancestry proportions for
1917 target individuals (columns) from source groups (rows), across the five source set analyses.
1918

1919 References

- 1920 1. Browning, B. L. & Browning, S. R. Detecting Identity by Descent and Estimating
1921 Genotype Error Rates in Sequence Data. *Am. J. Hum. Genet.* **93**, 840–851 (2013).
1922 2. Browning, S. R. & Browning, B. L. Accurate Non-parametric Estimation of Recent
1923 Effective Population Size from Segments of Identity by Descent. *Am. J. Hum. Genet.* **97**,
1924 404–418 (2015).
1925 3. Lawrence, M. *et al.* Software for Computing and Annotating Genomic Ranges. *PLOS*
1926 *Comput. Biol.* **9**, e1003118 (2013).
1927 4. Greenbaum, G., Rubin, A., Templeton, A. R. & Rosenberg, N. A. Network-based
1928 hierarchical population structure analysis for large genomic data sets. *Genome Res.* **29**,
1929 2020–2033 (2019).
1930 5. Csardi, G. & Nepusz, T. The igraph software package for complex network research.
1931 *InterJournal Complex Systems*, 1695 (2006).
1932 6. Traag, V. A., Waltman, L. & van Eck, N. J. From Louvain to Leiden: guaranteeing
1933 well-connected communities. *Sci. Rep.* **9**, 1–12 (2019).
1934 7. Lawson, D. J., Hellenthal, G., Myers, S. & Falush, D. Inference of Population
1935 Structure using Dense Haplotype Data. *PLoS Genet* **8**, e1002453 (2012).
1936 8. Hofmanová, Z. *et al.* Early farmers from across Europe directly descended from
1937 Neolithic Aegeans. *Proc. Natl. Acad. Sci.* **113**, 6886–6891 (2016).
1938 9. Hellenthal, G. *et al.* A Genetic Atlas of Human Admixture History. *Science* **343**, 747–
1939 751 (2014).
1940 10. Soetaert, K., Meersche, K. V. den & Oevelen, D. van. *limSolve: Solving Linear*
1941 *Inverse Models.* (2009).
1942 11. Eisenmann, S. *et al.* Reconciling material cultures in archaeology with genetic data:
1943 The nomenclature of clusters emerging from archaeogenomic analysis. *Sci. Rep.* **8**,
1944 13003 (2018).
1945 12. Leslie, S. *et al.* The fine-scale genetic structure of the British population. *Nature* **519**,
1946 309–314 (2015).
1947 13. Dorp, L. van *et al.* Evidence for a Common Origin of Blacksmiths and Cultivators in
1948 the Ethiopian Ari within the Last 4500 Years: Lessons for Clustering-Based Inference.
1949 *PLOS Genet.* **11**, e1005397 (2015).
1950

1951 3g) Selecting non-British individuals from the UK Biobank

1952 Will Barrie¹ and Dan Lawson²

1953 ¹Zoology Department, University of Cambridge, UK.

1954 ²School of Mathematics and Integrative Epidemiology Unit, University of Bristol, UK.

1955

1956 Introduction

1957 The UK Biobank (UKB) contains approximately 40,000 individuals not born in the UK.
1958 Because many of these individuals are admixed or British, we set up a pipeline to (1) exclude
1959 genetically British-like individuals and (2) select individuals of a typical genetic ancestral
1960 background for each country, in order to investigate the genetic contribution of each ancient
1961 ancestry to modern European, Asian and African populations.

1962 Methods

1963 For individuals from each country in Europe, Asia and Africa but not the UK, (Data-Field
1964 1647: Country of birth (UK/elsewhere) and Data-Field 20115: Country of Birth (non-UK
1965 origin)), we ran two density-based scans using Scikit-learn's ¹ DBSCAN method (Density-
1966 Based Spatial Clustering of Applications with Noise) on a distance matrix constructed using
1967 the first 18 PCs², weighted by their Eigenvalues. This algorithm finds cores of high density
1968 within a distance matrix, which can be any shape, and can include nearby non-core points.
1969 The eps parameter can be adjusted to determine how strict the clustering is. This is
1970 preferable to using visual PC cut-offs, for which it is difficult to include higher PCs; and k-
1971 means clustering, which assumes clusters are convex and all points must be clustered.

1972

1973 In the first scan, designed to remove individuals with British-like ancestry who were born
1974 abroad, individuals from a given country were combined with 8,000 random white British
1975 individuals, and the clustering algorithm was run on the combined data. Any individuals born
1976 abroad who clustered with the white British were excluded (eps=60). For countries that are
1977 very similar to Britain in ancestry (e.g. Germany, Denmark) this is a balance between
1978 excluding individuals who are genuinely British (very common in the 'German' samples) but
1979 not biasing the samples away from British-like ancestry.

1980

1981 In the second scan, the remaining individuals were clustered, and the largest cluster was
1982 chosen to represent a typical ancestry for that country. The appropriate eps value (i.e. how
1983 strict the clustering should be) is a reflection of the genetic diversity of a country, and so was

1984 adjusted manually to reflect this (Figure S3g.1). In a minority of cases, the major cluster was
1985 not the indigenous ancestral background, and so the second-largest was chosen (for
1986 example, in Kenya the largest cluster was individuals of Indian origin). All selections were
1987 visually verified. Countries that had no obvious main cluster (usually due to low sample
1988 numbers) were excluded; any country with 3 or fewer individuals was also excluded.

1989

1990 In order to select Irish individuals (Republic of Ireland and Northern Ireland), step 1 was
1991 skipped but step 2 was run with relatively tight parameters, in both cases excluding
1992 approximately 20% of individuals.

1993

1994 In order to test the effectiveness of the pipeline at selecting individuals of a similar ancestral
1995 background, we looked at the variance in the genome-wide painting proportions for each
1996 country. Countries with high variance would indicate recent admixture.

1997

1998 Results

1999 This pipeline selected 24,511 individuals from 126 countries. These selected samples were
2000 painted using a reference/donor panel of ancient individuals (Supplementary Note S3h).

2001

2002 The countries that had high variance in ancestry proportions among individuals (and
2003 therefore likely that the DBSCAN was not effective in choosing individuals of a similar
2004 ancestral background) were Kazakhstan, Yemen, Egypt, Seychelles. Results for these
2005 countries should be interpreted with caution.

2006

2007 Discussion

2008 The UKB represents an important source of data for white British people but also for people
2009 from other countries globally. Usually, researchers restrict themselves to the white British
2010 cohort, but here we develop a method to select individuals from other countries. This
2011 transforms the UKB from a resource that is informative about British ancestry to one that can
2012 be used to make inferences about populations worldwide.

2013

2014

Country	Number of individuals	Number in wb cluster	eps value	Final number selected
Kenya	1684	277	800	110

Netherlands	491	300	230	153
Switzerland	175	19	230	143
India	4012	358	230	3107
Belgium	158	75	230	70
Singapore	502	262	230	86
Palestine	60	13	700	28
Nigeria	1159	90	800	1016
Hungary	105	0	230	79
Czech Republic	126	2	230	107
Ghana	929	35	700	848
Sri Lanka	744	54	230	620
Egypt	313	82	600	223
Japan	266	12	230	242
Hong Kong	648	107	230	448
Germany	2136	1044	230	1045
Turkey	182	5	400	160
Iran	540	16	230	469
South Africa	1364	488	700	57
Angola	56	2	700	22
Cameroon	54	5	230	44
Pakistan	1439	47	230	1332
Zimbabwe	750	252	700	254
Channel Islands	121	94	230	24
Bangladesh	246	4	230	225
Tanzania	425	84	1000	25

Sierra Leone	230	5	800	199
Portugal	320	18	230	275
Uganda	616	73	700	126
Poland	637	2	230	619
China	413	26	230	371
Cyprus	328	114	230	160
Italy	821	16	230	789
Bulgaria	71	2	230	65
Israel	87	10	300	56
France	856	103	230	690
Malta	365	170	230	135
Myanmar (Burma)	124	23	400	31
Philippines	333	8	230	310
Iraq	337	17	400	298
Finland	158	2	230	154
Libya	110	42	800	35
Norway	134	19	230	104
Russia	159	0	300	124
Nepal	161	0	230	119
Denmark	231	137	230	86
Spain	355	2	230	339
Serbia/Montenegro	56	0	230	55
Algeria	92	1	600	68
Sicily	3	0	230	2
Afghanistan	112	1	500	103

Gibraltar	77	34	230	32
Lebanon	77	13	500	50
Sudan	117	21	800	61
Morocco	93	3	800	66
Greece	131	5	230	116
Austria	196	20	230	171
Ukraine	62	1	400	54
Congo	167	11	800	146
Lithuania	72	0	230	66
Vietnam	74	0	230	69
Romania	68	0	230	61
Malawi	111	25	800	10
Gambia	42	0	800	38
Equatorial Guinea	4	0	800	2
Thailand	104	6	300	87
Indonesia	62	7	400	35
Central African Republic	42	6	800	14
Sweden	216	9	230	196
Jordan	15	1	600	9
Croatia	46	0	230	45
Ethiopia	81	8	800	57
Somalia	119	2	800	78
Zambia	246	101	800	56
Tunisia	27	2	900	14
Rwanda	25	1	1000	19

Yemen	108	38	2000	56
Burundi	26	3	1000	19
Eritrea	54	4	800	44
Syria	31	0	800	27
Luxembourg	7	0	230	6
Cambodia	8	1	1500	7
Macau (Macao)	8	0	500	6
Seychelles	46	1	2000	23
Liberia	22	3	800	15
Kuwait	31	16	800	5
Taiwan	26	1	230	23
Niger	5	0	800	2
Georgia	4	0	400	3
Iceland	19	4	230	15
Macedonia	18	0	230	17
Ivory Coast	32	0	800	29
Mongolia	6	0	400	6
Kazakhstan	13	0	4000	13
Brunei	19	4	400	8
Latvia	53	0	300	52
Bosnia and Herzegovina	41	0	230	41
Guinea	11	0	800	6
Slovenia	11	0	230	11
Azerbaijan	6	0	500	5
Slovakia	35	0	230	30

Kyrgyzstan	4	0	2000	3
Estonia	15	0	230	14
Senegal	12	0	800	7
South Korea	26	0	230	24
Togo	10	0	230	9
Armenia	5	0	600	5
Albania	12	0	230	12
British Indian Ocean Territory	10	2	500	4
Kurdistan	2	0	600	2
North Korea	6	0	230	5
Laos	3	0	500	3
Lesotho	2	0	400	2
Serbia	2	0	230	2
Republic of Kosovo	6	0	230	6
Botswana	7	1	800	4
Uzbekistan	3	0	5000	3
Kashmir	3	0	300	3
Turkmenistan	1	0	230	1
Belarus	6	0	230	4
Tibet	1	0	230	1
Crete	1	0	230	1
Moldova	1	0	230	1
Tajikistan	1	0	230	1

2015 **Table S3g.1. Parameters for selection of individuals in the UKB born in a given**
2016 **country of a 'typical ancestral background'**. This shows the initial number of individuals
2017 coded as being born in a country; the number removed because they clustered with the
2018 white British; the eps value for selecting the main cluster; and the number of individuals in

2019 the final selected cluster. Countries with fewer than 3 individuals in the final cluster, or no
2020 obvious main cluster, were discarded. The eps value is dependent on the genetic diversity of
2021 the population being selected, and was chosen manually and visually checked. In most but
2022 not all cases the largest cluster was chosen.

2023 References

2024 1. Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *The Journal of Machine*
2025 *Learning Research* **12**, 2825–2830 (2011).

2026 2. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.
2027 *Nature* **562**, 203–209 (2018).

2028

2029 3h) Painting the UK BioBank

2030

2031 Will Barrie¹ and Dan Lawson²

2032 ¹Zoology Department, University of Cambridge, UK.

2033 ²School of Mathematics and Integrative Epidemiology Unit, University of Bristol, UK.

2034

2035 Introduction

2036 Here, we develop new methods to use Chromopainter ¹ on a biobank scale to ‘paint’ modern
2037 genomes from the UK Biobank (UKB) using ancient genomes, grouped into reference
2038 populations, as donors. Painting was done following the pipeline of Margaryan *et al.* ² based
2039 on GLOBETROTTER ³, and admixture proportions were estimated using Non-Negative
2040 Least squares. These results (technically most recent coalescences) are used as a proxy for
2041 ancestry. We store both genome-wide and local ancestry (i.e. per variant per individual)
2042 results.

2043

2044 Methods

2045 Painting pipeline introduction

2046 The process of painting consists of forming a reference/donor panel consisting of ancient
2047 individuals of as pure ancestry as possible, having undergone QC and clustering using

2048 fineSTRUCTURE. The target/recipient panel and reference/donor panel are then filtered for
2049 variants, merged, and the target panel is painted using the reference panel as donors.
2050

2051 Reference/donor panel formation

2052 We used imputed best guess haplotypes filtered for imputation information score
2053 (FORMAT/INFO) above 0.5. Samples were selected based on IBD-sharing, visual PCA
2054 inspection, and fineSTRUCTURE analysis (unsupervised clustering based on the coancestry
2055 matrix output of ChromoPainter; [Figure S3h.1](#)); low coverage, contaminated, and related
2056 individuals were excluded. The aim was to group samples into as pure 'source' populations
2057 as possible, while maintaining reasonable numbers in each population. We do not expect our
2058 filters for white British/non-british individuals to be perfect; furthermore, modelling modern
2059 Eurasians as a mixture of hunter-gatherer/Steppe/farmer is overly simplistic. Therefore, we
2060 also include ancient African and EastAsian reference populations to account for possible
2061 'non-European' ancestry.

2062

2063 Ultimately, 318 individuals split into ten reference populations were used ([Figure S3h.2](#),
2064 [Figure S3h.3](#), [Table S3h.1](#)): western hunter-gatherer (WHG), eastern hunter-gatherer (EHG),
2065 Caucasus hunter-gatherer (CHG), FarmerAnatolian, FarmerEarly, FarmerMiddle,
2066 FarmerLate, Yamnaya, African and EastAsian. Populations are characterised by
2067 preferentially copying from individuals within the population, as well as being biologically and
2068 historically meaningful. This dataset is henceforth called the "present aDNA dataset".

2069

2070 The farmers are split into four separate populations due to their differing behaviour as donors
2071 (columns) in the fineSTRUCTURE analysis ([Figure S3h.1](#)). There is a cline in their degree of
2072 WHG admixture that roughly correlates with age, while some samples also show Steppe
2073 admixture. Given the nature of the splits, the differences between these groups should be
2074 interpreted with caution, and for most downstream analysis these groups are merged into a
2075 'Farmer' ancestry.

2076

2077 Target/recipient panel formation

2078 We used white British individuals from the UKB as reported in Bycroft et al.⁴; these are
2079 individuals who self-reported as white British and have British-like ancestry according to
2080 PCA. We also used individuals from the UKB of a typical ancestral background selected by
2081 country of origin ([Supplementary Note S3g](#)). We used phased haplotype data, downloaded

2082 from <https://www.ukbiobank.ac.uk>. This totalled 408,884 white British individuals, and 24,511
2083 non-British individuals. This dataset is henceforth called the “UKB dataset”.
2084

2085 SNP selection and merging of the panels

2086 Due to computational considerations, the number of SNPs used in the painting was limited to
2087 those in the UKB Axiom Array; these SNPs were chosen to capture genome-wide variation,
2088 rare and coding variants, and variants relevant to specific phenotypes or regions of interest ⁴.
2089 The present aDNA dataset and UKB datasets were merged and filtered for these variants
2090 using QCTOOL v2 (https://www.well.ox.ac.uk/~gav/qctool_v2/), and then filtered to exclude
2091 variants with a minor allele frequency below 1% using bcftools
2092 (<http://samtools.github.io/bcftools/>), leaving a total of 549,323 SNPs across chromosomes 1-
2093 22.
2094

2095 Painting process

2096 ChromoPainter ¹ uses an approach premised on the observation that markers on the same
2097 chromosome are inherited together unless separated by recombination; at the population
2098 level, this results in linkage disequilibrium (LD) between close markers that reflect a shared
2099 history of descent. The haplotype-based algorithm of ChromoPainter aims to harness this
2100 information, detecting shared haplotypes to reconstruct phased recipient genomes as chunks
2101 ‘copied’ from donors.

2102

2103 Considering the genealogy of a single locus, we can identify one or more closest relatives to
2104 that locus, henceforth called ‘nearest neighbours’; if viewed as a genealogy, these are the
2105 other leaves of the tree underneath the first coalescence. Therefore at each locus of each
2106 haplotype, there exists one or more nearest neighbours. ChromoPainter aims to identify
2107 these using an approximate method based on that introduced by Li and Stephens ⁵: the
2108 Hidden Markov Model (HMM), which explicitly reconstructs the haplotype of a recipient/target
2109 individual as a series of chunks of genetic material donated by the other donor/reference
2110 individuals, using information on the types of the recipient and potential donor at each SNP.
2111 This approach is probabilistic, calculating the expectations of which haplotype acts as donor
2112 to a recipient as a function of position over an infinite number of paintings ¹. Although
2113 ChromoPainter was originally intended to use this information, in the form of a ‘co-ancestry
2114 matrix’, to ascertain fine-scale population structure and clustering (in the fineSTRUCTURE
2115 software package), the software can be used with pre-defined donor and recipient
2116 populations.

2117

2118 If the donor panel is formed of ancient individuals and the recipient individual is modern, the
2119 nearest neighbour should reflect some history of that locus. The ability of chromosome
2120 painting to accurately infer ancestry is expected to depend on the diversity between donor
2121 populations: more genetically similar populations and the algorithm will find it difficult to
2122 correctly identify the nearest neighbour(s). There is also the issue of 'masking' whereby
2123 haplotypes from older populations would have travelled through more recent populations
2124 before arriving in the modern population; this causes a genome-wide bias towards the more
2125 recent ancient populations, the effects of which are discussed below. Here, we use nearest
2126 neighbour as a proxy for local ancestry - i.e. which population that haplotype came from
2127 (which may not be a single unique population from our panel).

2128

2129 Chromosome painting cannot include the target of painting. Therefore, painting was done
2130 (following the pipeline of Margaryan et al. ² based on GLOBETROTTER ³ by leaving out one
2131 individual at random (chosen independently for each chromosome) from each other donor
2132 population for all donor individuals. Target individuals from the UK Biobank were painted by
2133 similarly removing one individual at random from all donor populations. This ensures that
2134 individuals from the reference and UK Biobank are exchangeable.

2135

2136 Once we had a well-chosen set of ancient populations from the present aDNA panel, each
2137 individual was repainted twice leaving out themselves as a possible donor: first to learn the
2138 painting parameters N_e and μ , and then to learn a genome-wide individual-specific donor-
2139 prior. For each of the reference populations, the average amount of genome received from
2140 each donor individual was learnt. We then painted the modern individuals in the UKB panel
2141 using the reference populations and the learnt parameters and priors.

2142

2143 The probability that each recipient copied each donor population at every SNP was recorded.
2144 The genome-wide information for each recipient was also stored, in the form of (i)
2145 chunkcounts, the number of chunks copied from each donor population and (ii) chunk
2146 lengths, the sum of the lengths of the chunks copied from each population, weighted by their
2147 copying probability. Admixture proportions were then estimated using Non-Negative Least
2148 Squares (NNLS).

2149

2150 Painting at biobank scale

2151 We used custom scripts to speed up this process (specifically reading from large phase
2152 files), to enable running for large numbers of recipients in parallel across multiple nodes, and

2153 to store the local copying probabilities in a memory-efficient format in real time (all scripts
2154 available at https://github.com/will-camb/Nero/tree/master/scripts/cp_panel_scripts). The
2155 total CPU time for painting the UKB panel was approximately 550,000 CPU hours.

2156 Results

2157 Ancestry-PCs relationship

2158 PCA is a dimensionality reduction technique that can be applied to genetic data, the results
2159 of which are useful as a means to visualise variation between individuals/groups, and are
2160 expected to reflect historical events that cause differences in ancestry due to drift, admixture
2161 etc. It is well established that PC1 vs PC2 vs PC3 generally separate African, European and
2162 East Asian populations. We ran multivariate linear regressions using ancestry components to
2163 predict UKB PCs ⁴. Previous work has shown that the main UKB PCs that reflect British
2164 population structure are PCs 5 and 9, describing variation between English, Scottish and
2165 Welsh ancestry, and PCs 11 and 14 which further separate structure within Wales and
2166 England ⁶.

2167

2168 We found significant correlations between ancestry components and PC4 (R-
2169 squared=0.553) and PC5 (R-squared=0.376), as well as PC1 (R-squared=0.165) and PC7
2170 (R-squared=0.130). We found that the high PC4 correlation with ancestry component was
2171 largely driven by a Steppe (Yamnaya/EHG) vs Farmer divide, both within Britain and
2172 internationally: high PC4 values are associated with high Steppe/low Farmer ancestry, while
2173 low PC4 values are associated with low Steppe/high Farmer ancestry.

2174

2175 Ancestry-geographic variation

2176 Within the British Isles, all individuals were painted with similar proportions from each
2177 reference population, as expected when measuring coalescence tracts rather than direct
2178 admixture tracts and after a long time since admixture events; but, the differences in copying
2179 proportions showed significant geographic heterogeneity. We ran multivariate linear
2180 regressions, using longitude and latitude of place of birth ("Place of birth in UK – east co-
2181 ordinate" and "Place of birth in UK – north co-ordinate") to predict NNLS ancestry fractions.
2182 We found significant correlations for Yamnaya ancestry (R-squared=0.081), Farmer ancestry
2183 (R-squared=0.066), CHG ancestry (R-squared=0.015), WHG ancestry (R-squared=0.007),
2184 African ancestry (R-squared=0.011) and EHG ancestry (R-squared=0.01, longitude only). To
2185 visualise this, we assigned individuals to a county based on their UKB place of birth data,

2186 and plotted the average admixture proportion per county for each ancestry, binned in ten
2187 equal interval quantiles using ArcGIS Online (www.arcgis.com; Figure 5, main text).

2188

2189 We found that Neolithic farmer ancestry was highest in southern and eastern England and
2190 lower in populations in Scotland, Wales and Cornwall. We found the opposite pattern in
2191 Yamnaya ancestry, representing the Steppe component, which has previously been shown
2192 to be higher in Scotland but not Wales ⁷; we found this was highest in the Outer Hebrides.
2193 This Farmer/Yamnaya dichotomy broadly reflects an Anglo-Saxon/'Celtic' distribution. We
2194 are unable to date when these subtle population structures arose, but note that the Neolithic
2195 Anatolian-related farmer ancestry is already present in the British and Roman Iron Age but
2196 lower in Saxon individuals (Extended Data Figure 3, main text), meaning these patterns
2197 cannot be explained just by Saxon-related ancestry. They are likely a result of pre-Roman
2198 migration between 1000 and 875 BC which resulted in a slight increase in Early Farmer
2199 ancestry in England and Wales but not Scotland ⁸, although we note our results show a
2200 marked difference between Wales/Cornwall and England too. We also found higher levels of
2201 WHG-related ancestry in central and Northern England.

2202

2203 Looking at a continent-wide level, the hunter-gatherer ancestries display distinct structure in
2204 modern populations (Figure 5, main text). WHG-related ancestry is highest in present-day
2205 individuals from the Baltic States, Belarus, Poland and Russia; EHG-related ancestry is
2206 highest in Mongolia, Finland, Estonia and Central Asia; and CHG-related ancestry is
2207 maximised in countries east of the Caucasus, in Pakistan, India, Afghanistan and Iran, in
2208 accordance with previous results ⁹. The CHG-related ancestry likely picks up both Caucasus
2209 hunter-gatherer and Iranian Neolithic signals, explaining the relatively high levels in south
2210 Asia ¹⁰. Consistent with expectations ^{11,12}, Neolithic Anatolian-related farmer ancestry is
2211 concentrated around the Mediterranean basin, with high levels in southern Europe, the Near
2212 East and North Africa, including the Horn of Africa but less in Northern Europe. A contrasting
2213 pattern was observed in Yamnaya-related ancestry decreasing from high levels in northern
2214 Europe, peaking in Ireland, Iceland, Norway and Sweden, but decreasing further south
2215 where Neolithic farmer ancestry still dominates. There is also evidence for its spread into
2216 southern Asia. These results provide a new level of detail on the modern distribution of
2217 ancient ancestries.

2218

2219 To better understand how countries varied in their ancestry proportions, we ran Scikit-learn's
2220 PCA ¹³ on the average admixture proportions per country. We then ran a hierarchical
2221 clustering algorithm on the first 4 PCs (explained variance=0.244), and built a dendrogram

2222 (Figure S3h.4)¹⁴. For further analysis, we excluded countries in clusters dominated by
2223 African or East Asian ancestry, leaving 80 countries.

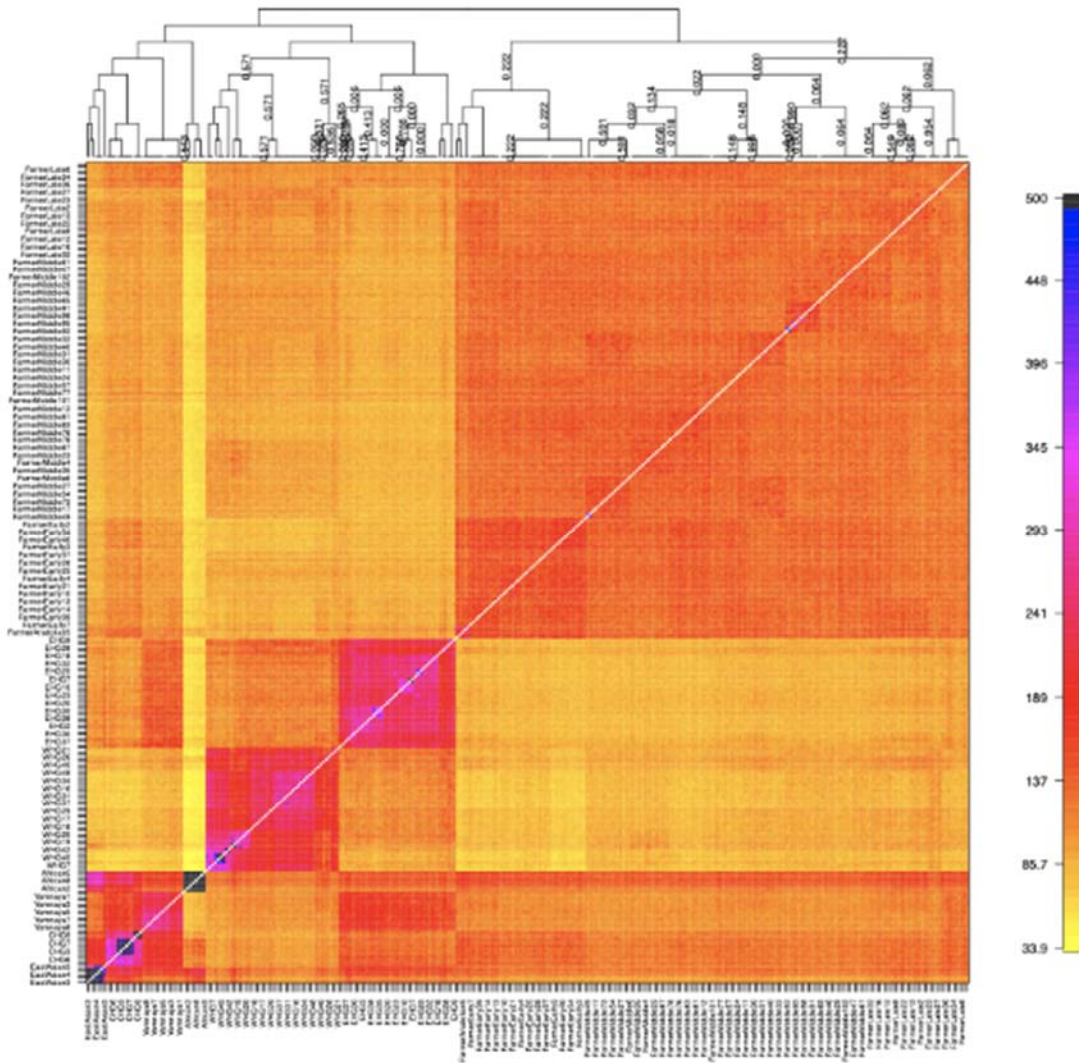
2224

2225 We used Sklearn's StandardScaler utility class to standardise each feature (zero mean and
2226 unit variance), then ran a PCA using the standardised average admixture proportions for
2227 each of the 80 remaining countries (<https://github.com/erdogant/distfit>), and plotted a biplot
2228 (PC1 vs PC2 with loadings for each feature plotted), which shows the correlations between
2229 ancestries (Figure S3h.5). This gives a more visual representation of how countries group by
2230 ancestry, and broadly reflects actual geography.

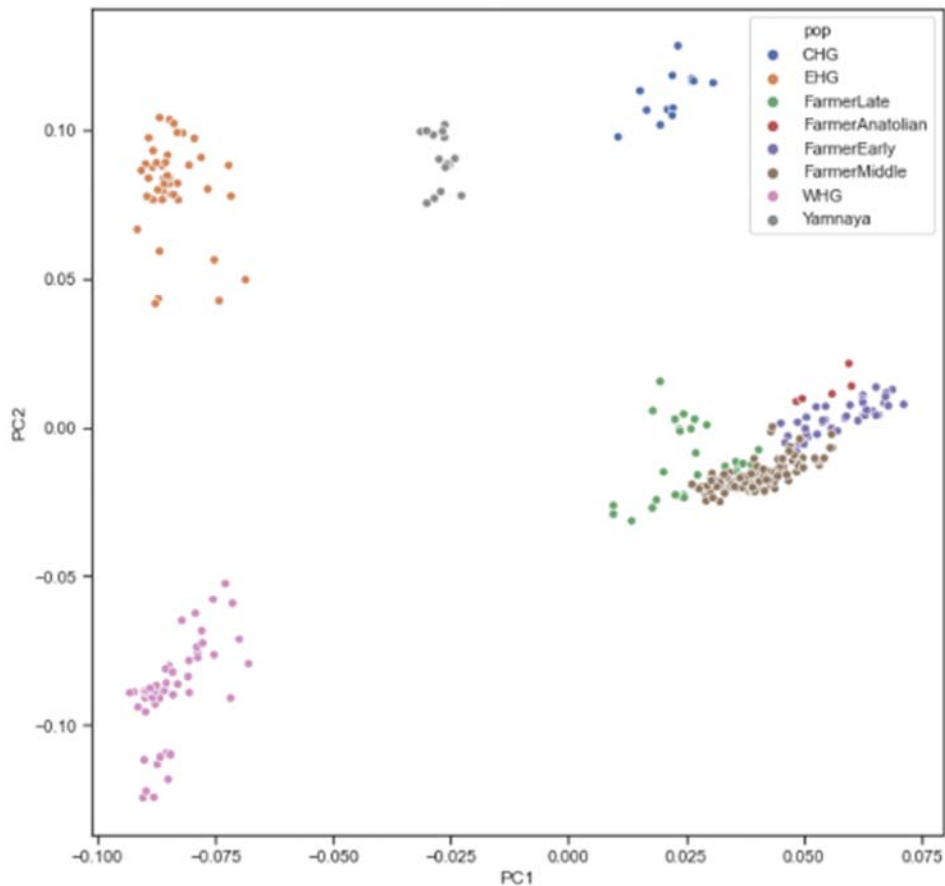
2231 Discussion

2232 The large ancient DNA panel established here combined with the UKB allows us to trace for
2233 the first time the fine-scale distribution of Mesolithic/Neolithic/Bronze Age ancestry
2234 components in modern British individuals, using DNA directly from ancient individuals. It also
2235 demonstrates the ancestry differences within an 'ethnic group' (white British) traditionally
2236 regarded as being relatively homogenous, highlighting the need for care over ancestry
2237 considerations when using resources like the UK Biobank.

2238



2240
 2241 **Figure S3h.1. Co-ancestry heatmap of selected ancient samples.** The output of
 2242 fineSTRUCTURE analysis of the ancient reference panel, showing copying proportions
 2243 between ancient populations (columns=donors, rows=recipients). There is a cline in Hunter-
 2244 Gatherer admixture in the Farmers, roughly correlating with age. For most downstream
 2245 analyses, the Farmer populations were merged.



2246
 2247
 2248
 2249
 2250
 2251
 2252
 2253
 2254
 2255

Figure S3h.2. PCA of ancient reference samples, coloured by assigned population. PC1 vs PC2 of a PCA of the ancient western Eurasian samples (excluding African and EastAsian), coloured by their assigned population used in the painting. As can be seen, populations are fairly distinct, with intermediate admixed individuals having been excluded. Some Farmers are admixed with Steppe and Hunter-Gatherer populations to differing degrees, but particularly among later individuals.



2256



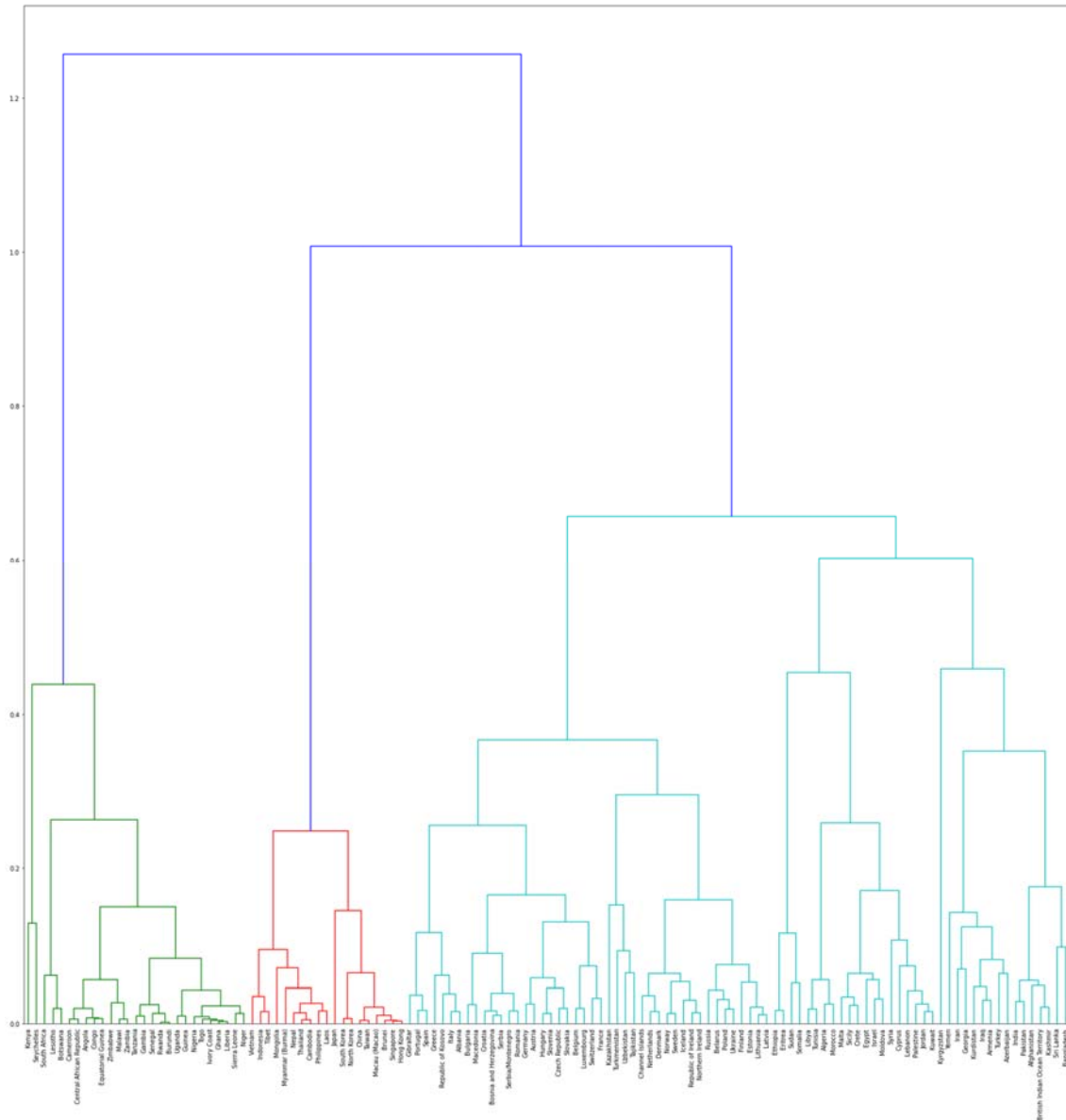
2257

2258

2259

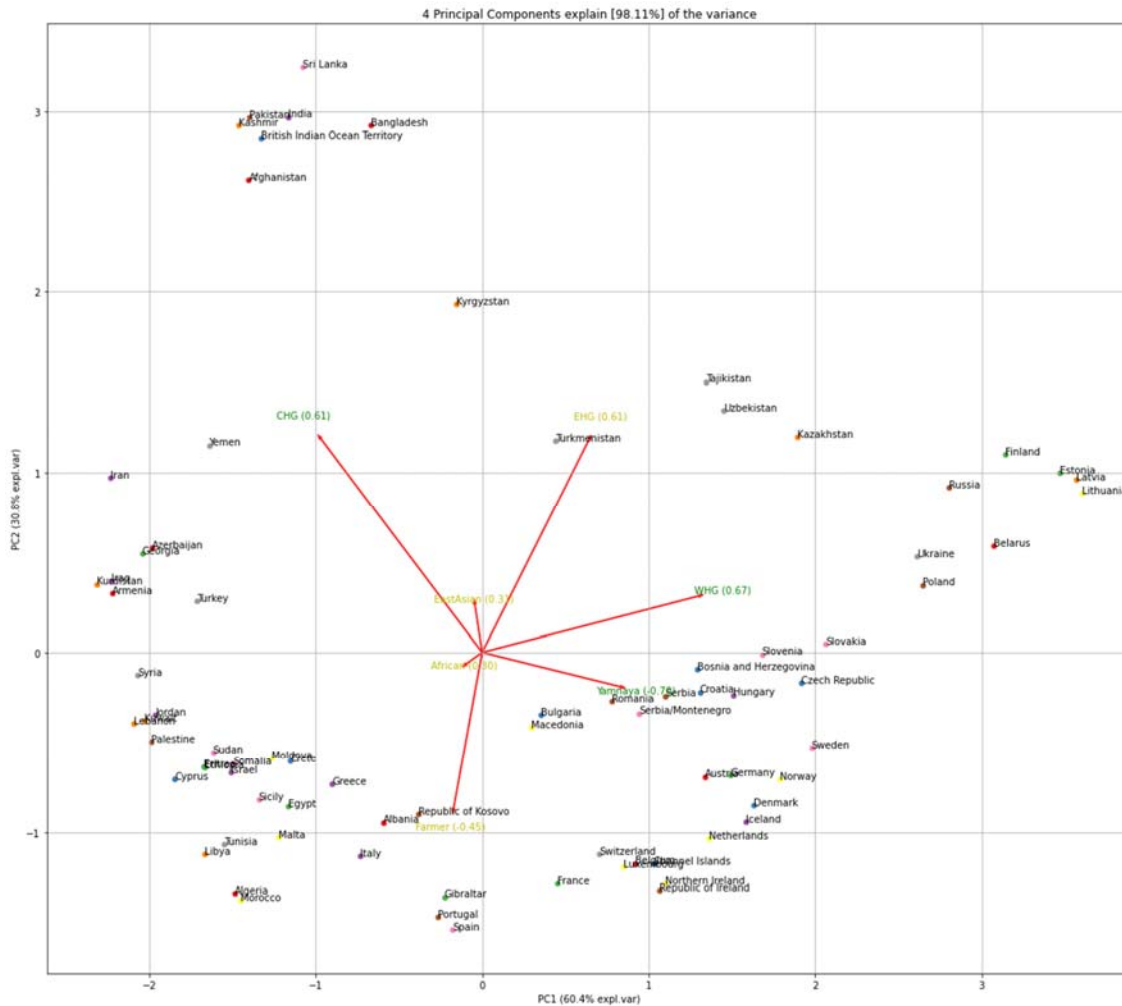
2260

Figure S3h.3. Maps of ancient sample locations coloured by assigned reference population (above) and age (below). Not showing African and East Asian samples.



2261
 2262
 2263
 2264
 2265

Figure S3h.4. Dendrogram based on hierarchical clustering of first 4 PCs of average admixture proportions per country. For further analysis, countries in clusters dominated by African (green) or East Asian (red) ancestry were dropped.



2266
2267
2268
2269
2270
2271

Figure S3h.5. PCA biplot of standardised average NLS admixture proportion per country, based on 80 countries in Europe, West/Southern Asia, the Middle East and North Africa.

referenceP population	sampleId	country	groupLabel	lati tude	long itud e	ageAv erage	coverage	s e x	clusterIBDFine
African	mfo01	South Africa	SouthAfrica_IronAge	- 28. 73	30.8 1	378.0	7.061854822 000000	X X	6.1_SouthAfrica_400B P
African	bab01	South Africa	SouthAfrica_Neolithi c	- 29. 54	31.2 2	2040.5	1.29865429	X Y	6.2_SouthAfrica_2000 BP_1000BP
African	ela01	South Africa	SouthAfrica_IronAge	- 28. 92	29.1 3	493.0	13.46552376	X X	6.1_SouthAfrica_400B P
African	110871	Camero on	Cameroon_Neolithic	5.8 6	10.0 8	7885.0	15.21262736	X Y	6.3_Cameroon_8000B P_3000BP
African	110873	Camero on	Cameroon_Neolithic	5.8 6	10.0 8	3065.0	3.276739876	X Y	6.3_Cameroon_8000B P_3000BP

African	I9133	South Africa	SouthAfrica_Neolithic	- 31.98	18.52	1970.0	2.0755013510000000	X Y	6.2_SouthAfrica_2000BP_1000BP
African	baa01	South Africa	SouthAfrica_Neolithic	- 29.54	31.22	1908.5	13.50021278	X Y	6.2_SouthAfrica_2000BP_1000BP
African	new01	South Africa	SouthAfrica_IronAge	- 27.76	29.92	417.5	10.89613659	X X	6.1_SouthAfrica_400BP
CHG	WC1	Iran	Iran_Neolithic	34.61	47.11	9218.5	10.4300754	X Y	2.1_Iran_10000BP_8500BP
CHG	AH4	Iran	Iran_Neolithic	34.19	48.37	9929.5	0.8674508960000000	X X	2.1_Iran_10000BP_8500BP
CHG	AH2	Iran	Iran_Neolithic	34.19	48.37	9930.5	0.6492785520000000	X Y	2.1_Iran_10000BP_8500BP
CHG	AH1	Iran	Iran_Neolithic	34.19	48.37	9900.0	1.161365185	X X	2.1_Iran_10000BP_8500BP
CHG	DA380	Turkmenistan	Turkmenistan_Neolithic_Namazga	37.6	59.33	5180.5	0.495449798	X X	2.1_Turkmenistan_7000BP_5000BP
CHG	DA381	Turkmenistan	Turkmenistan_Neolithic_Namazga	37.19	61.03	5181.0	0.83822953	X Y	2.1_Turkmenistan_7000BP_5000BP
CHG	NEO816	Iran	Iran_Neolithic	33.76	47.1	8700.0	0.940243388	X Y	2.1_Iran_10000BP_8500BP
CHG	NEO281	Georgia	Georgia_Mesolithic	42.22	43.32	9724.0	3.607878115	X Y	2.1_Caucasus_13000BP_10000BP
CHG	KK1	Georgia	Georgia_Mesolithic	42.28	43.28	9720.0	11.83484526	X Y	2.1_Caucasus_13000BP_10000BP
CHG	SATP	Georgia	Georgia_Mesolithic	42.38	42.59	13255.0	1.18417508	X Y	2.1_Caucasus_13000BP_10000BP
CHG	GD13a	Iran	Iran_Neolithic	34.45	48.12	9846.0	1.41728065	X X	2.1_Iran_10000BP_8500BP
CHG	DA383	Turkmenistan	Turkmenistan_Neolithic_Namazga	38.72	61.69	5150.0	0.7751095790000000	X X	2.1_Turkmenistan_7000BP_5000BP
CHG	NEO310	Turkmenistan	Turkmenistan_Neolithic	36.85	60.42	7150.0	1.278435989	X Y	2.1_Turkmenistan_7000BP_5000BP
EHG	Karelia	Russia	Russia_Mesolithic	61.65	35.65	8279.5	1.692885466	X Y	4.1_RussiaNW_11000BP_8000BP
EHG	NEO166	Russia	Russia_Neolithic	53.0	40.4	5668.0	2.1533629430000000	X Y	4.1_DonRiver_5800BP_5300BP
EHG	NEO167	Russia	Russia_Neolithic	53.0	40.4	5657.0	0.588921338	X X	4.1_DonRiver_5800BP_5300BP
EHG	NEO170	Russia	Russia_Neolithic	53.0	40.4	5562.0	0.4188968770000000	X X	4.1_DonRiver_5800BP_5300BP
EHG	NEO171	Russia	Russia_Neolithic	53.0	40.4	5835.0	0.818232704	X X	4.1_DonRiver_5800BP_5300BP
EHG	VK531	Norway	Norway_Neolithic	69.47	18	4350.0	1.4595637520000000	X Y	4.1_Norway_9300BP_4300BP
EHG	NEO173	Russia	Russia_Neolithic_Sredny	52.28	38.96	6345.0	0.689663925	X X	4.1_RussiaNW_7000BP_5000BP

EHG	NEO100	Russia	Russia_Mesolithic	51.57	53.68	9929.0	0.10793872200000000	X Y	4.1_RussiaNW_11000BP_8000BP
EHG	NEO88	Russia	Russia_Mesolithic	56.67	38.02	7871.0	2.6217337880000000	X Y	4.1_RussiaNW_11000BP_8000BP
EHG	H22	Norway	Norway_Mesolithic	58.83	6.33	9363.5	0.719978067	X X	4.1_Norway_9300BP_4300BP
EHG	stg001	Norway	Norway_Neolithic	67.76	14.85	5857.0	1.291071015	X Y	4.1_Norway_9300BP_4300BP
EHG	NEO87	Russia	Russia_Mesolithic	56.67	38.02	8259.0	0.183898205	X X	4.1_RussiaNW_11000BP_8000BP
EHG	NEO17	Norway	Norway_Mesolithic	58.06	7.74	9146.0	1.0991194990000000	X Y	4.1_Norway_9300BP_4300BP
EHG	Ukraine_N1	Ukraine	Ukraine_Neolithic	48.13	35.08	7250.0	0.167926522	X Y	4.1_Ukraine_10000BP_4000BP
EHG	Latvia_MN2	Latvia	Latvia_Neolithic_CC C	56.28	25.13	5965.0	1.147611647	X X	4.1_RussiaNW_7000BP_5000BP
EHG	NEO160	Russia	Russia_Neolithic	53.0	40.4	5269.0	1.326576858	X X	4.1_DonRiver_5800BP_5300BP
EHG	NEO178	Russia	Russia_Neolithic	56.78	40.45	5322.0	0.32433320600000000	X Y	4.1_RussiaNW_7000BP_5000BP
EHG	NEO163	Russia	Russia_Neolithic	53.0	40.4	5603.0	0.226384639	X Y	4.1_DonRiver_5800BP_5300BP
EHG	NEO180	Russia	Russia_Neolithic	56.78	40.45	5947.0	0.38518779200000000	X Y	4.1_RussiaNW_7000BP_5000BP
EHG	NEO687	Russia	Russia_Neolithic	57.58	58.2	5446.0	0.44089811100000000	X Y	4.1_RussiaNW_7000BP_5000BP
EHG	NEO560	Russia	Russia_Neolithic	60.41	38.93	7919.0	1.548541653	X Y	4.1_RussiaNW_11000BP_8000BP
EHG	NEO559	Russia	Russia_Neolithic	60.41	38.93	8268.0	1.228507174	X Y	4.1_RussiaNW_11000BP_8000BP
EHG	NEO557	Russia	Russia_Neolithic	60.41	38.93	7917.0	0.77765409100000000	X Y	4.1_RussiaNW_11000BP_8000BP
EHG	NEO556	Russia	Russia_Neolithic	60.41	38.93	7036.0	1.122201651	X X	4.1_RussiaNW_7000BP_5000BP
EHG	NEO555	Russia	Russia_Neolithic	60.41	38.93	8280.0	2.097005144	X Y	4.1_RussiaNW_11000BP_8000BP
EHG	NEO539	Russia	Russia_Mesolithic	59.7	39.5	10060.0	0.291758629	X X	4.1_RussiaNW_11000BP_8000BP
EHG	NEO179	Russia	Russia_Neolithic	56.78	40.45	5467.0	0.63935954700000000	X X	4.1_RussiaNW_7000BP_5000BP
EHG	NEO536	Russia	Russia_Mesolithic	59.7	39.5	9541.0	0.190177354	X Y	4.1_RussiaNW_11000BP_8000BP
EHG	NEO501	Ukraine	Ukraine_Mesolithic	48.2	35.22	10623.0	0.125455369	X Y	4.1_Ukraine_10000BP_4000BP
EHG	NEO202	Russia	Russia_Mesolithic	61.27	38.91	10884.0	2.217735963	X Y	4.1_RussiaNW_11000BP_8000BP
EHG	NEO197	Russia	Russia_Neolithic	56.78	40.45	5245.0	0.610890023	X Y	4.1_RussiaNW_7000BP_5000BP

EHG	NEO195	Russia	Russia_Neolithic	56.78	40.45	5749.0	0.59759596	X Y	4.1_RussiaNW_7000BP_5000BP
EHG	NEO174	Russia	Russia_Neolithic_Sredny	52.28	38.96	5306.0	1.1520597560000000	X X	4.1_DonRiver_5800BP_5300BP
EHG	NEO193	Russia	Russia_Neolithic	56.78	40.45	5453.0	0.232823778	X Y	4.1_RussiaNW_7000BP_5000BP
EHG	NEO184	Russia	Russia_Neolithic	56.78	40.44	5458.0	0.695336472	X X	4.1_RussiaNW_7000BP_5000BP
EHG	NEO185	Russia	Russia_Neolithic	56.78	40.45	7034.0	1.2132381190000000	X Y	4.1_RussiaNW_7000BP_5000BP
EHG	NEO194	Russia	Russia_Neolithic	56.78	40.45	5575.0	1.626304895	X Y	4.1_RussiaNW_7000BP_5000BP
EHG	NEO187	Russia	Russia_Neolithic	56.78	40.45	4940.0	0.16188839	X X	4.1_RussiaNW_7000BP_5000BP
EHG	Sidelkino	Russia	Russia_Mesolithic	54.53	51.11	11258.5	2.9967359070000000	X X	4.1_RussiaNW_11000BP_8000BP
EHG	NEO186	Russia	Russia_Neolithic	56.78	40.45	6922.0	0.36842406300000000	X Y	4.1_RussiaNW_7000BP_5000BP
EHG	NEO192	Russia	Russia_Neolithic	56.78	40.45	6841.0	0.269995424	X X	4.1_RussiaNW_7000BP_5000BP
EHG	NEO189	Russia	Russia_Neolithic	56.78	40.45	5648.0	0.614470507	X Y	4.1_RussiaNW_7000BP_5000BP
EastAsian	IK002	Japan	Japan_Jomon	34.65	137.14	2569.0	1.89109494	X X	3.1_Japan_3700BP_2600BP
EastAsian	DA45	Mongolia	Mongolia_IronAge_XiongNu	42.53	105.18	2095.0	9.039659899	X Y	3.1_SEAsia_4000BP_150BP
EastAsian	DA43	Mongolia	Mongolia_IronAge_XiongNu	42.53	105.18	2095.0	1.677926438	X Y	3.1_SEAsia_4000BP_150BP
EastAsian	DA39	Mongolia	Mongolia_IronAge_XiongNu	48.02	101.35	1948.0	2.087185526	X Y	3.1_SteppeCE_2000BP_700BP
EastAsian	DA38	Mongolia	Mongolia_IronAge_XiongNu	49.27	101.72	2124.5	2.85446773	X X	3.1_SteppeC_TianShan_2700BP_800BP
EastAsian	Funadomari_23	Japan	Japan_Jomon	45.38	141.04	3755.0	39.44124624	X X	3.1_Japan_3700BP_2600BP
EastAsian	Funadomari_5	Japan	Japan_Jomon	45.38	141.04	3755.0	3.8214367780000000	X Y	3.1_Japan_3700BP_2600BP
FarmerAnatolian	Bon001	Turkey	Anatolia_Neolithic_Aceramic	37.75	32.86	10032.0	0.16310050200000000	X Y	2.3_Anatolia_10000BP_8000BP
FarmerAnatolian	Bon002	Turkey	Anatolia_Neolithic_Aceramic	37.75	32.86	10078.0	6.692061812	X X	2.3_Anatolia_10000BP_8000BP
FarmerAnatolian	Bon004	Turkey	Anatolia_Neolithic_Aceramic	37.75	32.86	10076.0	0.242113723	X Y	2.3_Anatolia_10000BP_8000BP
FarmerAnatolian	Tep002	Turkey	Anatolia_Neolithic	38.17	34.49	8585.0	0.7072354540000000	X X	2.3_Anatolia_10000BP_8000BP
FarmerAnatolian	Tep004	Turkey	Anatolia_Neolithic	38.17	34.49	8295.0	0.467990841	X X	2.3_Anatolia_10000BP_8000BP
FarmerEarly	R3	Italy	Italy_Neolithic	41.96	13.54	7729.5	4.059641042	X X	2.3_EuropeS_8000BP_6000BP

FarmerEarly	R17	Italy	Italy_Neolithic	43.72	13.03	7223.5	0.56582164	X Y	2.3_EuropeS_8000BP_6000BP
FarmerEarly	R19	Italy	Italy_Neolithic	43.72	13.03	7233.0	0.52086802	X Y	2.3_EuropeS_8000BP_6000BP
FarmerEarly	R18	Italy	Italy_Neolithic	43.72	13.03	7298.5	0.6271842	X X	2.3_EuropeS_8000BP_6000BP
FarmerEarly	R16	Italy	Italy_Neolithic	43.72	13.03	7207.5	0.565514455	X X	2.3_EuropeS_8000BP_6000BP
FarmerEarly	R10	Italy	Italy_Neolithic	41.96	13.54	7629.0	1.3218194860000000	X X	2.3_EuropeS_8000BP_6000BP
FarmerEarly	R2	Italy	Italy_Neolithic	41.96	13.54	7984.0	3.7040638540000000	X X	2.3_EuropeS_8000BP_6000BP
FarmerEarly	R9	Italy	Italy_Neolithic	41.96	13.54	7496.0	4.04251971	X Y	2.3_EuropeS_8000BP_6000BP
FarmerEarly	MDV248	France	France_Neolithic_LBK	49.42	4.01	7015.5	0.12114648900000000	X Y	2.3_Europe_8500BP_5500BP
FarmerEarly	ROS45	France	France_Neolithic_Grossgartach	48.5	7.47	6642.5	0.26929279900000000	X Y	2.3_Europe_8500BP_5500BP
FarmerEarly	ROS78	France	France_Neolithic_Grossgartach	48.5	7.47	6550.0	0.435416491	X Y	2.3_Europe_8500BP_5500BP
FarmerEarly	Sch72-15	France	France_Neolithic_LBK	48.76	7.6	7036.5	0.235222248	X Y	2.3_Europe_8500BP_5500BP
FarmerEarly	Schw432	France	France_Neolithic_LBK	48.76	7.6	7100.0	0.148141996	X X	2.3_Europe_8500BP_5500BP
FarmerEarly	NEO137	Hungary	Hungary_Neolithic_Koros	46.42	20.33	7591.0	0.19840285800000000	X X	2.3_Europe_8500BP_5500BP
FarmerEarly	NEO140	Hungary	Hungary_Neolithic_Tisza	46.37	20.42	6718.0	0.14514291200000000	X Y	2.3_Europe_8500BP_5500BP
FarmerEarly	NEO145	Hungary	Hungary_Neolithic_Tisza	46.37	20.42	6744.0	0.21895615100000000	X X	2.3_Europe_8500BP_5500BP
FarmerEarly	NEO147	Hungary	Hungary_Neolithic_Tisza	46.37	20.42	6724.0	0.28457105900000000	X X	2.3_Europe_8500BP_5500BP
FarmerEarly	NEO695	Italy	Italy_Neolithic	43.08	13.06	7299.0	0.43141110300000000	X X	2.3_EuropeS_8000BP_6000BP
FarmerEarly	NEO674	Romania	Romania_Neolithic	44.9	22.43	5570.0	0.237000547	X Y	2.3_Europe_8500BP_5500BP
FarmerEarly	R8	Italy	Italy_Neolithic	41.96	13.54	7723.5	0.53180897	X X	2.3_EuropeS_8000BP_6000BP
FarmerEarly	kol6	CzechRepublic	Czech_Neolithic_Megalithic	50.03	15.2	6690.0	1.543500893	X X	2.3_EuropeCE_7000BP_5500BP
FarmerEarly	NEO655	Serbia	Serbia_Mesolithic	44.54	22.04	8668.0	0.22549267400000000	X X	2.3_Europe_8500BP_5500BP
FarmerEarly	PL_N36	Poland	Poland_Neolithic_BKG	50.67	21.38	6250.0	1.6699809730000000	X X	2.3_EuropeCE_7000BP_5500BP
FarmerEarly	NE5	Hungary	Hungary_Neolithic_ALP	47.17	20.83	7050.0	0.784933195	X Y	2.3_EuropeCE_7000BP_5500BP
FarmerEarly	PL_N31	Poland	Poland_Neolithic_BKG	52.61	18.8	6137.5	3.0038241510000000	X X	2.3_EuropeCE_7000BP_5500BP

FarmerEarly	NE1	Hungary	Hungary_Neolithic_LP	47.85	21.15	7138.5	18.42090683	XX	2.3_EuropeCE_7000BP_5500BP
FarmerEarly	Stuttgart	Germany	Germany_Neolithic	48.78	9.18	7140.0	16.19254244	XX	2.3_EuropeCE_7000BP_5500BP
FarmerEarly	Bar31	Turkey	Anatolia_Neolithic	40.3	29.61	8278.5	3.6490902710000000	XY	2.3_Anatolia_10000BP_8000BP
FarmerEarly	Bar8	Turkey	Anatolia_Neolithic	40.3	29.61	8071.0	7.171025264	XX	2.3_Anatolia_10000BP_8000BP
FarmerEarly	NE6	Hungary	Hungary_Neolithic_LBK	47.17	19.83	7051.5	0.937998868	XY	2.3_EuropeCE_7000BP_5500BP
FarmerEarly	PL_N25	Poland	Poland_Neolithic_BKG	52.61	18.8	6250.0	2.304346802	XX	2.3_EuropeCE_7000BP_5500BP
FarmerEarly	PL_N26	Poland	Poland_Neolithic_BKG	50.67	21.38	6151.0	2.146927132	XY	2.3_EuropeCE_7000BP_5500BP
FarmerEarly	PL_N19	Poland	Poland_Neolithic_FBC	52.62	18.96	5462.0	1.6841604730000000	XX	2.3_EuropeCE_7000BP_5500BP
FarmerEarly	Pal7	Greece	Greece_Neolithic	40.51	22.5	6351.0	1.265496659	XX	2.3_EuropeCE_7000BP_5500BP
FarmerEarly	PL_N28	Poland	Poland_Neolithic_BKG	52.61	18.8	6073.5	1.75499685	XY	2.3_EuropeCE_7000BP_5500BP
FarmerEarly	Rev5	Greece	Greece_Neolithic	40.33	22.56	8301.0	1.133693538	XX	2.3_EuropeS_8000BP_6000BP
FarmerEarly	Klei10	Greece	Greece_Neolithic	40.26	21.74	6062.5	2.047213233	XY	2.3_EuropeS_8000BP_6000BP
FarmerEarly	NE2	Hungary	Hungary_Neolithic_LP	47.52	21.59	7123.5	0.148202309	XX	2.3_Europe_8500BP_5500BP
FarmerEarly	PL_N27	Poland	Poland_Neolithic_BKG	52.61	18.8	6250.0	1.790218266	XY	2.3_EuropeCE_7000BP_5500BP
FarmerEarly	NE7	Hungary	Hungary_Neolithic_Lengyel	47.17	19.83	6374.0	0.909572155	XY	2.3_EuropeCE_7000BP_5500BP
FarmerLate	NEO119	France	France_Neolithic	44.47	4.77	4382.0	0.10885148	XX	2.2_EuropeSW_6000BP_3500BP
FarmerLate	NEO886	Denmark	Denmark_Neolithic	54.99	12.42	5457.0	0.271693685	XY	2.4_EuropeNE_5600BP_4600BP
FarmerLate	NEO925	Denmark	Denmark_Neolithic	54.77	10.68	4947.0	0.294332844	XX	2.4_EuropeNE_5600BP_4600BP
FarmerLate	NEO653	Spain	Iberia_BronzeAge	43.4	-4.71	3423.0	0.283734421	XX	2.2_EuropeSW_6000BP_3500BP
FarmerLate	TV3831	Portugal	Iberia_BronzeAge	37.94	-7.6	3550.0	0.9943057240000000	XY	2.2_Iberia_7300BP_3500BP
FarmerLate	COV20126	Spain	Iberia_Neolithic	37.41	-4.42	5588.0	0.303266123	XY	2.2_EuropeSW_6000BP_3500BP
FarmerLate	NEO896	Denmark	Denmark_Neolithic	54.97	12.49	5446.0	0.121712208	XX	2.4_EuropeNE_5600BP_4600BP
FarmerLate	NEO43	Denmark	Denmark_Neolithic	55.99	10.25	5067.0	0.108400524	XY	2.4_EuropeNE_5600BP_4600BP
FarmerLate	NEO39	Sweden	Sweden_Neolithic	55.57	13.04	5074.0	0.174678265	XY	2.4_EuropeNE_5600BP_4600BP

FarmerLate	TV32032extra	Portugal	Iberia_BronzeAge	37.94	-7.6	3550.0	0.865992675	X Y	2.2_Iberia_7300BP_3500BP
FarmerLate	NEO744	Denmark	Denmark_Neolithic	55.86	11.59	5333.0	0.220878469	X Y	2.4_EuropeNE_5600BP_4600BP
FarmerLate	MonteGato104	Portugal	Iberia_BronzeAge	38.02	-7.86	3535.0	1.236961027	X Y	2.2_Iberia_7300BP_3500BP
FarmerLate	NEO830	Italy	Italy_Neolithic	43.38	13.55	5393.0	0.104506976	X X	2.2_EuropeSW_6000BP_3500BP
FarmerLate	NEO757	Denmark	Denmark_Neolithic	55.9	11.12	5452.0	0.129920601	X X	2.4_EuropeNE_5600BP_4600BP
FarmerLate	atp9	Spain	Iberia_BronzeAge	42.35	-3.52	3634.0	0.419041583	X X	2.2_EuropeSW_6000BP_3500BP
FarmerLate	QUIN234	France	France_BronzeAge	43.3	1.96	3600.0	0.120622268	X X	2.2_EuropeSW_6000BP_3500BP
FarmerLate	ROS102	France	France_Neolithic_Grossgartach	48.5	7.47	6550.0	0.151456852	X Y	2.3_Europe_8500BP_5500BP
FarmerLate	NEO640	Poland	Poland_Neolithic_FBC	50.27	20.45	4902.0	0.20055724	X Y	2.4_EuropeNE_5600BP_4600BP
FarmerLate	NEO609	Portugal	Iberia_Neolithic	38.68	-9.16	4333.0	0.134584406	X Y	2.2_EuropeSW_6000BP_3500BP
FarmerLate	NEO599	Denmark	Denmark_Neolithic	55.55	11.68	5134.0	0.19019233	X Y	2.4_EuropeNE_5600BP_4600BP
FarmerLate	NEO597	Denmark	Denmark_Neolithic	55.59	11.57	5210.0	0.177413333	X Y	2.4_EuropeNE_5600BP_4600BP
FarmerLate	NEO595	Denmark	Denmark_Neolithic	55.79	11.29	5452.0	0.218556781	X Y	2.4_EuropeNE_5600BP_4600BP
FarmerLate	ans016	Sweden	Sweden_Neolithic	57.34	18.26	4646.0	0.34085699700000000	X Y	2.4_EuropeNE_5600BP_4600BP
FarmerLate	NEO945	Denmark	Denmark_Neolithic	56.49	9.83	5445.0	1.384502704	X X	2.4_EuropeNE_5600BP_4600BP
FarmerLate	NEO25	Denmark	Denmark_Neolithic	56.43	10.79	4956.0	0.35636724800000000	X X	2.4_EuropeNE_5600BP_4600BP
FarmerLate	ans005	Sweden	Sweden_Neolithic_Megalithic	57.34	18.26	5265.0	0.129055433	X X	2.4_EuropeNE_5600BP_4600BP
FarmerLate	NEO721	Spain	Iberia_Neolithic	40.44	-3.5	4170.0	0.35821067700000000	X X	2.2_Iberia_7300BP_3500BP
FarmerLate	NEO943	Denmark	Denmark_Neolithic	55.46	9.69	4614.0	1.754137916	X Y	2.4_EuropeNE_5600BP_4600BP
FarmerLate	NEO942	Denmark	Denmark_Neolithic	55.58	11.29	5491.0	0.890575166	X X	2.4_EuropeNE_5600BP_4600BP
FarmerLate	NEO753	Denmark	Denmark_Neolithic	55.77	12.21	5531.0	0.163364465	X X	2.4_EuropeNE_5600BP_4600BP
FarmerLate	esp005	Spain	Iberia_BronzeAge_Cogotas	42.0	-1	3370.0	2.457932356	X Y	2.2_Iberia_7300BP_3500BP
FarmerLate	pir001	Spain	Iberia_BronzeAge_Argar	37.89	-4.78	3725.0	0.21560743100000000	X Y	2.2_EuropeSW_6000BP_3500BP
FarmerLate	por004	Spain	Iberia_Neolithic	42.35	-3.52	4955.0	0.129932409	X Y	2.2_EuropeSW_6000BP_3500BP

FarmerLate	san216	Spain	Iberia_Neolithic	42.69	-2.73	5665.5	0.200001027	X Y	2.2_EuropeSW_6000BP_3500BP
FarmerLate	ans003	Sweden	Sweden_Neolithic_Megalithic	57.34	18.26	5250.0	0.135654055	X X	2.4_EuropeNE_5600BP_4600BP
FarmerLate	ValeOuro10207	Portugal	Iberia_BronzeAge	38.06	-8.11	3550.0	0.248412062	X X	2.2_EuropeSW_6000BP_3500BP
FarmerMiddle	NEO121	France	France_Neolithic	44.47	4.77	4531.0	0.530892193	X Y	2.2_EuropeSW_6000BP_3500BP
FarmerMiddle	NEO28	Denmark	Denmark_Neolithic	55.91	12.31	5459.0	0.922300627	X Y	2.4_EuropeNE_5600BP_4600BP
FarmerMiddle	NEO36	Sweden	Sweden_Neolithic	55.57	13.04	5097.0	2.384739487	X X	2.4_EuropeNE_5600BP_4600BP
FarmerMiddle	NEO29	Denmark	Denmark_Neolithic	55.13	10.9	5489.0	0.530600133	X X	2.4_EuropeNE_5600BP_4600BP
FarmerMiddle	NEO23	Denmark	Denmark_Neolithic	55.6	11.31	5533.0	3.33772768	X Y	2.4_EuropeNE_5600BP_4600BP
FarmerMiddle	WET370	France	France_Neolithic	48.06	7.3	5521.0	0.172865405	X Y	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	NEO866	Denmark	Denmark_Neolithic	54.87	11.84	5456.0	1.521286722	X Y	2.4_EuropeNE_5600BP_4600BP
FarmerMiddle	NEO891	Denmark	Denmark_Neolithic	55.73	12.1	5661.0	0.595471133	X X	2.4_EuropeNE_5600BP_4600BP
FarmerMiddle	ans017	Sweden	Sweden_Neolithic_Megalithic	57.34	18.26	5080.0	7.08351096	X Y	2.4_EuropeNE_5600BP_4600BP
FarmerMiddle	CO1	Hungary	Hungary_Neolithic_Baden	47.17	19.83	4750.0	0.8595327420000000	X X	2.3_EuropeS_8000BP_6000BP
FarmerMiddle	RISE489	Italy	Italy_Neolithic	45.26	10.38	4693.0	0.50758911	X Y	2.2_EuropeSW_6000BP_3500BP
FarmerMiddle	NEO641	Poland	Poland_Neolithic_FBC	50.27	20.45	5132.0	0.26184454	X X	2.3_EuropeS_8000BP_6000BP
FarmerMiddle	NEO630	UK	Britain_Neolithic	58.73	-2.94	4898.0	0.211902615	X Y	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	NEO627	UK	Britain_Neolithic	58.73	-2.94	5132.0	0.4495226710000000	X Y	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	NEO626	UK	Britain_Neolithic	58.73	-2.94	5082.0	0.3719390520000000	X Y	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	NEO624	UK	Britain_Neolithic	58.73	-2.94	4897.0	2.099933966	X X	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	NEO717	UK	Britain_Neolithic	58.73	-2.94	4893.0	0.481287489	X Y	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	NEO935	Denmark	Denmark_Neolithic	55.56	12.02	5187.0	5.027509563	X Y	2.4_EuropeNE_5600BP_4600BP
FarmerMiddle	NEO933	Denmark	Denmark_Neolithic	55.25	10.75	5337.0	0.522015892	X Y	2.4_EuropeNE_5600BP_4600BP
FarmerMiddle	Aveline_1	UK	Britain_Neolithic	51.32	-2.75	5489.5	0.782428624	X X	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	NEO142	Hungary	Hungary_Neolithic_Tisza	46.37	20.42	6641.0	0.895629723	X X	2.3_Europe_8500BP_5500BP

FarmerMiddle	BurnGround	UK	Britain_Neolithic	51.84	-1.85	5770.0	0.410850116	X Y	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	CaveHa3_1	UK	Britain_Neolithic	54.07	-2.29	5264.0	0.112689685	X Y	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	Coldrum_1	UK	Britain_Neolithic	51.32	0.37	5430.0	0.547611201	X Y	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	Embo_1	UK	Britain_Neolithic	57.91	-4	5050.0	0.127757081	X X	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	Fussels_Lodge_1	UK	Britain_Neolithic	51.09	-1.73	5657.5	0.73024687	X X	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	Jubilee_cave	UK	Britain_Neolithic	54.08	-2.27	5462.5	0.116653425	X Y	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	Kelco_cave	UK	Britain_Neolithic	54.07	-2.29	5536.0	0.115486462	X Y	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	Ballynahatty	Ireland	Ireland_Neolithic	54.54	-5.96	5131.5	10.0157214	X X	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	PL_N18	Poland	Poland_Neolithic_FBC	52.62	18.96	5462.5	1.909978283	X X	2.4_EuropeNE_5600BP_4600BP
FarmerMiddle	PSS4693	France	France_Neolithic_Noyen	48.52	3.6	5438.5	0.740638747	X X	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	PL_N38	Poland	Poland_Neolithic_GAC	52.61	18.9	5033.0	1.7720028230000000	X X	2.4_Poland_5000BP_4700BP
FarmerMiddle	Carsington_pasture_1	UK	Britain_Neolithic	53.08	-1.64	5538.5	9.748563794000000	X Y	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	PL_N20	Poland	Poland_Neolithic_FBC	52.62	18.96	5462.0	0.8023952120000000	X Y	2.4_EuropeNE_5600BP_4600BP
FarmerMiddle	Mor6	France	France_Neolithic_LBK	48.82	7.63	7036.0	0.16119936100000000	X X	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	ros3	Sweden	Sweden_Neolithic_FBC	60.26	16.41	4955.0	0.38093967100000000	X X	2.4_EuropeNE_5600BP_4600BP
FarmerMiddle	RISE1161	Poland	Poland_Neolithic_GAC	48.7	21.2	4757.0	1.366034772	X X	2.4_Poland_5000BP_4700BP
FarmerMiddle	RISE1165	Poland	Poland_Neolithic_GAC	48.7	21.2	4742.5	2.189830631	X Y	2.4_Poland_5000BP_4700BP
FarmerMiddle	RISE1166	Poland	Poland_Neolithic_GAC	48.7	21.2	4907.5	3.058742732	X X	2.4_Poland_5000BP_4700BP
FarmerMiddle	RISE1249	Poland	Poland_Neolithic_GAC	50.2	21.4	4736.5	1.0104430690000000	X X	2.4_Poland_5000BP_4700BP
FarmerMiddle	RISE1248	Poland	Poland_Neolithic_GAC	50.2	21.4	4725.0	0.7995153760000000	X X	2.4_Poland_5000BP_4700BP
FarmerMiddle	RISE1246	Poland	Poland_Neolithic_GAC	50.2	21.4	4715.0	0.549865829	X X	2.4_Poland_5000BP_4700BP
FarmerMiddle	RISE1252	Poland	Poland_Neolithic_GAC	50.8	21.5	4725.0	0.43094765800000000	X Y	2.4_Poland_5000BP_4700BP
FarmerMiddle	RISE1250	Poland	Poland_Neolithic_GAC	50.6	21.7	4725.0	0.493666575	X Y	2.4_Poland_5000BP_4700BP
FarmerMiddle	RISE1254	Poland	Poland_Neolithic_GAC	51.1	17.1	4725.0	0.437402634	X Y	2.4_Poland_5000BP_4700BP

FarmerMiddle	Gok2	Sweden	Sweden_Neolithic_FBC	58.18	13.41	4850.0	1.220935916	XX	2.4_EuropeNE_5600BP_4600BP
FarmerMiddle	NEO847	UK	Britain_Neolithic	51.7	-2.3	5463.0	1.777915755	XY	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	ros005	Sweden	Sweden_Neolithic_FBC	60.26	16.41	4740.0	0.886088897	XY	2.4_EuropeNE_5600BP_4600BP
FarmerMiddle	PEI 2.00	France	France_Neolithic_Campaniforme	43.14	2.25	4385.5	0.302517672	XY	2.2_EuropeSW_6000BP_3500BP
FarmerMiddle	R1014	Italy	Italy_Neolithic	41.37	13.29	4950.0	0.615266206	XY	2.2_Italy_7000BP_4000BP
FarmerMiddle	R104	Italy	Italy_LateAntiquity	41.89	12.48	1450.0	0.879014223	XY	2.2_Italy_7000BP_4000BP
FarmerMiddle	RISE1159	Poland	Poland_Neolithic_GAC	48.7	21.2	4730.0	27.46258284	XX	2.4_Poland_5000BP_4700BP
FarmerMiddle	RISE1170	Poland	Poland_Neolithic_GAC	48.7	21.2	4748.5	3.79009955	XX	2.4_Poland_5000BP_4700BP
FarmerMiddle	RISE1168	Poland	Poland_Neolithic_GAC	48.7	21.2	4676.0	18.93418868	XY	N/A
FarmerMiddle	mur	Spain	Iberia_Neolithic_Almagra	42.35	-3.52	7136.0	3.4675614490000000	XY	2.2_Iberia_7300BP_3500BP
FarmerMiddle	lai001	UK	Britain_Neolithic	59.13	-3.05	5180.0	0.225980982	XY	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	mid001	UK	Britain_Neolithic_Megalithic	59.13	-3.05	5450.0	0.282625993	XY	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	mid002	UK	Britain_Neolithic_Megalithic	57.75	-3.92	5180.0	0.2557871540000000	XY	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	prs002	Ireland	Ireland_Neolithic_Megalithic	54.25	-8.56	5675.0	5.870842597	XX	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	prs003	Ireland	Ireland_Neolithic_Megalithic	54.25	-8.56	5600.0	0.2237872080000000	XY	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	prs006	Ireland	Ireland_Neolithic_Megalithic	54.25	-8.56	5330.0	0.2635244420000000	XX	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	prs009	Ireland	Ireland_Neolithic_Megalithic	54.25	-8.56	5310.0	7.571866381	XY	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	NEO823	Italy	Italy_BronzeAge	40.88	16.73	4665.0	0.404577629	XY	2.2_Italy_7000BP_4000BP
FarmerMiddle	prs010	Ireland	Ireland_Neolithic_Megalithic	54.25	-8.56	5530.0	0.2320574680000000	XY	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	prs013	Ireland	Ireland_Neolithic_Megalithic	54.25	-8.56	5320.0	4.952996849	XY	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	prs016	Ireland	Ireland_Neolithic_Megalithic	54.25	-8.56	5560.0	8.754512586	XY	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	ans008	Sweden	Sweden_Neolithic_Megalithic	57.34	18.26	5135.0	2.027487661	XY	2.4_EuropeNE_5600BP_4600BP
FarmerMiddle	CB13	Spain	Iberia_Neolithic_Cardial	41.37	1.89	7348.0	0.931947851	XX	2.2_Iberia_7300BP_3500BP
FarmerMiddle	atp016	Spain	Iberia_Neolithic	42.35	-3.52	5039.5	13.20832827	XX	2.2_Iberia_7300BP_3500BP
FarmerMiddle	atp12-1420	Spain	Iberia_Neolithic	42.35	-3.52	4895.5	2.528221016	XY	2.2_Iberia_7300BP_3500BP

FarmerMiddle	c40331	Spain	Iberia_Neolithic	37.37	-4.25	5649.5	0.293459982	X Y	2.2_Iberia_7300BP_3500BP
FarmerMiddle	prs012	Ireland	Ireland_Neolithic_Megalithic	54.25	-8.56	5660.0	0.25153112700000000	X Y	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	LugarCanto44	Portugal	Iberia_Neolithic	39.42	-8.82	5950.0	2.016550504	X X	2.2_Iberia_7300BP_3500BP
FarmerMiddle	RISE1241	Poland	Poland_Neolithic_GAC	50.6	21.7	4752.5	0.859098485	X Y	2.4_Poland_5000BP_4700BP
FarmerMiddle	R22	Italy	Italy_Neolithic	40.81	8.44	3895.5	0.776104455	X X	2.2_Italy_7000BP_4000BP
FarmerMiddle	BERG157-2	France	France_Neolithic_BORSMichelsberg	43.22	2.41	6050.0	0.346274491	X X	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	BERG157-7	France	France_Neolithic_BORSMichelsberg	43.22	2.41	6131.5	0.267013925	X Y	2.2_EuropeSW_6000BP_3500BP
FarmerMiddle	CabecoArruda117B	Portugal	Iberia_Neolithic	39.11	-8.66	5050.0	0.376607974	X Y	2.2_Iberia_7300BP_3500BP
FarmerMiddle	BLP10	France	France_Neolithic_Michelsberg	49.39	3.74	6052.0	0.18285663600000000	X Y	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	BUCH2	France	France_Neolithic_Cemy	48.24	4.11	6250.0	0.36476039200000000	X Y	2.2_EuropeSW_6000BP_3500BP
FarmerMiddle	NEO812	France	France_Neolithic_Cardial	43.32	2.42	6545.0	6.542360034	X X	2.2_Iberia_7300BP_3500BP
FarmerMiddle	CRE20D	France	France_Neolithic_ChasseenAncien	43.21	3.13	6151.0	0.256045807	X X	2.2_EuropeSW_6000BP_3500BP
FarmerMiddle	LugarCanto42	Portugal	Iberia_Neolithic	39.42	-8.82	5950.0	3.006333862	X X	2.2_Iberia_7300BP_3500BP
FarmerMiddle	LU339	Portugal	Iberia_Neolithic	41.71	-6.93	6797.5	4.60334026	X Y	2.2_Iberia_7300BP_3500BP
FarmerMiddle	LD270	Portugal	Iberia_Neolithic	41.71	-6.93	6336.0	4.064587193	X Y	2.2_Iberia_7300BP_3500BP
FarmerMiddle	LD1174	Spain	Iberia_Neolithic	37.41	-4.42	6415.0	3.558801721	X X	2.2_Iberia_7300BP_3500BP
FarmerMiddle	CabecoArruda122A	Portugal	Iberia_Neolithic	39.11	-8.66	5050.0	1.782958508	X Y	2.2_Iberia_7300BP_3500BP
FarmerMiddle	CovaMoura364	Portugal	Iberia_Neolithic	38.75	-9.22	4900.0	0.794100402	X Y	2.2_Iberia_7300BP_3500BP
FarmerMiddle	Es97-1	France	France_Neolithic_Michelsberg	50.92	1.71	6004.5	0.294790665	X Y	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	NEO790	Denmark	Denmark_Neolithic	55.71	12.27	5663.0	0.685227719	X Y	2.4_EuropeNE_5600BP_4600BP
FarmerMiddle	CovaMoura9B	Portugal	Iberia_Neolithic	38.75	-9.22	4900.0	2.611737333	X X	2.2_Iberia_7300BP_3500BP
FarmerMiddle	R6	Italy	Italy_Neolithic	41.96	13.54	7159.5	0.604714196	X Y	2.2_Italy_7000BP_4000BP
FarmerMiddle	bal004	UK	Britain_Neolithic_Megalithic	57.77	-3.9	5190.0	1.56791032700000000	X X	2.2_EuropeAtlantic_7000BP_5000BP
FarmerMiddle	R24	Italy	Italy_Neolithic	40.81	8.44	5450.0	0.549967737	X X	2.2_Italy_7000BP_4000BP
FarmerMiddle	R25	Italy	Italy_Neolithic	40.81	8.44	4150.0	0.53976087	X X	2.2_Italy_7000BP_4000BP

FarmerMiddle	R26	Italy	Italy_Neolithic	40.81	8.44	4150.0	0.525953541	X Y	2.2_Italy_7000BP_400BP
FarmerMiddle	R27	Italy	Italy_Neolithic	40.81	8.44	4150.0	0.70739547	X Y	2.2_Italy_7000BP_400BP
FarmerMiddle	R29	Italy	Italy_Neolithic	40.81	8.44	4150.0	0.559161215	X Y	2.2_Italy_7000BP_400BP
FarmerMiddle	R28	Italy	Italy_Neolithic	40.81	8.44	4150.0	0.728827575	X X	2.2_Italy_7000BP_400BP
FarmerMiddle	Dolmen Ansião 96B	Portugal	Iberia_Neolithic	39.75	-8.81	5450.0	1.962153759	X Y	2.2_Iberia_7300BP_350BP
FarmerMiddle	R4	Italy	Italy_Neolithic	41.96	13.54	4865.0	3.6761982310000000	X Y	2.2_Italy_7000BP_4000BP
FarmerMiddle	R5	Italy	Italy_Neolithic	41.96	13.54	4839.5	1.5029056810000000	X X	2.2_Italy_7000BP_4000BP
FarmerMiddle	LugarCanto41	Portugal	Iberia_Neolithic	39.42	-8.82	5950.0	1.06714512	X X	2.2_Iberia_7300BP_350BP
FarmerMiddle	BERG02-2	France	France_Neolithic_BORSMichelsberg	43.22	2.41	5870.0	0.344146565	X X	2.2_EuropeSW_6000BP_3500BP
WHG	NEO855	Denmark	Denmark_Mesolithic	56.4	10.72	6302.0	1.3829770000000000	X Y	4.2_Denmark_10500BP_6000BP
WHG	NEO856	Denmark	Denmark_Mesolithic	56.37	10.64	6777.0	0.56205807	X X	4.2_Denmark_10500BP_6000BP
WHG	NEO679	Sweden	Sweden_Mesolithic	55.39	13.48	6834.0	0.164673359	X X	4.2_Denmark_10500BP_6000BP
WHG	NEO683	Denmark	Denmark_Mesolithic	55.4	9.83	7529.0	1.81852777	X X	4.2_Denmark_10500BP_6000BP
WHG	NEO938	Spain	Iberia_Mesolithic	43.4	-4.71	7878.0	0.4751514570000000	X X	4.2_Iberia_9000BP_700BP
WHG	NEO694	Spain	Iberia_Mesolithic	38.73	-0.46	9217.0	0.2847588340000000	X Y	4.2_Iberia_9000BP_700BP
WHG	NEO853	Denmark	Denmark_Mesolithic	55.55	10.62	6047.0	1.964862968	X X	4.2_Denmark_10500BP_6000BP
WHG	NEO852	Denmark	Denmark_Mesolithic	56.03	10.26	6308.0	0.189591791	X X	4.2_Denmark_10500BP_6000BP
WHG	NEO733	Denmark	Denmark_Mesolithic	55.77	11.39	6824.0	1.3166819510000000	X X	4.2_Denmark_10500BP_6000BP
WHG	NEO791	Denmark	Denmark_Mesolithic	55.33	11.15	7048.0	2.492448215	X Y	4.2_Denmark_10500BP_6000BP
WHG	NEO941	Denmark	Denmark_Mesolithic	56.71	10.17	6372.0	0.135296816	X X	4.2_Denmark_10500BP_6000BP
WHG	NEO751	Denmark	Denmark_Mesolithic	56.87	9.22	6343.0	0.297822065	X Y	4.2_Denmark_10500BP_6000BP
WHG	NEO932	Denmark	Denmark_Mesolithic	55.25	11.23	7499.0	2.7601622000000000	X X	4.2_Denmark_10500BP_6000BP
WHG	NEO749	Denmark	Denmark_Mesolithic	55.85	12.56	7070.0	1.905133435	X Y	4.2_Denmark_10500BP_6000BP
WHG	NEO747	Denmark	Denmark_Mesolithic	55.85	12.56	6729.0	0.2494273580000000	X Y	4.2_Denmark_10500BP_6000BP
WHG	NEO745	Denmark	Denmark_Mesolithic	55.85	12.56	6790.0	0.447895875	X Y	4.2_Denmark_10500BP_6000BP
WHG	NEO960	Denmark	Denmark_Mesolithic	55.58	11.58	5926.0	0.150141897	X Y	4.2_Denmark_10500BP_6000BP

WHG	NEO759	Denmark	Denmark_Mesolithic	55.4	12.37	9028.0	2.948103712	X Y	4.2_Denmark_10500BP_6000BP
WHG	sylltholm	Denmark	Denmark_Mesolithic	54.65	11.35	7709.5	2.291787968	X X	4.2_Denmark_10500BP_6000BP
WHG	NEO648	Spain	Iberia_Mesolithic	43.4	-4.71	7539.0	1.844283075	X Y	4.2_Iberia_9000BP_7000BP
WHG	Cheddar_man	UK	Britain_Mesolithic	51.28	-2.77	10300.0	2.054202123	X Y	4.2_EuropeW_13500BP_8000BP
WHG	PL_N22	Poland	Poland_Neolithic_BKG	52.61	18.9	6291.0	1.491551035	X X	4.2_EuropeE_8600BP_6000BP
WHG	KO1	Hungary	Hungary_Neolithic_Koros	47.56	20.72	7660.0	1.014600016	X Y	4.2_EuropeE_8600BP_6000BP
WHG	Canes1	Spain	Iberia_Mesolithic	43.36	-4.72	7115.0	1.6461215340000000	X X	4.2_Iberia_9000BP_7000BP
WHG	Chan	Spain	Iberia_Mesolithic	42.73	-7.03	9131.0	5.008215765	X X	4.2_Iberia_9000BP_7000BP
WHG	Bichon	Switzerland	Switzerland_Mesolithic	47.1	6.87	13665.0	7.692393112	X Y	4.2_EuropeW_13500BP_8000BP
WHG	Loschbour	Luxembourg	Luxembourg_Mesolithic	49.81	6.4	8050.0	18.23029647	X Y	4.2_EuropeW_13500BP_8000BP
WHG	Brana	Spain	Iberia_Mesolithic	42.91	-5.38	7815.0	3.019525774	X Y	4.2_Iberia_9000BP_7000BP
WHG	R11	Italy	Italy_Mesolithic	41.96	13.54	11908.0	0.957641603	X Y	4.2_Italy_15000BP_9000BP
WHG	NEO669	Serbia	Serbia_Mesolithic	44.56	22.03	7950.0	0.23932044	X X	4.2_EuropeE_8600BP_6000BP
WHG	R7	Italy	Italy_Mesolithic	41.96	13.54	10681.5	3.153769086	X Y	4.2_Italy_15000BP_9000BP
WHG	ST3	Italy	Italy_Mesolithic	37.85	14.7	14800.0	0.475034886	X Y	4.2_Italy_15000BP_9000BP
WHG	PER1150503	France	France_Mesolithic	45.77	0.33	9067.0	0.315139903	X X	4.2_EuropeW_13500BP_8000BP
WHG	PER3023	France	France_Mesolithic	45.77	0.33	9067.0	0.161333797	X X	4.2_EuropeW_13500BP_8000BP
WHG	R15	Italy	Italy_Mesolithic	41.96	13.54	9124.5	3.070164483	X Y	4.2_Italy_15000BP_9000BP
WHG	NEO91	Denmark	Denmark_Mesolithic	55.39	12.31	9122.0	1.176549838	X Y	4.2_Denmark_10500BP_6000BP
WHG	NEO646	Spain	Iberia_Mesolithic	43.4	-4.71	8273.0	1.590267827	X X	4.2_Iberia_9000BP_7000BP
WHG	NEO645	Denmark	Denmark_Mesolithic	55.91	11.09	5870.0	0.211986752	X Y	4.2_Denmark_10500BP_6000BP
WHG	NEO598	Denmark	Denmark_Mesolithic	55.95	11.9	6075.0	0.7272044320000000	X X	4.2_Denmark_10500BP_6000BP
WHG	NEO19	Denmark	Denmark_Mesolithic	56.27	10.47	8163.0	3.262849417	X Y	4.2_Denmark_10500BP_6000BP
WHG	NEO586	Denmark	Denmark_Mesolithic	56.37	10.57	7031.0	0.201188562	X Y	4.2_Denmark_10500BP_6000BP
WHG	NEO583	Denmark	Denmark_Mesolithic	56.37	10.57	6981.0	0.1763990890000000	X X	4.2_Denmark_10500BP_6000BP
WHG	NEO570	Denmark	Denmark_Mesolithic	56.4	10.72	6369.0	2.861753026	X X	4.2_Denmark_10500BP_6000BP

WHG	NEO589	Denmark	Denmark_Mesolithic	55.33	11.15	7478.0	7.410700578000000	X Y	4.2_Denmark_10500BP_6000BP
WHG	NEO254	Denmark	Denmark_Mesolithic	55.4	10.13	10463.0	0.41962998	X Y	4.2_Denmark_10500BP_6000BP
WHG	NEO123	Denmark	Denmark_Mesolithic	54.96	11.85	8182.0	0.286386228	X X	4.2_Denmark_10500BP_6000BP
WHG	NEO122	Denmark	Denmark_Mesolithic	54.96	11.85	8146.0	0.564966744	X X	4.2_Denmark_10500BP_6000BP
WHG	NEO568	Denmark	Denmark_Mesolithic	56.81	9.18	6586.0	1.981486947000000	X Y	4.2_Denmark_10500BP_6000BP
Yamnaya	poz81	Poland	Poland_Neolithic_CWC	52.29	17.55	4705.0	1.92879788	X Y	1.2_EuropeNE_4800BP_3000BP
Yamnaya	RISE509	Russia	Siberia_BronzeAge_Afanasiovo	54.36	90.92	4732.0	4.52834127	X X	1.2_Steppe_5000BP_4300BP
Yamnaya	RISE511	Russia	Siberia_BronzeAge_Afanasiovo	54.36	90.92	4744.0	5.20403929	X X	1.2_Steppe_5000BP_4300BP
Yamnaya	RISE546	Russia	Russia_BronzeAge_Yamnaya	46.54	43.7	4850.0	0.125905828	X Y	1.2_Steppe_5000BP_4300BP
Yamnaya	RISE547	Russia	Russia_BronzeAge_Yamnaya	46.54	43.7	4710.5	0.686466601000000	X Y	1.2_Steppe_5000BP_4300BP
Yamnaya	RISE548	Russia	Russia_BronzeAge_Yamnaya	46.54	43.7	4850.0	0.910878358000000	X Y	1.2_Steppe_5000BP_4300BP
Yamnaya	RISE550	Russia	Russia_BronzeAge_Yamnaya	46.56	43.68	4934.5	0.440260727	X Y	1.2_Steppe_5000BP_4300BP
Yamnaya	RISE555	Russia	Russia_BronzeAge	48.72	44.5	4627.0	0.237337432	X Y	1.2_Steppe_5000BP_4300BP
Yamnaya	Yamnaya	Kazakhstan	Kazakhstan_BronzeAge_Yamnaya	49.13	75.85	4902.5	26.39165529	X Y	1.2_Steppe_5000BP_4300BP
Yamnaya	MJ-09	Ukraine	Ukraine_BronzeAge_Catacomb	47.43	34.27	4285.5	0.199475487000000	X X	1.2_Steppe_5000BP_4300BP
Yamnaya	MJ-06	Ukraine	Ukraine_BronzeAge_Yamnaya	49.32	35.37	4629.5	0.161999877	X X	1.2_EuropeNE_4800BP_3000BP
Yamnaya	NEO175	Russia	Russia_Neolithic_Sredny	52.28	38.96	4607.0	0.416698274	X Y	1.2_Steppe_5000BP_4300BP
Yamnaya	Latvia_LN1	Latvia	Latvia_Neolithic_CWC	56.28	25.13	4833.0	0.197755635	X X	1.2_EuropeNE_4800BP_3000BP
Yamnaya	RISE552	Russia	Russia_BronzeAge_Yamnaya	46.62	43.33	4446.0	2.458824579	X Y	1.2_Steppe_5000BP_4300BP
Yamnaya	RISE240	Russia	Russia_BronzeAge_Yamnaya	46.58	43.68	4706.0	0.173195772	X X	1.2_Steppe_5000BP_4300BP

2272 **Table S3h.1. Metadata and grouping of ancient individuals into reference populations.**

2273

2274 References

2275 1. Lawson, D. J., Hellenthal, G., Myers, S. & Falush, D. Inference of population structure
2276 using dense haplotype data. *PLoS Genet.* **8**, e1002453 (2012).

2277 2. Margaryan, A. *et al.* Population genomics of the Viking world. *Nature* **585**, 390–396
2278 (2020).

- 2279 3. Hellenthal, G. *et al.* A genetic atlas of human admixture history. *Science* **343**, 747–
2280 751 (2014).
- 2281 4. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.
2282 *Nature* **562**, 203–209 (2018).
- 2283 5. Li, N. & Stephens, M. Modeling linkage disequilibrium and identifying recombination
2284 hotspots using single-nucleotide polymorphism data. *Genetics* **165**, 2213–2233 (2003).
- 2285 6. Sarmanova, A., Morris, T. & Lawson, D. J. Population stratification in GWAS meta-
2286 analysis should be standardized to the best available reference datasets. *bioRxiv*
2287 2020.09.03.281568 (2020) doi:10.1101/2020.09.03.281568.
- 2288 7. Galinsky, K. J., Loh, P.-R., Mallick, S., Patterson, N. J. & Price, A. L. Population
2289 Structure of UK Biobank and Ancient Eurasians Reveals Adaptation at Genes
2290 Influencing Blood Pressure. *Am. J. Hum. Genet.* **99**, 1130–1139 (2016).
- 2291 8. Patterson, N. *et al.* Large-scale migration into Britain during the Middle to Late
2292 Bronze Age. *Nature* (2021) doi:10.1038/s41586-021-04287-4.
- 2293 9. Sikora, M. *et al.* The population history of northeastern Siberia since the Pleistocene.
2294 *Nature* **570**, 182–188 (2019).
- 2295 10. Shinde, V. *et al.* An Ancient Harappan Genome Lacks Ancestry from Steppe
2296 Pastoralists or Iranian Farmers. *Cell* **179**, 729–735.e10 (2019).
- 2297 11. Hofmanová, Z. *et al.* Early farmers from across Europe directly descended
2298 from Neolithic Aegeans. *Proc. Natl. Acad. Sci. U. S. A.* **113**, 6886–6891 (2016).
- 2299 12. Feldman, M. *et al.* Late Pleistocene human genome suggests a local origin for
2300 the first farmers of central Anatolia. *Nat. Commun.* **10**, 1218 (2019).
- 2301 13. Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *J. Mach. Learn.*
2302 *Res.* **12**, 2825–2830 (2011).
- 2303 14. Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing
2304 in Python. *Nat. Methods* **17**, 261–272 (2020).

2305

2306

2307 3i) Building a population history model of Europeans and using
2308 it to assign local ancestry to all haplotypes in modern and
2309 ancient European samples

2310

2311

Alice Pearson¹, Richard Durbin^{1,2}

2312

2313

¹Department of Genetics, University of Cambridge, UK.

2314

²Wellcome Sanger Institute, Wellcome Genome Campus, UK.

2315

2316 Introduction

2317

We built a quantitative admixture graph model that represents the major ancestry flows

2318

contributing to modern European genomes over the last 50,000 years. Within this model we

2319

placed chromosome sampling points from appropriate populations and times that resemble

2320

European and West Asian samples from the imputed present aDNA dataset. Based on the

2321

model, we developed an approach using genealogical nearest neighbours ¹ in tree

2322

sequences inferred using the Relate program ² to estimate the path backwards in time

2323

through the population structure taken by each modern and ancient haplotype at each

2324

position in the genome. Throughout we evaluate our methods by comparing results from the

2325

real data to those from simulations of the model.

2326 Method and Results

2327 Model of population structure

2328

We first filtered the present aDNA dataset to include 1015 ancient genomes most relevant to

2329

modern European genetic history. This included samples that have a West Eurasian

2330

archeological location and lying on EHG-WHG-CHG-Farmer clines in Principle Component

2331

Analysis and excluded very old West Eurasian samples and archaics. These were merged

2332

with all 503 present-day European 1000 Genomes Project samples to make up our dataset,

2333

totalling 1518 diploid samples. Figure **S3i.1** shows a schematic of our model that describes

2334

the evolution of population structure in Europe during the last 50 ky. We used parameters

2335

subjectively estimated from Principal Component Analysis and Jones *et al.* ³, and

2336

constructed a simulator for this model in msprime (a coalescent simulator) ⁴. Shortly after

2337 the expansion of anatomically modern humans into Eurasia there is an early population split
2338 45 kya between the Northern Europeans (NE) who continued travelling north west into
2339 Europe and West Asians (WA) who stayed more locally in the Levant and South Caucasus
2340 area. The WA population then splits to form the Caucasus Hunter-Gatherers (CHGs) and the
2341 Anatolian Farmers (Ana) 24kya. 18kya the Western Hunter-Gatherers (WHGs) and Eastern
2342 Hunter-Gatherers (EHGs) begin to diverge within the NE. At that point, the four genetic
2343 strands that make up present day European ancestry are distinct in our model and
2344 subsequently admixture between these components describes the formation of modern
2345 European gene pool; we show the formation of the Neolithic Farmers (Neo) from admixture
2346 between WHGs and the Ana 6kya, the formation of the Yamnaya Steppe people (Yam) from
2347 admixture of EHG and CHGs during the early Bronze Age 5.4 kya and finally the formation
2348 of the Bronze Age gene pool 4.2kya as a two-way mix between the Yam and the Neo. From
2349 the Bronze Age up to present-day, there is simply exponential growth in population size.

2350 The sampling distribution in the model is based on the numbers and average ages of sample
2351 groups in the present aDNA data set, where Iron Age, Viking and Late Antiquity samples
2352 younger than 2500 years BP are grouped with the present day genomes and those older are
2353 grouped with Bronze Age samples.

2354

2355 Local Ancestry Using Tree Sequences:

2356 In the absence of recombination, a set of sample chromosomes from different genomes are
2357 related to each other by a single tree, identifiable by the sharing of derived mutations. A
2358 recombination event between any two chromosomes changes the tree topology. As you
2359 move from one end of the sample chromosomes to the other, a sequence of changing trees
2360 is observed, each encoding the genealogy of a segment of DNA (sample haplotypes) and
2361 each tree change reflecting one or more recombination events. The Relate software ² aims
2362 to infer the true underlying tree sequence from genotype data. We worked with three types of
2363 tree sequence data; tree sequences simulated directly from the model using msprime (model
2364 simulated), tree sequences inferred by Relate from genotype data simulated from the model
2365 (Relate simulated) and tree sequences inferred by Relate from the real genotype data
2366 (Relate MesoNeo).

2367 With ancient samples and admixture events, the first coalescence alone is insufficient for
2368 understanding the full ancestry of a given haplotype. This is because the first coalescence
2369 event may occur at a time younger than the age of some sampled groups, in which case the
2370 older sampled individuals could not be found as the closest relatives and therefore the true
2371 local ancestry of haplotypes may not be correctly established. Similarly, with some sampled
2372 populations formed via admixture of other sampled populations, the closest relatives to a
2373 haplotype may by chance be from the admixed population alone even if the age of first
2374 coalescence is old enough to capture all sampled groups, in which case again the true local
2375 ancestry of a haplotype may not be found. The path that a haplotype takes backwards in
2376 time from the present day to the root of the model is more informative about its local ancestry
2377 as its relationship to all relevant historical and admixing populations is established. If we
2378 assume an infinite-sites model, all alleles only appear once in history by mutation and all
2379 take a single path from the present day to the root.

2380 The method involves taking a focal haplotype and its marginal tree from the tree sequence
2381 describing its genealogical relationship to all other sample haplotypes. From that haplotype
2382 we traverse up the tree, jumping to successive parent nodes towards the root i.e parent,
2383 grandparent, great-grandparent etc. We want to identify what population each of these nodes
2384 (ancestors) occurs in, which in turn describes what path the haplotype has taken backwards
2385 in time. In simulated tree sequences, we can simply obtain (using tskit) the population
2386 identity of each node but in tree sequences inferred with Relate this is not possible. We
2387 therefore adapted the concept of Genealogical Nearest Neighbours ¹. At each node we
2388 record its age and the number of reference samples belonging to each ancestral group in the
2389 leaves below that node as a proportion of all the reference samples in the leaves below that
2390 node (not including the focal leaf itself or leaves seen at previously analysed nodes further
2391 down the tree). I refer to the distribution of reference ancestry proportions and age at a node
2392 as GNN_x (Genealogical Nearest Neighbours at *x*) where *x* is the *x*th node examined towards
2393 the root. The ordered collection of all GNN_x distributions of all *x* nodes examined during a
2394 tree traversal describes the population identity of each node and therefore the path that the
2395 focal haplotype has taken to the root. The key is to consider the GNNs together, not
2396 independently.

2397 We therefore first define the set of ancient reference samples from our dataset that are
2398 genetically diagnostic of each of the 7 ancestral groups. These samples form tight clusters

2399 within broader ancestry clouds in the PCA. In order to assign paths to millions of sample
 2400 lineages we implemented a supervised machine learning method using the Python Keras
 2401 package with a TensorFlow backend ⁵. The input to the network is the ordered collection of
 2402 GNN distributions plus the node age for the first five informative nodes traversed towards the
 2403 root from a single sample haplotype, configured as a 5 x 8 matrix. Informative nodes are
 2404 those that have at least one leaf from the reference set of ancient samples. If we reach the
 2405 root of the tree in less than five informative nodes, then the remaining rows of the matrix are
 2406 filled with -15 as padding. The output of the network is a numerical label 0-6 characterising
 2407 the path taken by the haplotype backwards in time. These paths can be seen clearly in
 2408 Figure S3i.1 and are described in Table S3i.1. Given the sampling distribution, there is not
 2409 always an informative node in the area of the model where the four ancestries are distinct,
 2410 200-600 generations ago. Labels 5, 6 and 0 capture these situations where we have limited
 2411 information to assign the full path and instead only a partial path or none at all.

Path label	Populations of inheritance
0	None can be determined
1	Neolithic Farmers/Anatolian/West Asian
2	Yamnaya/CHG/West Asian
3	Neolithic Farmers/WHG/Northern European
4	Yamnaya/EHG/Northern European
5	Northern European
6	West Asian

2412 **Table S3i.1:** The populations that carried a haplotype through time by inheritance, ordered
 2413 backwards in time that make up each path.

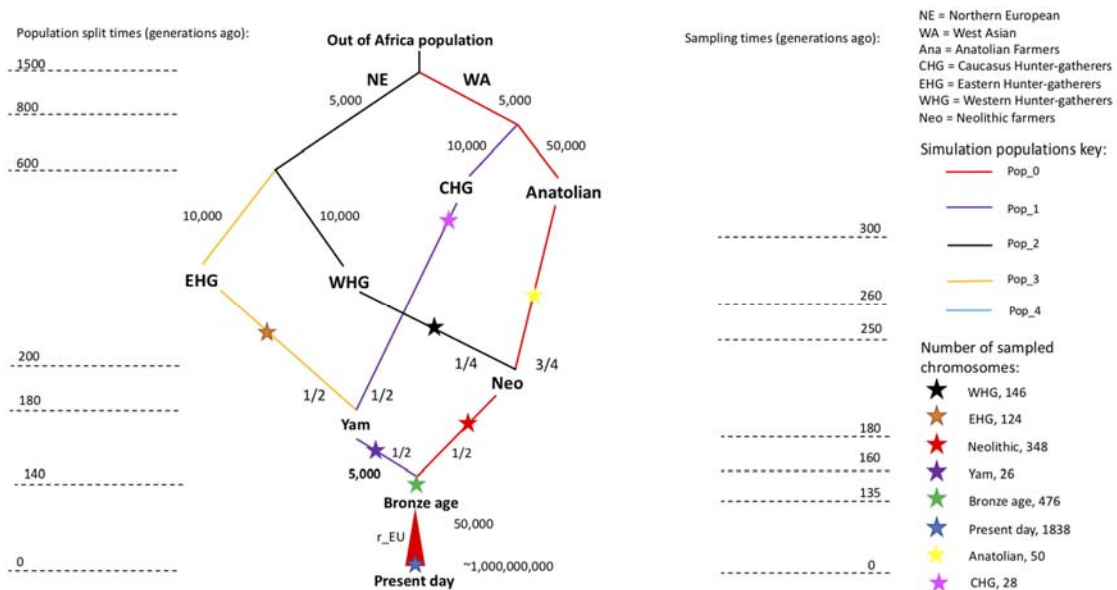
2414 A large amount of training data can be generated by simulation of tree sequences from the
 2415 model. We simulated a tree sequence and corresponding VCF file of chromosome 3 and
 2416 applied Relate to the VCF file to obtain a Relate simulated tree sequence. We then extracted
 2417 GNN distributions, traversing from all Present-day, Bronze Age, Neo and Yam samples at
 2418 evenly spaced trees from the Relate simulated tree sequence. Path labels at corresponding
 2419 sites were extracted from the model simulated tree sequence. WHG, EHG, CHG, and Ana
 2420 can be given paths 1-4 from their population identity alone and therefore are not involved in
 2421 the network training. Our total training set consisted of 4,000,000 GNN matrices and labels.

2422 85% of true labels were full paths 1-4. The network was trained using a categorical cross-
 2423 entropy function, Adam optimisation and a batch size of 30. Training took place over 30
 2424 epochs. For testing we generated another Relate Simulated tree sequence and extracted
 2425 1,000,000 Relate GNNs and simulated true labels. The network displayed 75% accuracy.
 2426 Figure **S3i.B** is a confusion matrix comparing the classed labels to the true labels.

2427 Relate tree sequences were inferred for all chromosomes of our merged dataset, using a
 2428 mutation rate of 1.25×10^{-8} per bp per generation and fine-scale human recombination maps.
 2429 We assigned paths for all samples at every site using the trained network. We annotated the
 2430 merged VCF files with the path assignments as a FORMAT tag with ID=AP standing for
 2431 Ancestral Path.

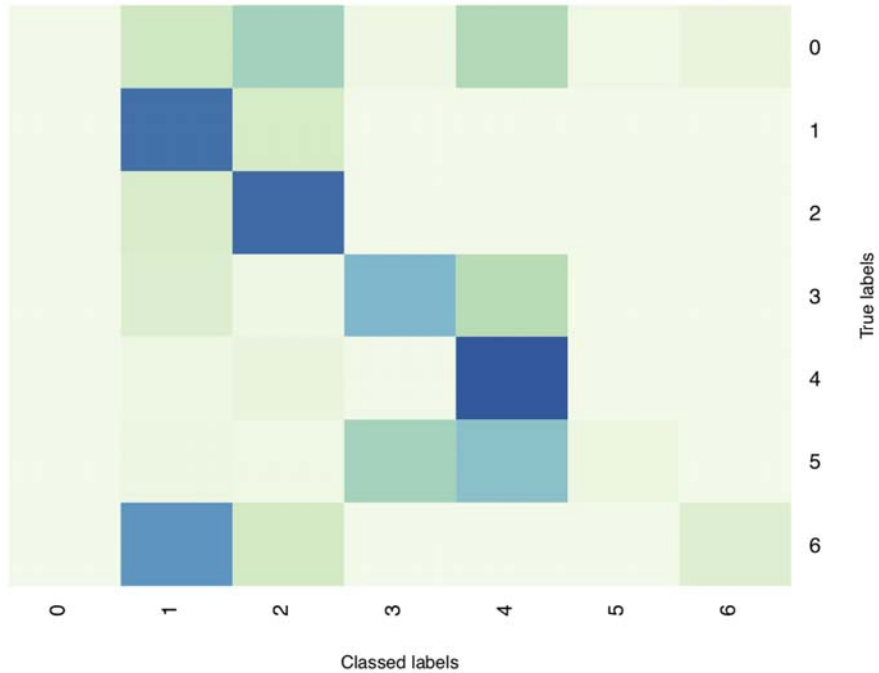
2432 Admixture fractions for hybrid groups can be estimated from the local ancestry assignments
 2433 as the proportion of sites taking each of the relevant paths out of all sites in all samples of a
 2434 hybrid group. We calculated this genome-wide from the VCF files containing the new AP
 2435 format tag. The results are shown in Figure **S3i.3**.

2436 Figures:



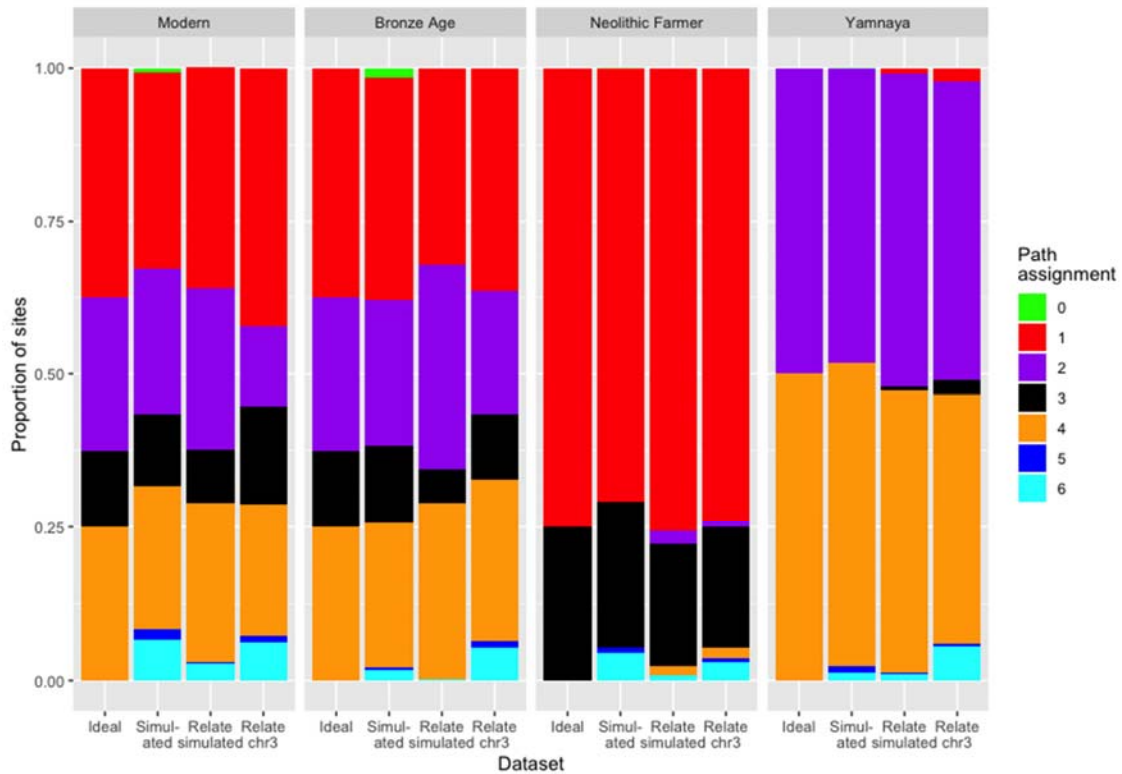
2437
 2438
 2439 **Figure S3i.1:** A schematic of the model of population structure in Europe that was built in
 2440 msprime. Moving down the figure is forwards in time and the population split times and
 2441 admixture times are given in generations ago. Coloured lines represent the four populations
 2442 declared in the simulation that extend through time. Sampled populations and times are
 2443 marked with a star and the number of chromosomes sampled is given in the key.

2444
2445
2446
2447



2448
2449
2450
2451
2452

Figure S3i.2: Heatmap of the confusion matrix between the true labels and classed labels for Relate simulated test data. This is normalised by the sum of rows i.e of the true labels, how many the classifier correctly predicted.



2453
 2454
 2455
 2456
 2457
 2458
 2459
 2460
 2461
 2462
 2463
 2464
 2465
 2466
 2467

Figure S3i.3: Admixture fractions displayed as bars divided by the proportion of all sites assigned each path 0-6. For each admixed population there are four bars from four different tree sequences: *Ideal* is our best guess at the admixture fractions given previously published results and shows the fractions we used as input in msprime simulation. *Simulated* are the proportions of sites assigned each path extracted from the msprime simulated tree sequences which we used to train the neural, therefore partial paths are included where the node distribution means full paths cannot be determined. *Relate simulated* are the proportions of sites classed as each path in a Relate tree sequence inferred from simulated data, when GNNs are extracted from this tree sequence and classified as paths by the neural network. *Relate chr3* are the proportions of sites classed as each path in a Relate tree sequence inferred from the real (ancient genomic) chromosome 3 data, when GNNs are extracted from this tree sequence and classified as paths by the neural network.

2468 **References**

2469 1. Kelleher, J. *et al.* Inferring whole-genome histories in large population datasets. *Nat.*
 2470 *Genet.* **51**, 1330–1338 (2019).
 2471 2. Speidel, L., Forest, M., Shi, S. & Myers, S. R. A method for genome-wide genealogy
 2472 estimation for thousands of samples. *Nat. Genet.* **51**, 1321–1329 (2019).
 2473 3. Jones, E. R. *et al.* Upper Palaeolithic genomes reveal deep roots of modern

- 2474 Eurasians. *Nat. Commun.* **6**, 1–8 (2015).
- 2475 4. Kelleher, J., Etheridge, A. M. & McVean, G. Efficient Coalescent Simulation and
2476 Genealogical Analysis for Large Sample Sizes. *PLoS Comput. Biol.* **12**, e1004842
2477 (2016).
- 2478 5. Chollet, F. *et al.* Keras. <https://keras.io> (2015).

4) Natural selection and trait evolution

2479

2480

4a) Estimating allele frequency trajectories of trait-associated variants

2481

2482

2483 Evan K. Irving-Pease¹, Aaron J. Stern², Rasmus Nielsen^{1,2}, Fernando Racimo¹

2484

2485 ¹Lundbeck Foundation GeoGenetics Centre, GLOBE Institute, University of Copenhagen,
2486 Copenhagen, Denmark

2487

²Department of Integrative Biology, University of California, Berkeley

2488

2489 Introduction

2490 Genome-wide association studies (GWAS) of present-day human populations have identified
2491 large numbers of genetic variants associated with complex traits. However, the extent to
2492 which these variants have been under positive selection during recent human evolution is
2493 unclear. We aimed to model the allele frequency trajectories and selection coefficients of
2494 GWAS variants through time, using genomic data from both present-day populations and
2495 ancient individuals sampled across West Eurasia during the Holocene. We used the software
2496 *CLUES*¹, which supports inference of allele frequency trajectories from marginal trees
2497 sampled from a reconstruction of the ancestral recombination graph (ARG)² for a set of
2498 genomic sequences, in combination with genotype likelihoods from serially sampled ancient
2499 DNA (aDNA).

2500

2501 To account for population structure in our samples, we applied a novel chromosome painting
2502 technique ([Supplementary Note S3i](#)). This technique is based on inference of a sample's
2503 nearest neighbours in the marginal trees of an ARG that contains labelled individuals. In our
2504 case, the labelling corresponds to ancestral populations that predate the main episodes of
2505 admixture in West Eurasia ([Supplementary Note S3i](#)). This method allows us to accurately
2506 assign ancestral population labels to haplotypes found in both ancient and present-day
2507 individuals. By conditioning our selection analyses on these haplotype backgrounds, we can
2508 infer the selection trajectories of GWAS risk alleles in a manner that is approximately
2509 invariant to change in the admixture proportions through time. These ancestry specific allele
2510 trajectories reveal many novel aspects about the dynamic interplay between selection and
2511 admixture in West Eurasia throughout the Holocene.

2512 Methods

2513 The computational pipeline to perform all analyses was written in the `snakemake` workflow
2514 management system ³. For a full list of all the software and versions used, see [Table S4a.1](#).
2515 The directed acyclic graph (DAG) of the computational pipeline is shown in [Figure S4a.1](#). All
2516 pipeline code, custom scripts and a `conda` environment to replicate the analyses are
2517 available in the GitHub repository (https://github.com/ekirving/mesoneo_paper).

2518 SNP Ascertainment

2519 GWAS SNPs

2520 We ascertained -a list of GWAS targets by downloading version v1.0.2 (2020-06-04) of the
2521 NHGRI-EBI GWAS Catalog ⁴; containing 187,403 GWAS associations for 3,735 traits. To
2522 account for the varying significance thresholds used in the 4,007 published studies included
2523 in the catalogue, we restricted our analysis to SNPs with a genome-wide significance
2524 threshold of $p < 5e-8$. We further filtered the catalogue to retain only single-nucleotide
2525 polymorphisms (SNPs) with a valid dbSNP Reference SNP identifier (rsID); resulting in
2526 121,795 GWAS associations for 70,224 rsIDs. For each of the retained associations, we
2527 retrieved the trait ontology hierarchy by querying the EMBL-EBI Ontology Lookup Service ⁵.
2528 We then queried the Ensembl REST API ⁶, to retrieve metadata about each rsID; including
2529 chromosome and position in the GRCh37 assembly, ancestral allele, and nearest genes.

2530 Control SNPs

2531 To determine the extent to which GWAS variants are enriched for selection, we paired each
2532 GWAS SNP with a unique “Control SNP”. Control SNPs were ascertained by selecting all
2533 biallelic SNPs within the imputed dataset ([Supplementary Note 2](#)) and excluding any that fell
2534 within +/- 50 kb of a GWAS SNP or a gene region. Gene annotations for GRCh37 were
2535 downloaded from Ensembl (release 87) ⁷. Control SNPs were grouped into bins based on
2536 their derived allele frequency (DAF), rounded to the nearest 1%, and paired randomly
2537 (without replacement) with GWAS SNPs in the same chromosome and DAF bin.

2538 Simulated Neutral SNPs

2539 To measure the effects of demography on the modelled allele frequency trajectories, we
2540 frequency paired GWAS SNPs with neutral SNPs, simulated under the demography used to
2541 train the chromosome painting model ([Supplementary Note 3i](#)). Neutral simulations were
2542 performed with *msprime* ⁸, for genomes of length 198 Mbp (i.e., the approximate length of
2543 chr3), with sample sizes and ages broadly equivalent to those in the empirical aDNA dataset.

2544

2545 1000G ARG

2546 We built genome-wide genealogies for all samples in the 1000 Genomes Project (1000G)
2547 Phase 3 release ⁹ using the software *Relate* (1.1.3) ².

2548 Data pre-processing

2549 Prior to inference of the ARG, we converted VCF files into HAPS format, removed all non-
2550 biallelic SNPs, polarised SNPs against the ancestral allele calls from the Ensembl Compara
2551 71 database (ens-staging2:3306) ¹⁰, filtered sites using the 1000 Genomes StrictMask
2552 (20140520) and generated SNP annotations using *Relate*.

2553 ARG Inference

2554 We jointly inferred genome-wide genealogies for all 1000G Phase 3 samples, using *Relate*,
2555 assuming a mutation rate of 1.25e-8 and an Ne of 30,000. From this ARG, we extracted
2556 subtrees containing samples belonging to three European (EUR) populations: (i) Finnish in
2557 Finland (FIN); (ii) British in England and Scotland (GBR); and (iii) Toscani in Italia (TSI). We
2558 used these subtrees to jointly reinfer branch lengths and to infer a population size history for
2559 the EUR metapopulation. Lastly, we remapped all SNPs which had been pruned during
2560 inference of the ARG onto the branch-length calibrated EUR subtrees.

2561 Modifications to CLUES

2562 Here we describe several modifications made to *CLUES* (see the GitHub wiki page for more
2563 information <https://github.com/35ajstern/clues/wiki>).

2564

2565 ARG sampling using Relate

2566 Instead of sampling ARGs using *ARGweaver* ¹¹, we sample ARGs using *Relate* ² for
2567 scalability reasons. *Relate* differs from *ARGweaver* in that it assumes a continuous-time
2568 coalescent process (vs discrete-time for *ARGweaver*); hence we modified the hidden Markov
2569 model (HMM) used by *CLUES* to (1) take time steps every generation, vs over a smaller
2570 number (~10-50) of timesteps; and (2) within time steps, the probability density of
2571 coalescence is calculated using the approach of ¹² and ¹³, vs the discrete-time lines-of-
2572 descent approach used in ¹¹ and ¹.

2573

2574 Ancient DNA samples

2575 We also introduced a new feature that allows the user to specify a time series of ancient
2576 genotype likelihoods which are incorporated into the HMM. We incorporate these samples
2577 by, for a given timestep t , including (i.e. multiplying by) a Binomial($n=2$, $p=X_t$) emission
2578 probability for each ancient sampled during timestep t in the HMM, where X_t is the latent
2579 allele frequency during timestep t . In this particular application, we supplanted genotype
2580 likelihoods with genotype posterior probabilities (which should be identical under a uniform
2581 prior); we confirmed through tests that this did not yield any systematic biases.

2582 Selection Analysis

2583 CLUES with Modern 1000G data

2584 For each modelled SNP, we used `Relate` to draw 100 samples from the MCMC posterior
2585 distribution of trees at that locus. Trees were sampled assuming a mutation rate of $1.25e-8$
2586 ^{2,14} and using the population size history from the EUR calibrated subtrees.

2587

2588 We ran `CLUES` to infer allele frequency trajectories and selection coefficients from modern
2589 1000G data, using: (i) the 100 sampled trees from `Relate` (`--times`); (ii) the inferred EUR
2590 population size history (`--coal`); (iii) with trajectories polarised the by the derived allele (`--`
2591 `A1`); (iv) a terminal frequency equal to the DAF of each SNP in EUR (`--popFreq`); and (v)
2592 constrained selection to a single epoch spanning the last 15,000 years (`--timeBins`).

2593

2594 We ran these models for all GWAS and Control SNPs present in the imputed dataset for
2595 which `Relate` was able to confidently map a mutation to the inferred trees ($n=73,232$).

2596 CLUES with aDNA Time Series

2597 We also ran `CLUES` in an alternative mode, excluding the modern ARG data, and replacing
2598 them with aDNA time series data, using: (i) the time series of aDNA genotype probabilities (`-`
2599 `-ancientSamps`) (ii) the inferred EUR population size history (`--coal`); (iii) a terminal
2600 frequency equal to the DAF of each SNP in EUR (`--popFreq`); and (iv) constraining
2601 selection to a single epoch spanning the last 15,000 years (`--timeBins`).

2602

2603 We ran these models for all GWAS and Control SNPs present in the imputed dataset,
2604 irrespective of their mappability in the `Relate` analyses ($n=73,988$), as well as for all
2605 Simulated SNPs ($n=11,665$).

2606 CLUES with aDNA Ancestral Paintings

2607 We also ran four additional models for each GWAS and Control SNP, in which we
2608 conditioned the time series of aDNA genotype probabilities on one of four specific ancestral
2609 haplotype pathways (Supplementary Note S3i):

2610

- 2611 1. ANA (Anatolian Farmers -> Neolithic)
- 2612 2. CHG (Caucasus Hunter-gatherers -> Yamnaya)
- 2613 3. WHG (Western Hunter-gatherers -> Neolithic)
- 2614 4. EHG (Eastern Hunter-gatherers -> Yamnaya)

2615

2616 In a minority of cases, the chromosome painting model assigned a haplotype to one of two
2617 basal pathways: (i) North European ancestry (WHG and EHG); and (ii) West Asian ancestry
2618 (CHG and ANA). In such cases, we included these haplotypes in the conditional analyses for
2619 both downstream pathways.

2620

2621 All other particulars of these models were identical to the earlier aDNA analyses, except that
2622 genotypes were passed to *CLUES* in haploid mode (`--ancientHaps`), even when both
2623 haplotypes in an individual shared the same painting.

2624

2625 For the simulated SNPs, we ran these analyses on two different datasets: (i) a simulation
2626 labelled with the true pathways of each haplotype; and (ii) a simulation in which the pathways
2627 were inferred by the chromosome painting model.

2628

2629 Reference and mapping bias filters

2630 To address issues of mapping biases that may cause artifactual changes in allele frequency
2631 at individual sites, we also constructed a causal model for distinguishing direct effects of age
2632 on allele frequency from indirect effects mediated by read depth, read length, and/or error
2633 rates (Supplementary Note S4b). We then filtered out SNPs in which more than half of the
2634 signal for temporal allele frequency change was driven by non-biological artefacts ($0.5 \geq F_j <$
2635 1.0). Furthermore, we implemented an additional mapping bias test, in which we compared
2636 the inferred present-day frequency of all SNPs based on (i) a *CLUES* model of the aDNA
2637 time-series data which was conditioned on the present-day frequency of each variant in the
2638 three EUR populations (see above); and (ii) a simpler *CLUES* model containing the aDNA
2639 time-series data alone. We then filtered out all SNPs in which the addition of the modern
2640 data both increased the significance of the selection test and resulted in an absolute present-

2641 day frequency difference of > 0.1 between the two models. We also filtered out SNPs in
2642 which the observed pattern of genotypes in modern individuals was inconsistent with the
2643 marginal trees inferred from the surrounding haplotypes, as determined by `Relate`².

2644 Genome-wide selection

2645 We aggregated the results of all *CLUES* models (n=73,988) and converted the likelihood-
2646 ratio scores into p-values, using the chi-squared distribution with one degree of freedom. To
2647 identify genome-wide selection peaks, we used *harvester* (0.1)¹⁵ with a relaxed input filter (`-`
2648 `inlimit 1e-3`) and a minimum peak height equal to the Bonferroni corrected p-value
2649 threshold (`-peak-limit 5.87`) for the number of tests in each grouping, after firstly
2650 applying a positional correction for the sparseness of the SNP ascertainment (because the
2651 published GWAS SNPs had already been pruned for linkage disequilibrium in their original
2652 studies). We then merged any adjacent peaks that were less than 100 kb apart in the same
2653 chromosome.

2654

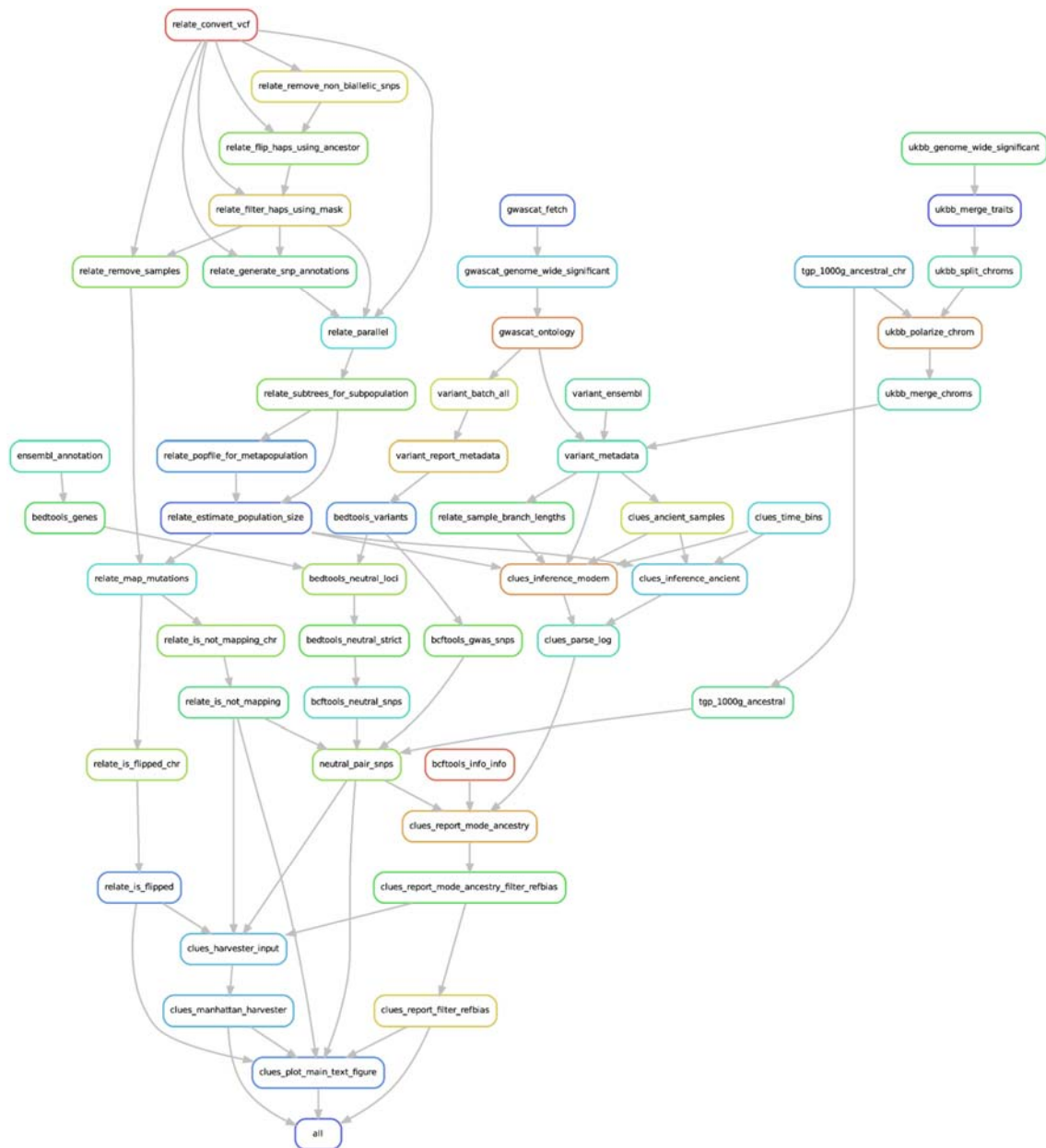
2655 We also obtained a list of putatively selected SNPs and their p-values, inferred in an earlier
2656 aDNA study¹⁶. To ascertain novel sweep loci in our study, we used *harvester* to infer
2657 peaks from the published p-values. Additionally, we used *CLUES* to infer new p-values for all
2658 genome-wide significant SNPs (p<5e-8; n=381) from the earlier study¹⁶, as not all SNPs
2659 found to be significant in that study were present in our GWAS/Control ascertainment.

2660

2661 **Table S4a.1.** Software and versions used in the allele trajectory pipeline

Software	Version	URL	Reference
bcftools	1.10.2	https://github.com/samtools/bcftools	17
bedtools	2.29.2	https://github.com/arg5x/bedtools2	18
biopython	1.76	https://github.com/biopython/biopython	19
clues	36cb7de	https://github.com/35ajstern/clues	1
conda	4.9.0	https://github.com/conda/conda	20
harvester	0.1	https://genomics.ut.ee/en/tools/manhattan-harvester	15
msprime		https://github.com/tskit-dev/msprime	8
numpy	1.17.0	https://github.com/numpy/numpy	21
pandas	1.0.4	https://github.com/pandas-dev/pandas	22

pysam	0.15.3	https://github.com/pysam-developers/pysam	23
python	3.6.7	https://www.python.org	24
r-base	3.6.1	https://www.r-project.org/	25
r-bedr	1.0.7	https://github.com/cran/bedr	26
r-dplyr	0.8.0.1	https://github.com/tidyverse/dplyr	27
r-ggplot2	3.1.1	https://github.com/tidyverse/ggplot2	28
r-ggrastr	0.2.1	https://github.com/VPetukhov/ggrastr	29
r-ggrepel	0.8.2	https://github.com/slowkow/ggrepel	30
r-ggribes	0.5.1	https://github.com/wilkelab/ggribes	31
r-stringr	1.4.0	https://github.com/tidyverse/stringr	32
relate	1.1.3	https://myersgroup.github.io/relate	2
scipy	1.4.1	https://github.com/scipy/scipy	33
snakemake	5.12.3	https://github.com/snakemake/snakemake	3



2663
2664
2665
2666

Figure S4a.1. Directed acyclic graph (DAG) of the allele frequency trajectory analysis pipeline

2667 Results

2668 Selection in 1000G EUR

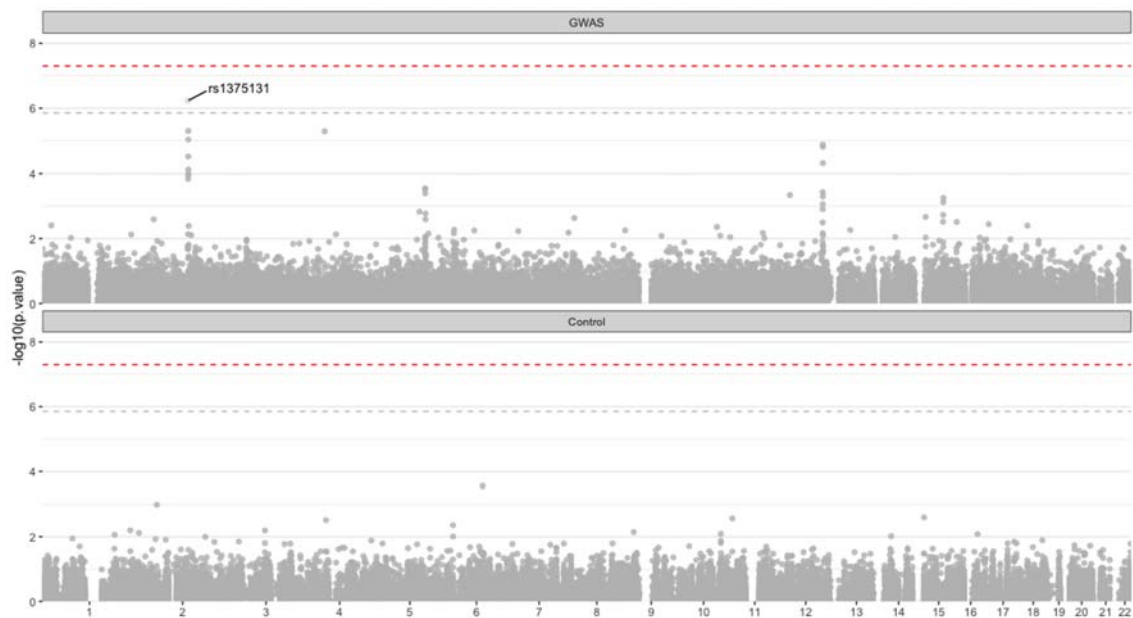
2669 *CLUES* analysis of all GWAS (n=35,592) and Control group SNPs (n=35,592) in the 1000G
2670 Phase 3 populations FIN, GBR, and TSI identified zero genome-wide significant SNPs
2671 ($p < 5e-8$), and only one GWAS group SNP with a p-value below the Bonferroni corrected
2672 significance threshold ($p < 1.40e-6$) (see [Fig. S4a.2](#); Supplemental [Table XV](#)). Despite the

2673 overall lack of genome-wide significant SNPs, the GWAS group was significantly enriched for
2674 evidence of selection when compared to the Control group (Wilcoxon signed-rank test, p-
2675 value<2.2e-16).

2676

2677 The only significant SNP was rs1375131 (chr2:135954797), an intron variant in *ZRANB3*,
2678 which is associated with mosquito bite size³⁴. Non-significant SNPs within the surrounding
2679 peak region include the lactase persistence SNP rs4988235 (*MCM6*; p=9.3e-6), which has
2680 been widely reported as a target of strong selection in West Eurasians^{35,36}. Among the
2681 GWAS SNPs, there was limited evidence for non-significant selection peaks, with partial
2682 overlap between these regions and those previously reported as genome-wide significant in
2683 other studies. Among the Control SNPs, no evidence of selection was identified.

2684



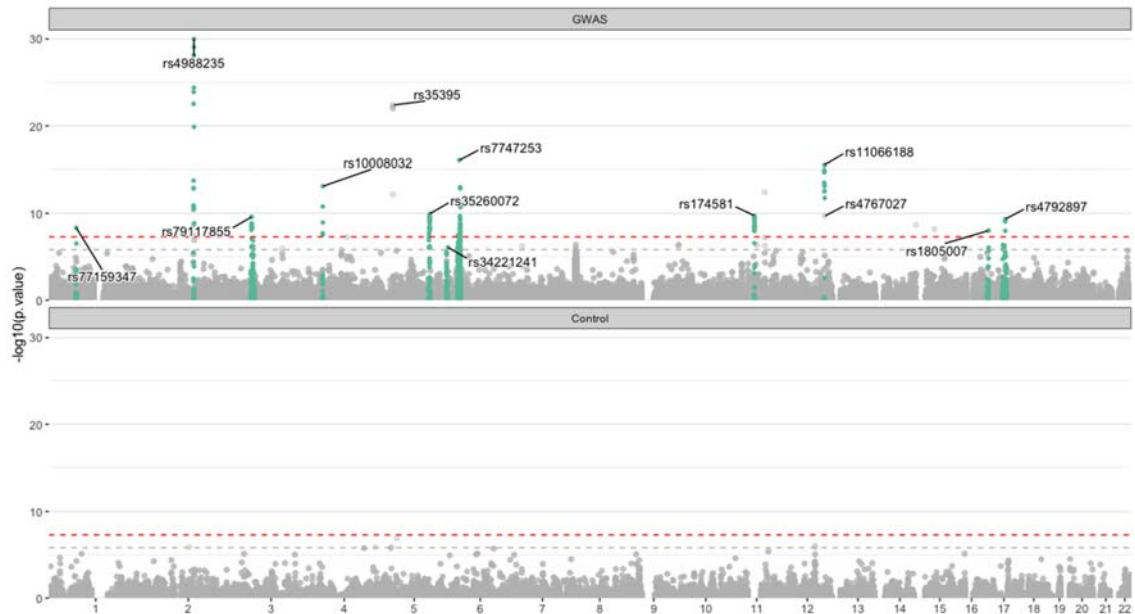
2685

2686 **Figure S4a.2.** Manhattan plot of the p-values from running CLUES on an ARG containing all samples
2687 in FIN, GBR, and TSI from 1000G Phase 3, for (a) GWAS SNPs from the GWAS Catalog; and (b)
2688 Control SNPs, frequency paired with the GWAS SNPs.

2689 Selection in aDNA Time Series

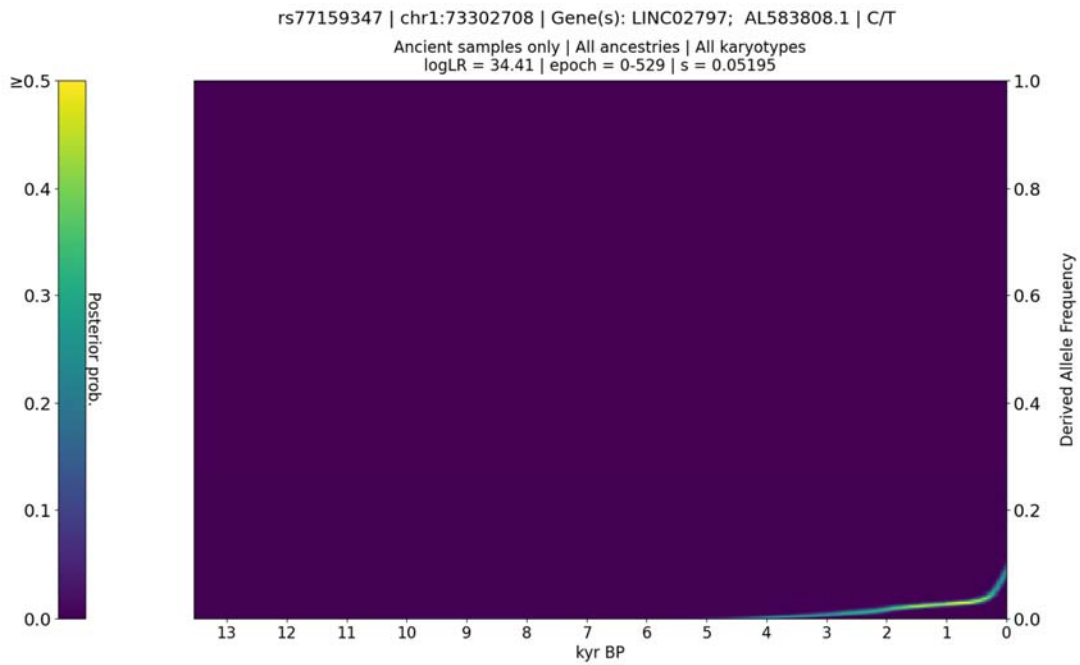
2690 *CLUES* analysis of all GWAS (n=33,323) and Control group SNPs (n=33,323) in the aDNA
2691 time-series dataset identified 127 genome-wide significant SNPs (p<5e-8); 127 in the GWAS
2692 group and 0 in the Control group. Using a Bonferroni corrected significance threshold, we
2693 detected 207 significant SNPs (p < 1.50e-6); 202 in the GWAS group (97.6%) and 5 in the
2694 Control group. Within the GWAS group, we identified 11 Bonferroni corrected significant
2695 selection peaks (see [Fig. S4a.3](#); [Supplementary Table XVI](#)), of which 6 overlapped with

2696 those previously characterised in ¹⁶. No significant selection peaks were detected in the
2697 Control group.
2698



2699 **Figure S4a.3.** Manhattan plot of the p-values from running *CLUES* on an aDNA time series
2700 from all West Eurasian samples in the imputed dataset, for (a) GWAS SNPs from the GWAS
2701 Catalog; and (b) Control SNPs, frequency paired with the GWAS SNPs.
2702
2703
2704

2705 Peak 1: LINC02797, AL583808.1



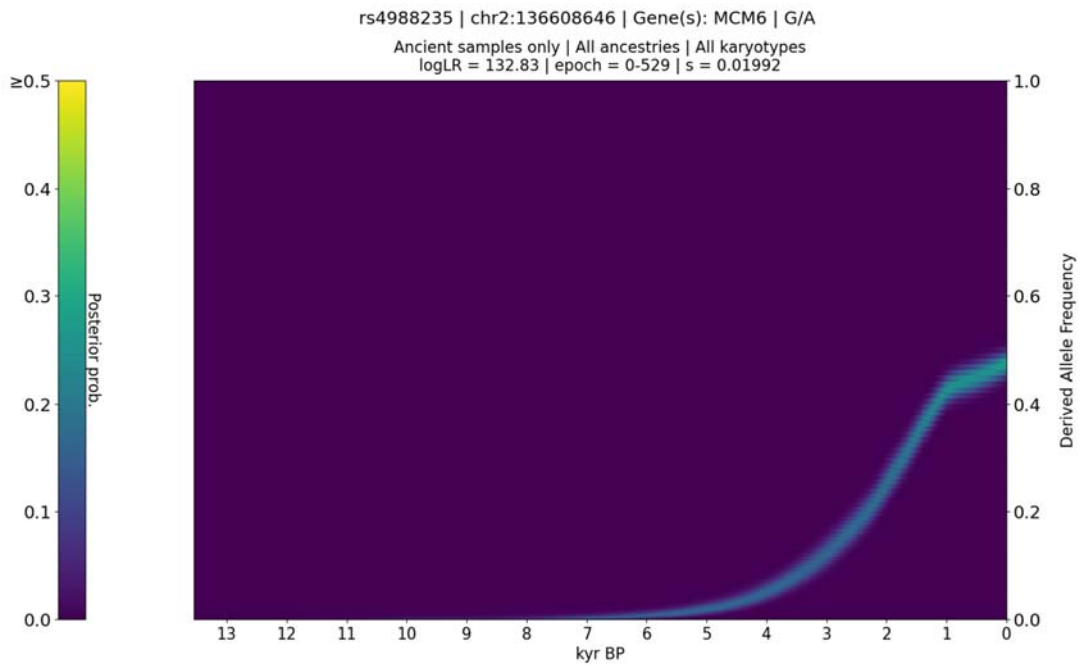
2706 T: Household income (MTAG) (PMID: 31844048)

2707 **Figure S4a.4.** CLUES plot of the aDNA time series analysis for all West Eurasian samples in
2708 the imputed dataset, showing the posterior probability of the derived allele frequency
2709 trajectory for rs77159347, the most significant SNP in the selection peak spanning
2710 chr1:72480859-73978570.

2711 The first peak spanned the region chr1:72480859-73978570, with the most significant
2712 SNP being rs77159347 (*LINC02797*, *AL583808.1*; $p=4.48e-09$; $s=0.052$), associated
2713 with household income (MTAG)³⁷.

2714

2715 Peak 2: MCM6



A: Blood protein levels (PMID: 29875488)
A: Body mass index (PMID: 26426971)
A: Hip circumference (PMID: 25673412)

G: Blood protein levels (PMID: 30072576)
G: Body mass index (PMID: 30108127)

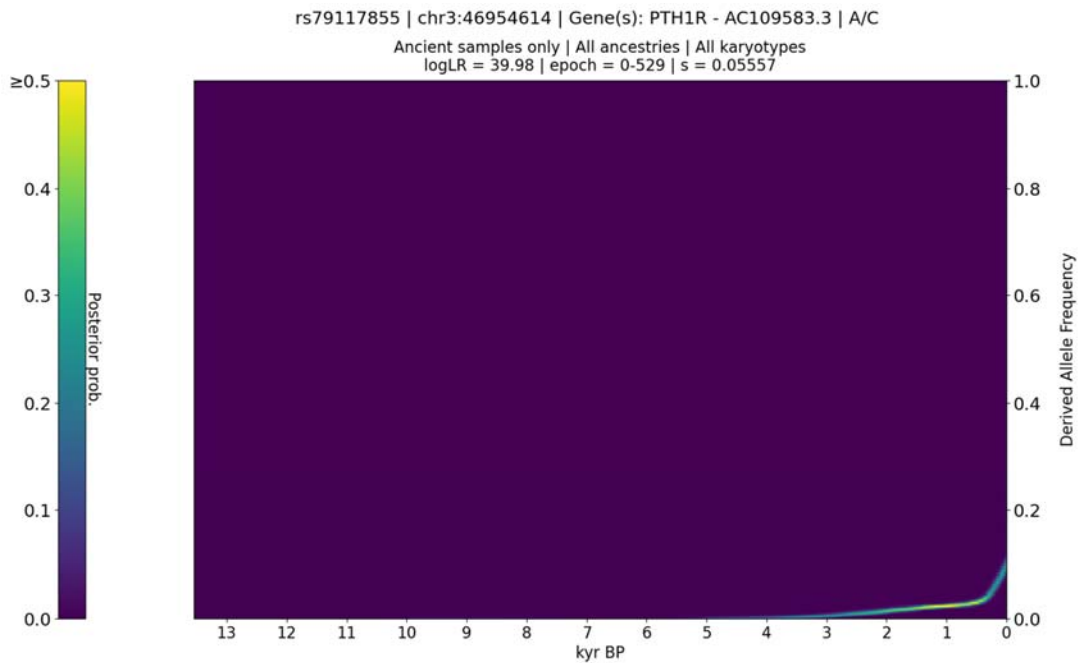
2716

2717 **Figure S4a.5.** CLUES plot of the aDNA time series analysis for all West Eurasian samples in
2718 the imputed dataset, showing the posterior probability of the derived allele frequency
2719 trajectory for rs4988235, the most significant SNP in the selection peak spanning
2720 chr2:135407409-137512400.

2721 The second peak spanned the region chr2:135407409-137512400, with the most
2722 significant SNP being rs4988235 (*MCM6*; $p=9.86e-31$; $s=0.0199$), associated with
2723 lactase persistence; blood protein levels; body mass index; and hip circumference^{38–42}.

2724

2725 Peak 3: PTH1R - AC109583.3



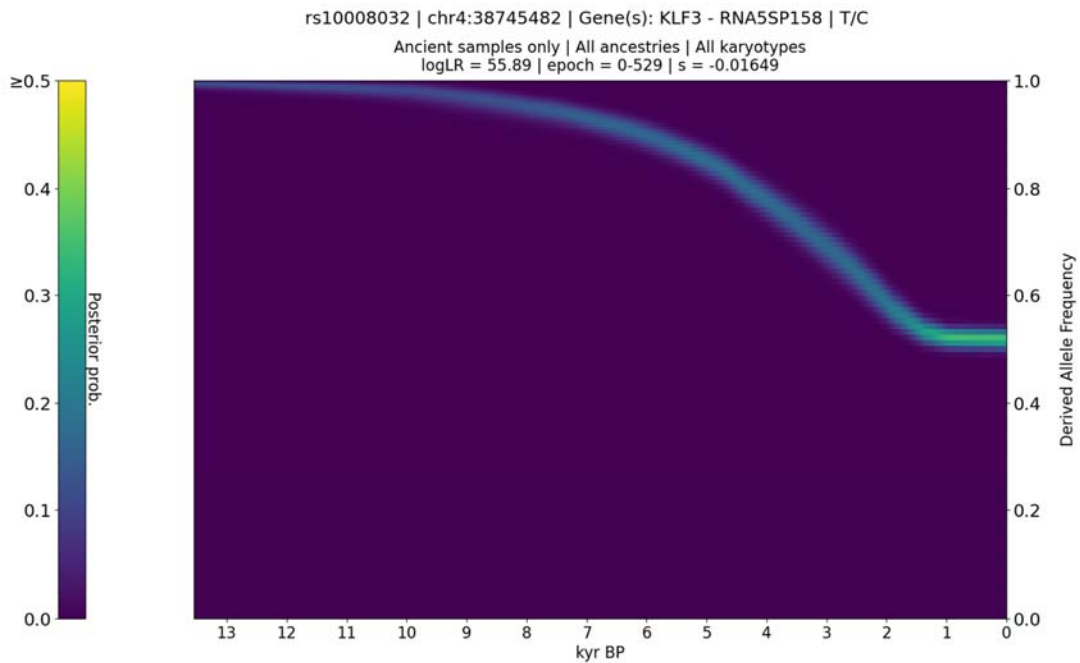
2726 ? : Blood protein levels (PMID: 28915241)

2727 **Figure S4a.6.** CLUES plot of the aDNA time series analysis for all West Eurasian samples in
2728 the imputed dataset, showing the posterior probability of the derived allele frequency
2729 trajectory for rs79117855, the most significant SNP in the selection peak spanning
2730 chr3:44526639-53850000.

2731 The third peak spanned the region chr3:44526639-53850000, with the most significant SNP
2732 being rs79117855 (*PTH1R* - AC109583.3; $p=2.57e-10$; $s=0.0556$), associated with blood
2733 protein levels⁴³.

2734

2735 Peak 4: KLF3 - RNA5SP158



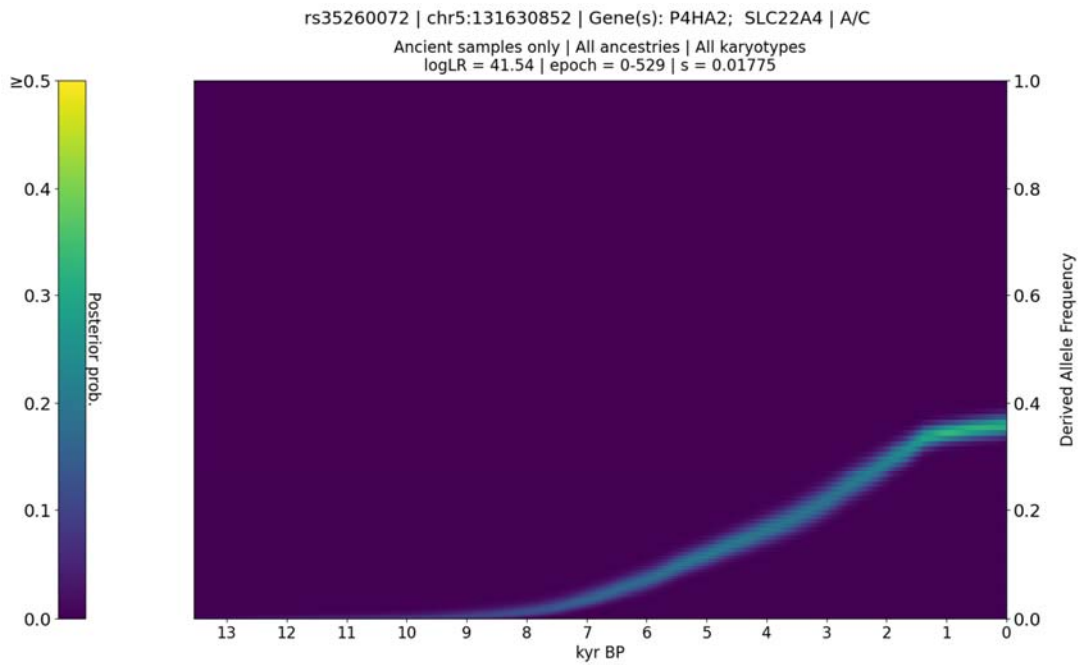
2736 ? : Allergic disease (asthma, hay fever or eczema) (PMID: 29785011)

2737 **Figure S4a.7.** CLUES plot of the aDNA time series analysis for all West Eurasian samples in
2738 the imputed dataset, showing the posterior probability of the derived allele frequency
2739 trajectory for rs10008032, the most significant SNP in the selection peak spanning
2740 chr4:38593259-38815500.

2741 The fourth peak spanned the region chr4:38593259-38815500, with the most significant SNP
2742 being rs10008032 (*KLF3 - RNA5SP158*; $p=7.66e-14$; $s=-0.0165$), associated with allergic
2743 disease (asthma, hay fever or eczema) ⁴⁴.

2744

2745 Peak 5: P4HA2, SLC22A4



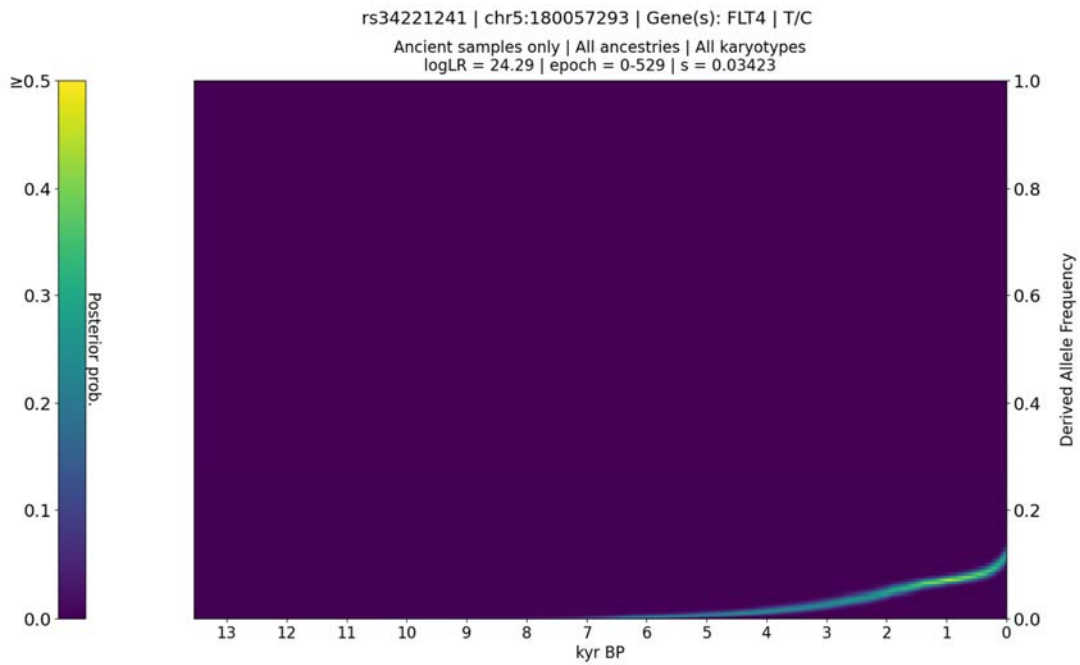
2746 C: Itch intensity from mosquito bite (PMID: 28199695)

2747 **Figure S4a.8.** CLUES plot of the aDNA time series analysis for all West Eurasian samples in
2748 the imputed dataset, showing the posterior probability of the derived allele frequency
2749 trajectory for rs35260072, the most significant SNP in the selection peak spanning
2750 chr5:128016159-132349650.

2751 The fifth peak spanned the region chr5:128016159-132349650, with the most significant
2752 SNP being rs35260072 (*P4HA2*, *SLC22A4*; $p=1.15e-10$; $s=0.0177$), associated with itch
2753 intensity from mosquito bite³⁴.

2754

2755 Peak 6: FLT4



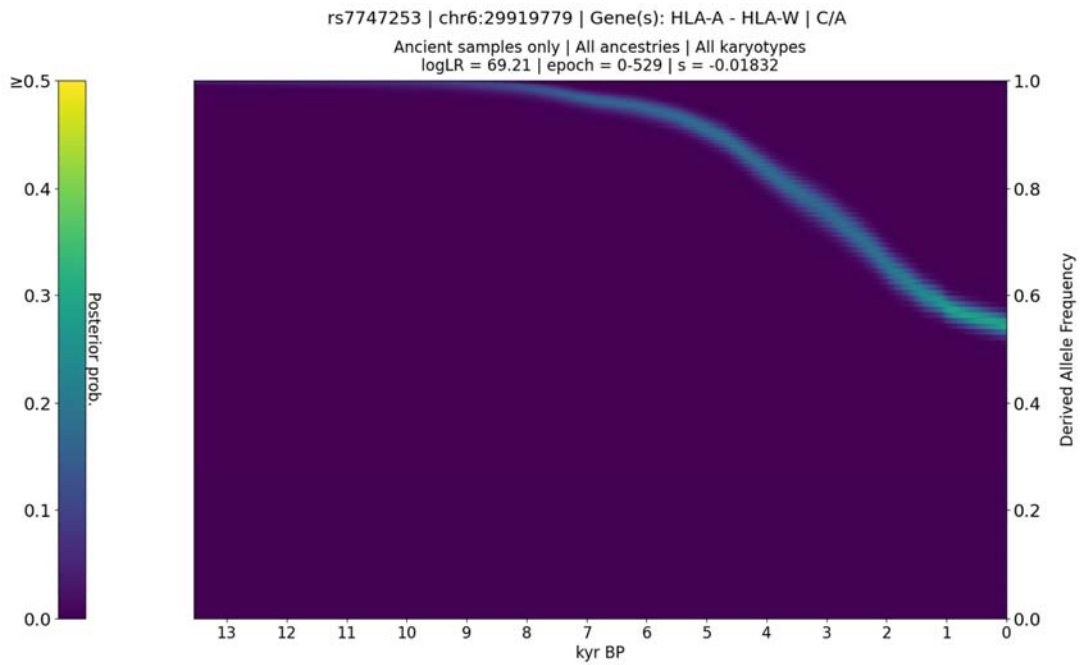
2756 C: Blood protein levels (PMID: 29875488)

2757 **Figure S4a.9.** CLUES plot of the aDNA time series analysis for all West Eurasian samples in
2758 the imputed dataset, showing the posterior probability of the derived allele frequency
2759 trajectory for rs34221241, the most significant SNP in the selection peak spanning
2760 chr5:176653519-180661980.

2761 The sixth peak spanned the region chr5:176653519-180661980, with the most significant
2762 SNP being rs34221241 (*FLT4*; $p=8.29e-07$; $s=0.0342$), associated with blood protein levels
2763 ⁴⁰.

2764

2765 Peak 7: HLA-A - HLA-W



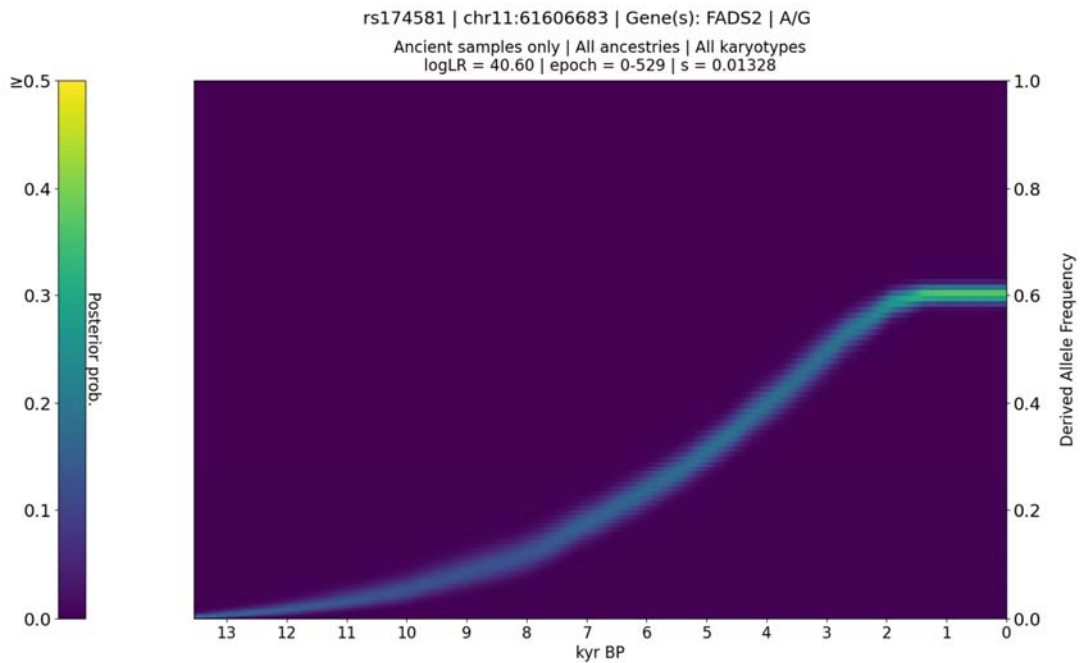
2766 C: Heel bone mineral density (PMID: 30598549)

2767 **Figure S4a.10.** CLUES plot of the aDNA time series analysis for all West Eurasian samples
2768 in the imputed dataset, showing the posterior probability of the derived allele frequency
2769 trajectory for rs7747253, the most significant SNP in the selection peak spanning
2770 chr6:25236639-33535470.

2771 The seventh peak spanned the region chr6:25236639-33535470, with the most significant
2772 SNP being rs7747253 (*HLA-A - HLA-W*; $p=8.86e-17$; $s=-0.0183$), associated with heel bone
2773 mineral density⁴⁵.

2774

2775 Peak 8: FADS2



2776 A: Male-pattern baldness (PMID: 28196072)

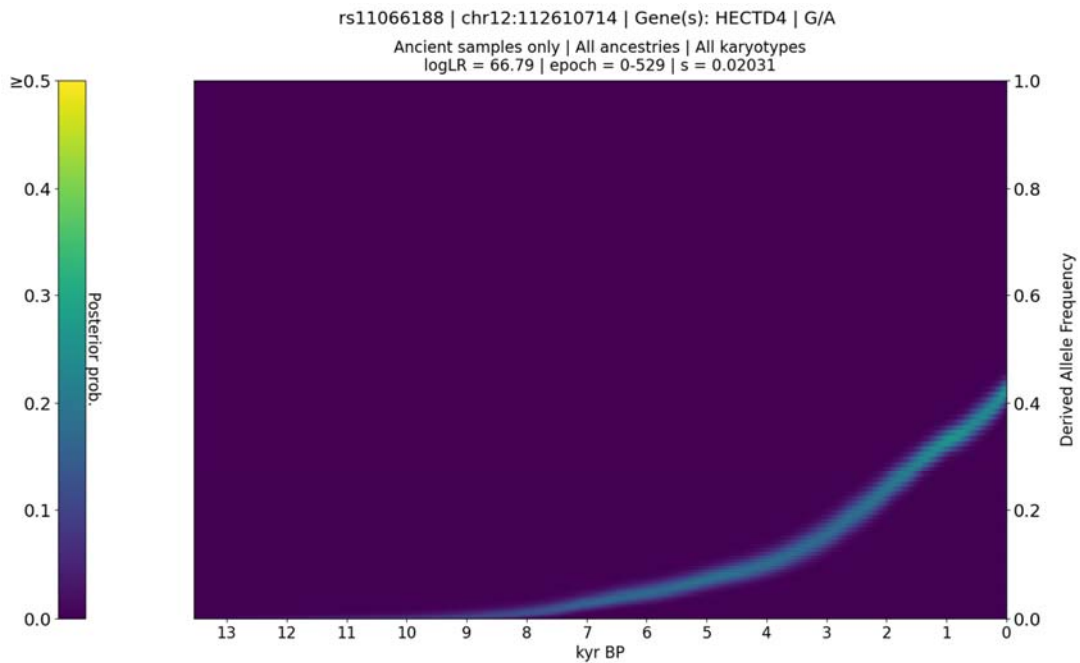
A: Serum metabolite ratios in chronic kidney disease (PMID: 29545352)

2777 **Figure S4a.11.** CLUES plot of the aDNA time series analysis for all West Eurasian samples
2778 in the imputed dataset, showing the posterior probability of the derived allele frequency
2779 trajectory for rs174581, the most significant SNP in the selection peak spanning
2780 chr11:61543499-61706010.

2781 The eighth peak spanned the region chr11:61543499-61706010, with the most significant
2782 SNP being rs174581 (*FADS2*; $p=1.87e-10$; $s=0.0133$), associated with male-pattern
2783 baldness; and serum metabolite ratios in chronic kidney disease ^{46,47}.

2784

2785 Peak 9: HECTD4



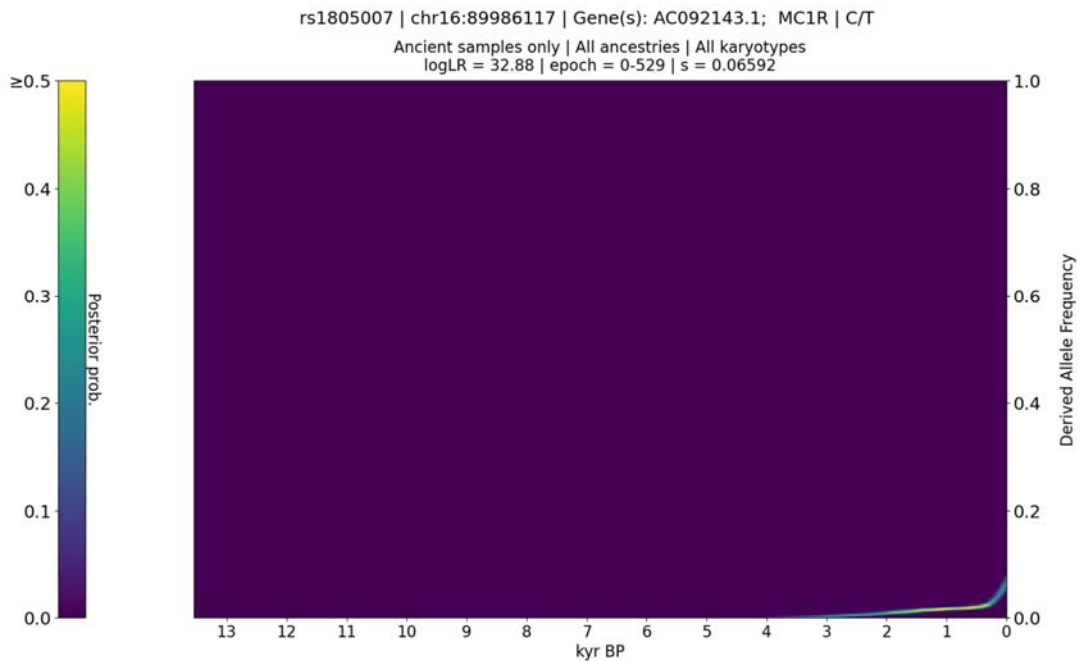
2786 ? : Celiac disease and Rheumatoid arthritis (PMID: 26546613)

2787 **Figure S4a.12.** CLUES plot of the aDNA time series analysis for all West Eurasian samples
2788 in the imputed dataset, showing the posterior probability of the derived allele frequency
2789 trajectory for rs11066188, the most significant SNP in the selection peak spanning
2790 chr12:111833789-113137570.

2791 The ninth peak spanned the region chr12:111833789-113137570, with the most significant
2792 SNP being rs11066188 (*HECTD4*; $p=3.02e-16$; $s=0.0203$), associated with celiac disease
2793 and rheumatoid arthritis⁴⁸.

2794

2795 Peak 10: AC092143.1, MC1R



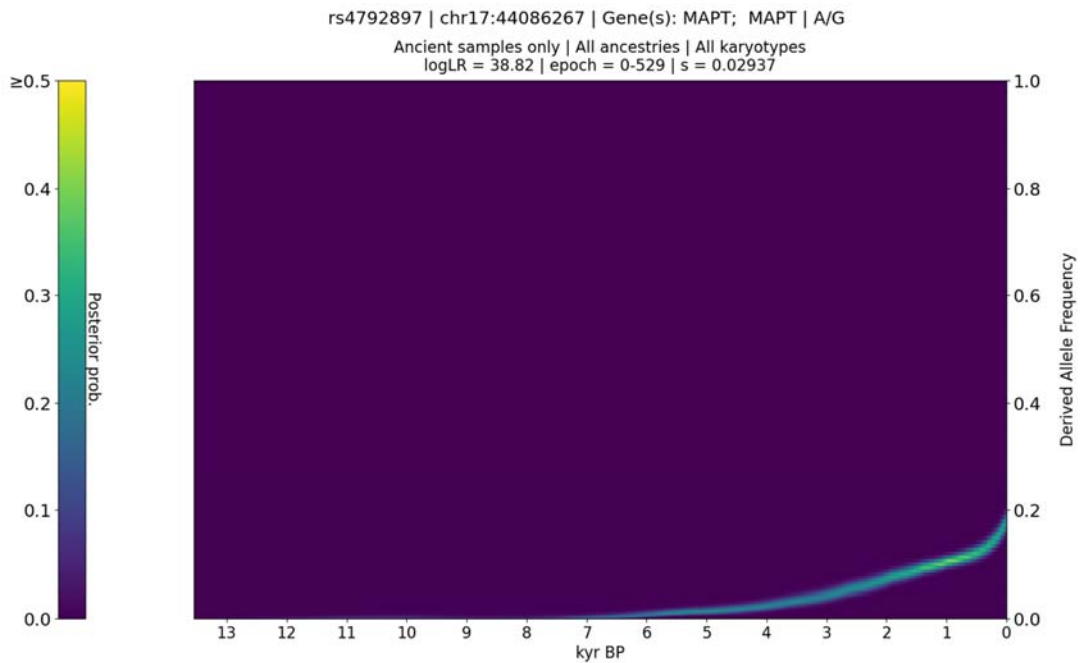
- ?: Balding type 1 (PMID: 30595370)
- ?: Hair morphology traits (PMID: 30166351)
- ?: Skin pigmentation traits (PMID: 30166351)
- A: Brown vs. black hair color (PMID: 30531825)
- C: Hair color (PMID: 23548203)
- C: Non-melanoma skin cancer (PMID: 23548203)
- C: Sunburns (PMID: 23548203)
- C: Tanning (PMID: 23548203)
- T: Basal cell carcinoma (PMID: 31174203; 27539887; 21700618)
- T: Blond vs. brown hair color (PMID: 17952075)
- T: Blond vs. brown/black hair color (PMID: 30531825)
- T: Cutaneous squamous cell carcinoma (PMID: 27424798)
- T: Freckles (PMID: 17952075)
- T: Hair color (PMID: 29662168)
- T: Keratinocyte cancer (MTAG) (PMID: 31174203)
- T: Melanoma (PMID: 28212542)
- T: Red vs non-red hair color (PMID: 17952075)
- T: Red vs. brown/black hair color (PMID: 30531825)
- T: Skin sensitivity to sun (PMID: 17952075)
- T: Squamous cell carcinoma (PMID: 31174203)

2796 **Figure S4a.13.** CLUES plot of the aDNA time series analysis for all West Eurasian samples
 2797 in the imputed dataset, showing the posterior probability of the derived allele frequency
 2798 trajectory for rs1805007, the most significant SNP in the selection peak spanning
 2799 chr16:86009759-90084560.
 2800

2801 The tenth peak spanned the region chr16:86009759-90084560, with the most significant
 2802 SNP being rs1805007 (*AC092143.1*, *MC1R*; $p=9.83e-09$; $s=0.0659$), associated with balding
 2803 type 1; basal cell carcinoma; blond vs. brown hair colour; blond vs. brown/black hair colour;
 2804 brown vs. black hair colour; cutaneous squamous cell carcinoma; freckles; hair colour; hair
 2805 morphology traits; keratinocyte cancer (MTAG); melanoma; non-melanoma skin cancer; red
 2806 vs non-red hair colour; red vs. brown/black hair colour; skin pigmentation traits; skin
 2807 sensitivity to sun; squamous cell carcinoma; sunburns; and tanning^{49–59}.

2808

2809 Peak 11: MAPT



2810 A: Snoring (PMID: 30804565)

2811 **Figure S4a.14.** CLUES plot of the aDNA time series analysis for all West Eurasian samples
2812 in the imputed dataset, showing the posterior probability of the derived allele frequency
2813 trajectory for rs4792897, the most significant SNP in the selection peak spanning
2814 chr17:36615969-47588760.

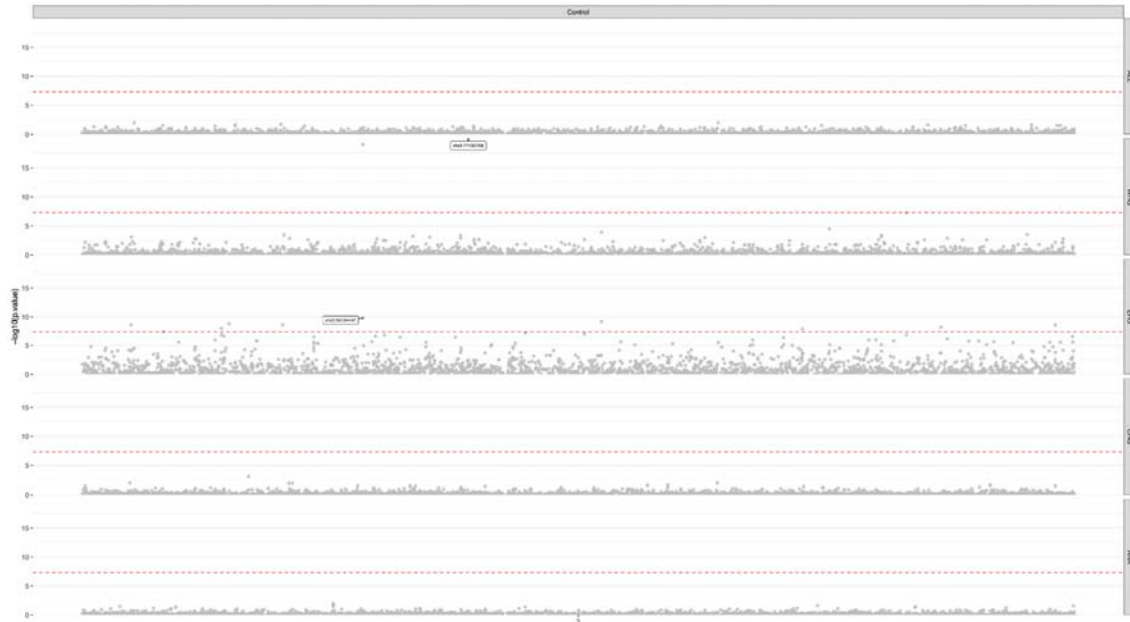
2815 The eleventh peak spanned the region chr17:36615969-47588760, with the most significant
2816 SNP being rs4792897 (*MAPT*; $p=4.65e-10$; $s=0.0294$), associated with snoring⁶⁰.

2817
2818
2819
2820
2821

2822 Selection in simulations with Ancestral Paintings

2823 *CLUES* analysis of all frequency paired simulated SNPs in both the true paths and inferred
2824 paths simulations detected no genome-wide significant sweep loci, indicating that the false
2825 positive rate of our ancestry stratified selection analysis is low.

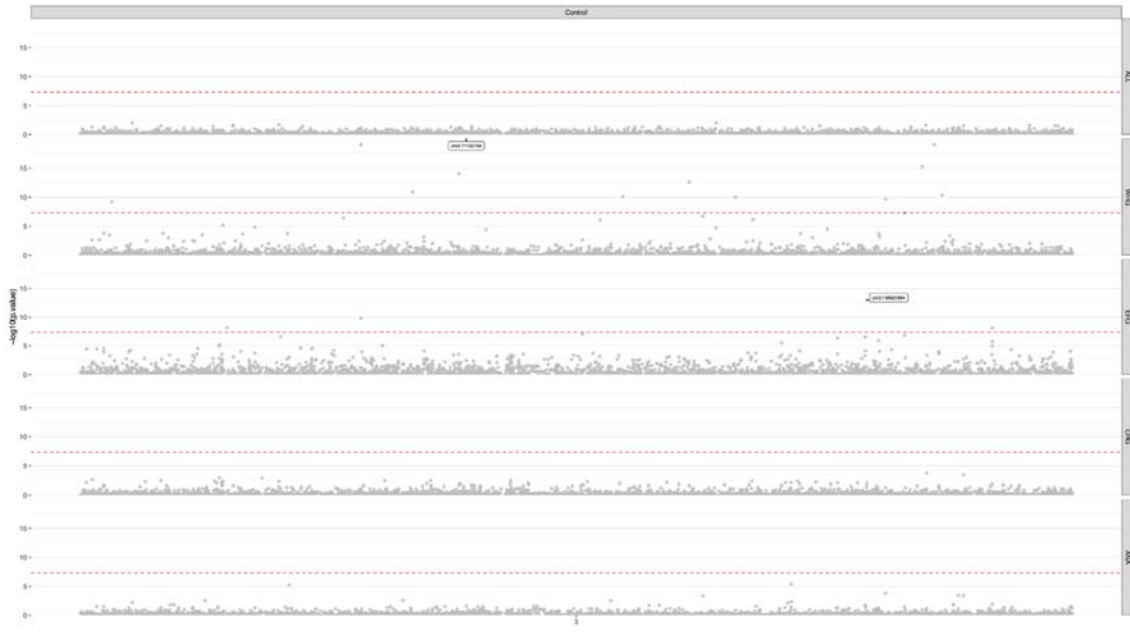
2826



2827

2828 **Figure S4a.15.** Manhattan plot of the p-values from running *CLUES* on a neutral simulation
2829 of chr3, using the true simulated paths of each ancestry painting. The first row shows results
2830 for all ancient samples considered in aggregate, and each subsequent row shows the results
2831 conditional on one of the four specific ancestral paintings: ANA (Anatolian Farmers), CHG
2832 (Caucasus Hunter-gatherers), WHG (Western Hunter-gatherers) and EHG (Eastern Hunter-
2833 gatherers).

2834

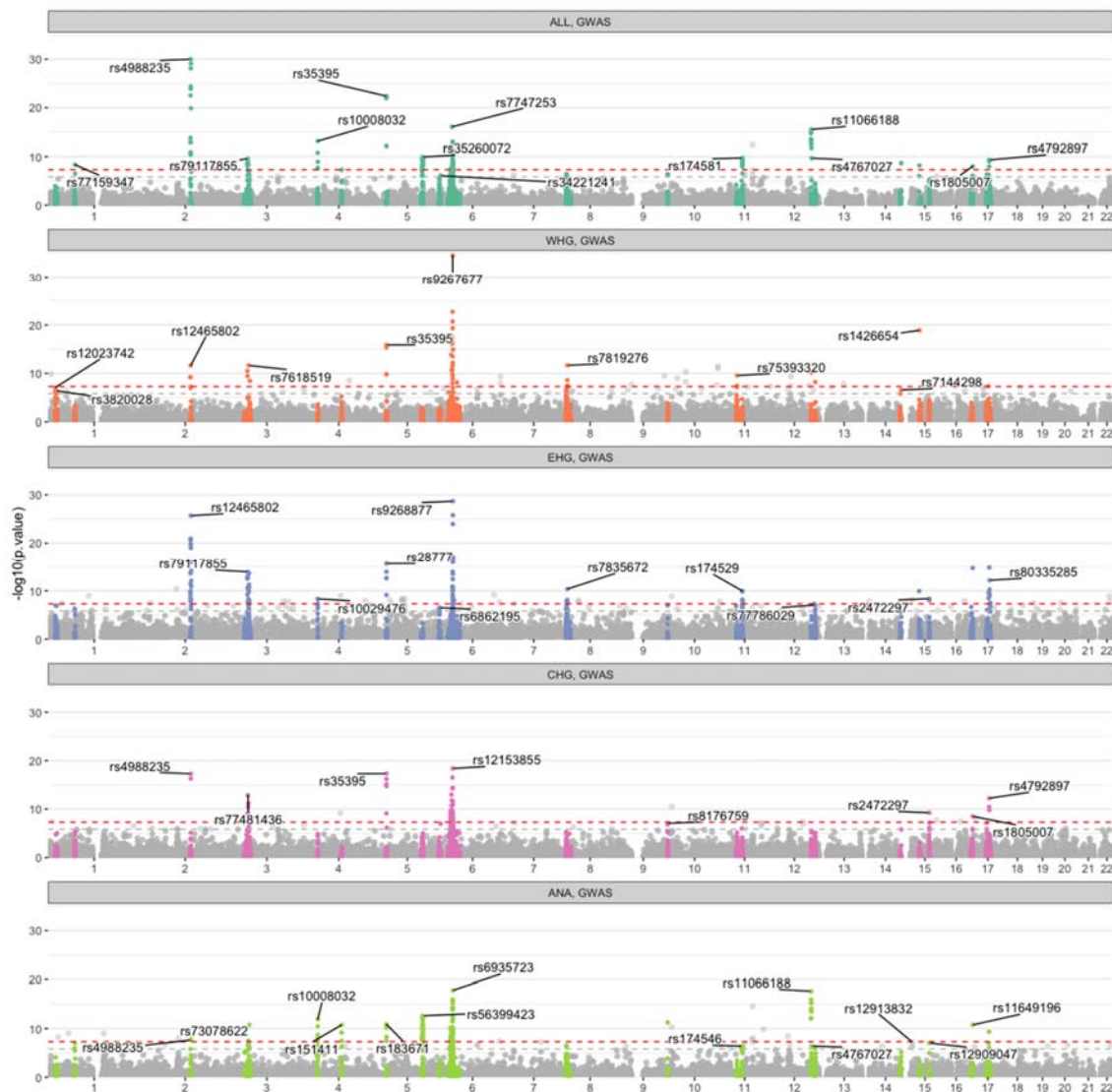


2835
 2836
 2837
 2838
 2839
 2840
 2841
 2842

Figure S4a.16. Manhattan plot of the p-values from running *CLUES* on a neutral simulation of chr3, using the inferred paths of each ancestry painting. The first row shows results for all ancient samples considered in aggregate, and each subsequent row shows the results conditional on one of the four specific ancestral paintings: ANA (Anatolian Farmers), CHG (Caucasus Hunter-gatherers), WHG (Western Hunter-gatherers) and EHG (Eastern Hunter-gatherers).

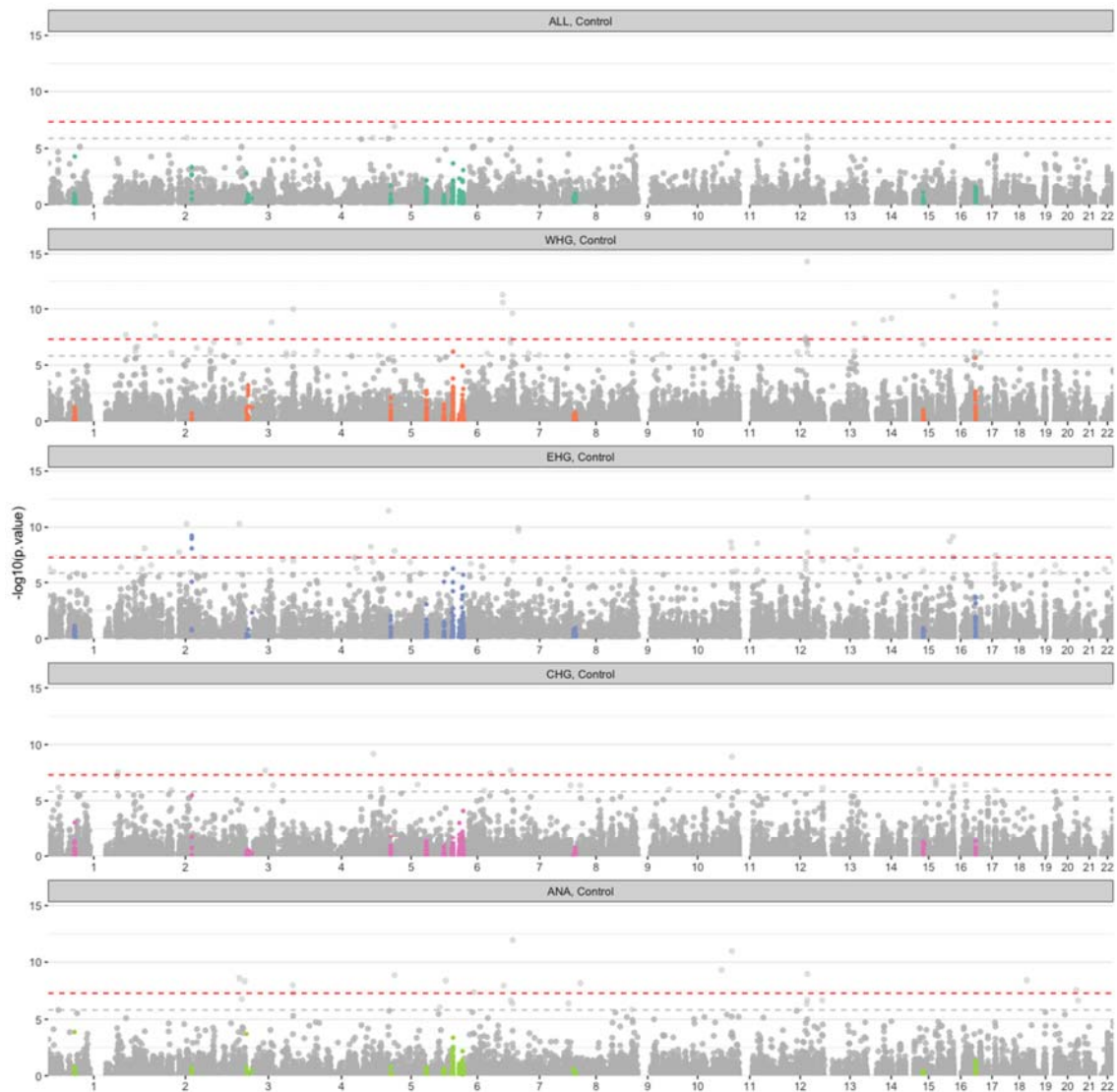
2843 Selection in aDNA with Ancestral Paintings

2844 CLUES analysis of all GWAS (n=33,323) and Control group SNPs (n=33,323) in the aDNA
2845 with Ancestral Paintings dataset identified 409 genome-wide significant SNPs ($p < 5e-8$); 346
2846 in the GWAS group and 63 in the Control group. Using a Bonferroni corrected significance
2847 threshold, we detected 758 significant SNPs ($p < 1.50e-06$); 593 in the GWAS group
2848 (78.23%) and 165 in the Control group. Within the GWAS group, we identified 21 non-
2849 overlapping Bonferroni corrected significant selection peaks across all ancestries (see Fig.
2850 S4a.17; Supplementary Table XVI).
2851



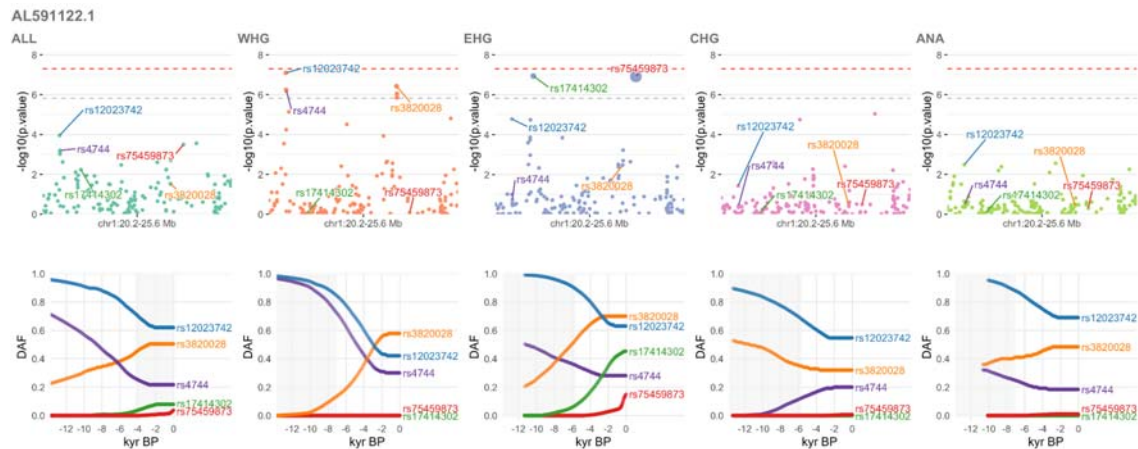
2852
2853 **Figure S4a.17.** Manhattan plot of the p-values from running CLUES on an aDNA time series
2854 conditioned on ancestry paintings from all West Eurasian samples in the imputed dataset for
2855 GWAS SNPs from the GWAS Catalog. The first row shows results for all ancient samples
2856 considered in aggregate, and each subsequent row shows the results conditional on one of

2857 the four specific ancestral paintings: ANA (Anatolian Farmers), CHG (Caucasus Hunter-
2858 gatherers), WHG (Western Hunter-gatherers) and EHG (Eastern Hunter-gatherers).
2859



2860
2861 **Figure S4a.18.** Manhattan plot of the p-values from running *CLUES* on an aDNA time series
2862 conditioned on ancestry paintings from all West Eurasian samples in the imputed dataset for
2863 Control SNPs, frequency paired with the GWAS SNPs. The first row shows results for all
2864 ancient samples considered in aggregate, and each subsequent row shows the results
2865 conditional on one of the four specific ancestral paintings: ANA (Anatolian Farmers), CHG
2866 (Caucasus Hunter-gatherers), WHG (Western Hunter-gatherers) and EHG (Eastern Hunter-
2867 gatherers).
2868

2869 Peak 1: AL591122.1



2870

2871

Figure S4a.19. Selection at the *AL591122.1* locus, spanning chr1:20236729-25570080.

2872

Results for the pan-ancestry analysis (ALL) plus the four marginal ancestries: Western hunter-gatherers (WHG), Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers (CHG) and Anatolian farmers (ANA). Row one shows zoomed Manhattan plots of the p-values for each ancestry (significant SNPs sized by their selection coefficients), and row two shows allele trajectories for the top SNPs across all ancestries.

2877

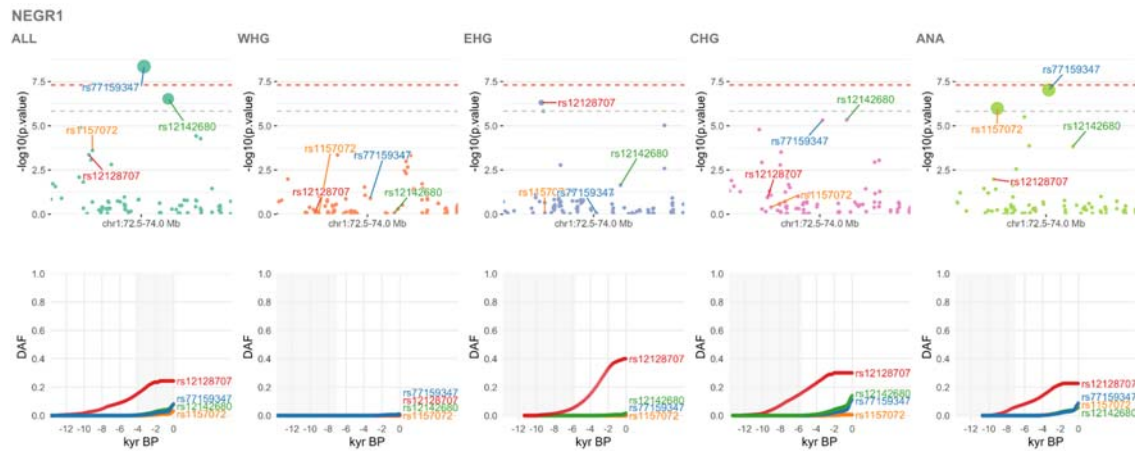
This peak spanned the region chr1:20236729-25570080, with significant SNPs including:

2878

- rs12023742 (*PLA2G2E - RN7SL304P*; WHG; $p=8.10e-08$; $s=-0.0136$), associated with group IIA secretory phospholipase A2 levels in individuals with elevated hsCRP⁶¹.
- rs17414302 (*PINK1-AS, PINK1*; EHG; $p=1.16e-07$; $s=0.0185$), associated with Household income (MTAG); Intelligence (MTAG)^{37,62}.
- rs75459873 (*MIR378F - H3P1*; EHG; $p=1.21e-07$; $s=0.0437$), associated with Psychotic experience (distressing)⁶³.
- rs3820028 (*E2F2*; WHG; $p=3.75e-07$; $s=0.0144$), associated with Heel bone mineral density⁵⁸.
- rs4744 (*PLA2G2A*; WHG; $p=5.67e-07$; $s=-0.0128$), associated with Blood protein levels⁴¹.

2889

2890 Peak 2: NEGR1



2891

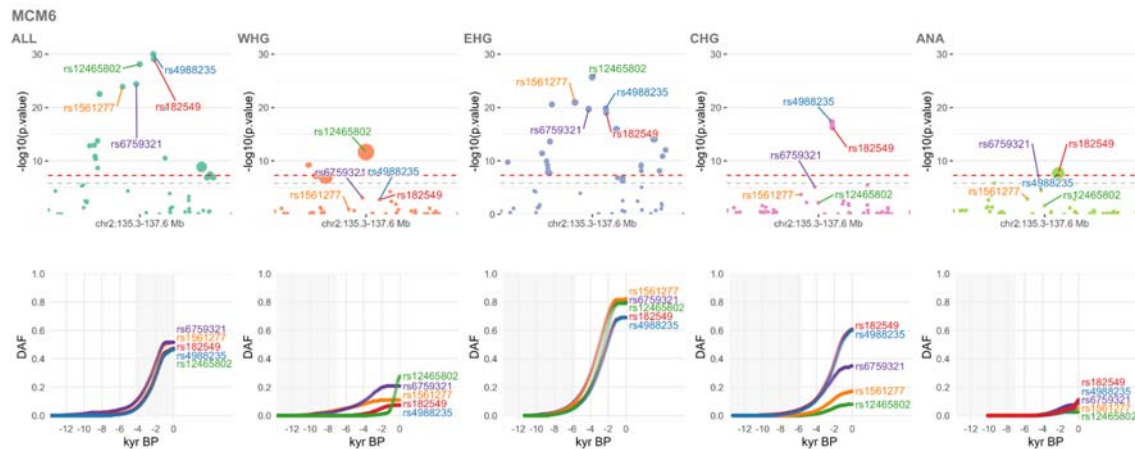
2892 **Figure S4a.20.** Selection at the *NEGR1* locus, spanning chr1:72480859-73978570. Results
 2893 for the pan-ancestry analysis (ALL) plus the four marginal ancestries: Western hunter-
 2894 gatherers (WHG), Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers (CHG) and
 2895 Anatolian farmers (ANA). Row one shows zoomed Manhattan plots of the p-values for each
 2896 ancestry (significant SNPs sized by their selection coefficients), and row two shows allele
 2897 trajectories for the top SNPs across all ancestries.

2898 This peak spanned the region chr1:72480859-73978570, with significant SNPs including:

- 2899 • rs77159347 (*LINC02797,AL583808.1*; ALL; $p=4.48e-09$; $s=0.052$), associated with
 2900 Household income (MTAG)³⁷.
- 2901 • rs12142680 (*KRT8P21 - RN7SKP19*; ALL; $p=2.99e-07$; $s=0.0437$), associated with
 2902 Educational attainment (years of education)⁶⁴.
- 2903 • rs12128707 (*NEGR1*; EHG; $p=4.93e-07$; $s=0.0182$), associated with Cognitive
 2904 performance; Cognitive performance (MTAG); Intelligence; Intelligence (MTAG)
 2905 62,65,66.
- 2906 • rs1157072 (*NEGR1*; ANA; $p=1.04e-06$; $s=0.0498$), associated with Household
 2907 income (MTAG); Intelligence (MTAG)^{37,62}.

2908

2909 Peak 3: MCM6



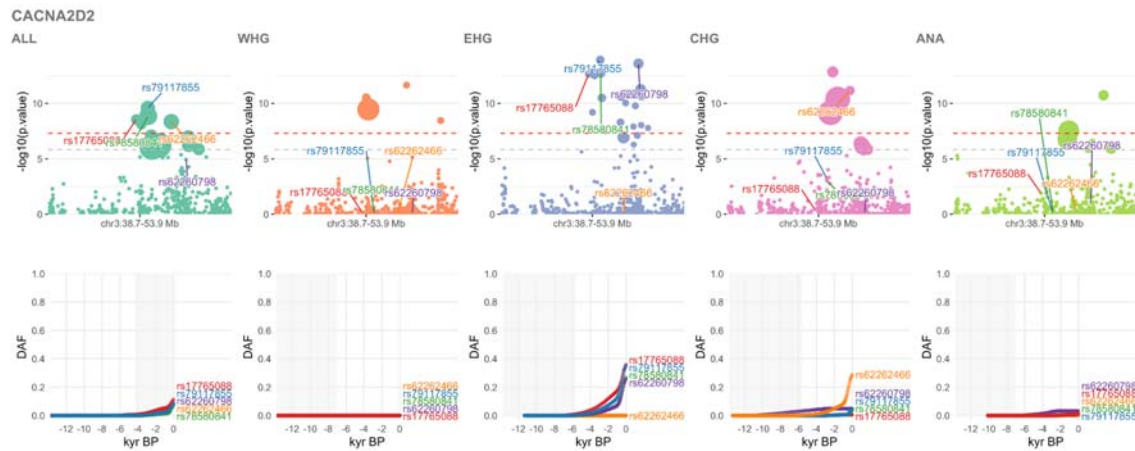
2910 **Figure S4a.21.** Selection at the *MCM6* locus, spanning chr2:135300859-137564020. Results
 2911 for the pan-ancestry analysis (ALL) plus the four marginal ancestries: Western hunter-
 2912 gatherers (WHG), Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers (CHG) and
 2913 Anatolian farmers (ANA). Row one shows zoomed Manhattan plots of the p-values for each
 2914 ancestry (significant SNPs sized by their selection coefficients), and row two shows allele
 2915 trajectories for the top SNPs across all ancestries.
 2916

2917 This peak spanned the region chr2:135300859-137564020, with significant SNPs including:

- 2918 • rs4988235 (*MCM6*; ALL; $p=9.86e-31$; $s=0.0199$), associated with Lactase
 2919 persistence; Blood protein levels; Body mass index; Hip circumference³⁸⁻⁴².
- 2920 • rs182549 (*MCM6*; ALL; $p=8.44e-30$; $s=0.0198$), associated with 1,5-anhydroglucitol
 2921 levels⁶⁷.
- 2922 • rs12465802 (*R3HDM1*; ALL; $p=7.71e-29$; $s=0.0196$), associated with Blood protein
 2923 levels; Mosquito bite size; Urinary metabolite levels in chronic kidney disease^{34,40,68}.
- 2924 • rs6759321 (*R3HDM1*; ALL; $p=4.07e-25$; $s=0.0188$), associated with Hand grip
 2925 strength⁶⁹.
- 2926 • rs1561277 (*ZRANB3*; ALL; $p=1.24e-24$; $s=0.0188$), associated with Hip
 2927 circumference³⁸.

2928

2929 Peak 4: CACNA2D2



2930

2931

Figure S4a.22. Selection at the *CACNA2D2* locus, spanning chr3:38723219-53850000.

2932

Results for the pan-ancestry analysis (ALL) plus the four marginal ancestries: Western

2933

hunter-gatherers (WHG), Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers

2934

(CHG) and Anatolian farmers (ANA). Row one shows zoomed Manhattan plots of the p-

2935

values for each ancestry (significant SNPs sized by their selection coefficients), and row two

2936

shows allele trajectories for the top SNPs across all ancestries.

2937

This peak spanned the region chr3:38723219-53850000, with significant SNPs including:

2938

- rs79117855 (*PTH1R* - *AC109583.3*; EHG; $p=1.18e-14$; $s=0.0305$), associated with Blood protein levels ⁴³.

2939

2940

- rs62260798 (*GNAI2*; EHG; $p=2.52e-14$; $s=0.0373$), associated with Morning person ⁷⁰.

2941

2942

- rs78580841 (*CCDC12*; EHG; $p=1.68e-13$; $s=0.037$), associated with Chronotype ⁷⁰.

2943

- rs17765088 (*CCR9*, *LZTFL1*; EHG; $p=2.22e-13$; $s=0.0257$), associated with

2944

Macrophage inflammatory protein 1b levels ⁷¹.

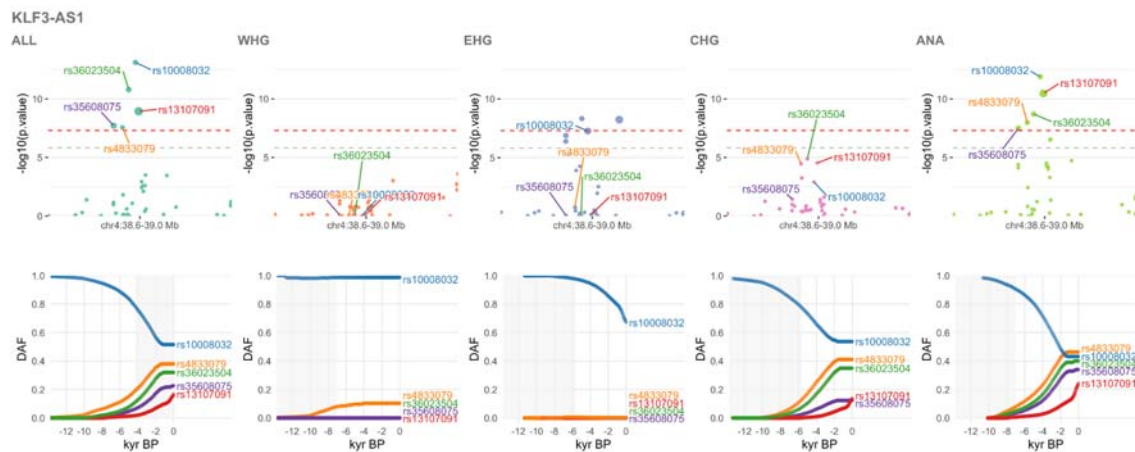
2945

- rs62262466 (*ARIH2*; CHG; $p=6.98e-12$; $s=0.0351$), associated with Morning person ⁷⁰.

2946

2947

2948 Peak 5: *KLF3-AS1*



2949

2950

Figure S4a.23. Selection at the *KLF3-AS1* locus, spanning chr4:38593259-38966080.

2951

Results for the pan-ancestry analysis (ALL) plus the four marginal ancestries: Western

2952

hunter-gatherers (WHG), Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers

2953

(CHG) and Anatolian farmers (ANA). Row one shows zoomed Manhattan plots of the p-

2954

values for each ancestry (significant SNPs sized by their selection coefficients), and row two

2955

shows allele trajectories for the top SNPs across all ancestries.

2956

This peak spanned the region chr4:38593259-38966080, with significant SNPs including:

2957

- rs10008032 (*KLF3 - RNA5SP158*; ALL; $p=7.66e-14$; $s=-0.0165$), associated with Allergic disease (asthma, hay fever or eczema) ⁴⁴.

2958

- rs36023504 (*KLF3*; ALL; $p=1.62e-11$; $s=0.0179$), associated with Body mass index; Red cell distribution width ⁵⁸.

2959

- rs13107091 (*RNA5SP158 - TLR10*; ANA; $p=3.41e-11$; $s=0.0281$), associated with Atopic asthma ⁷².

2960

- rs4833079 (*KLF3-AS1*; ANA; $p=1.02e-08$; $s=0.0158$), associated with Body mass index ⁴².

2961

- rs35608075 (*LINC02278 - KLF3-AS1*; ALL; $p=1.95e-08$; $s=0.0208$), associated with Male-pattern baldness ⁷³.

2962

2963

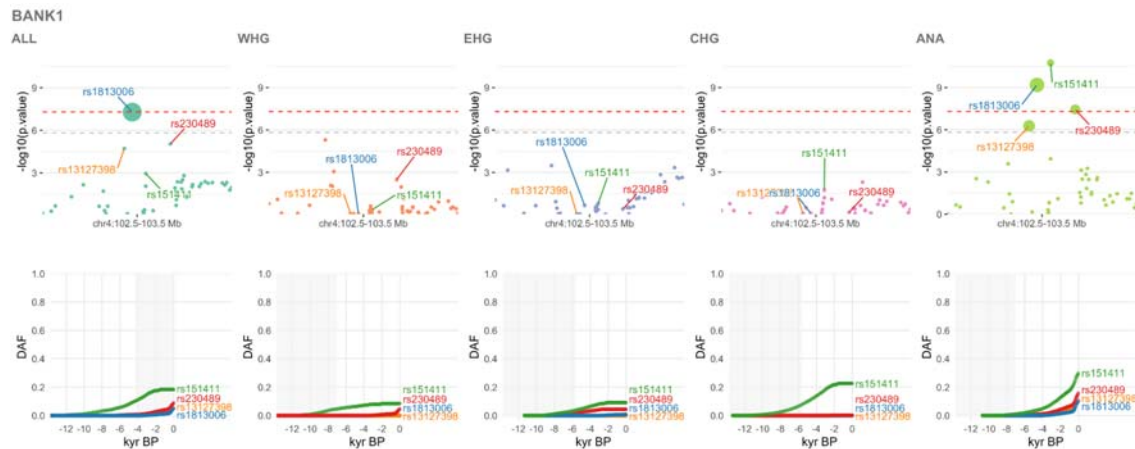
2964

2965

2966

2967

2968 Peak 6: BANK1



2969

2970 **Figure S4a.24.** Selection at the *BANK1* locus, spanning chr4:102507109-103525350.

2971

2972 Results for the pan-ancestry analysis (ALL) plus the four marginal ancestries: Western

2973

2974 hunter-gatherers (WHG), Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers

2975

2976 (CHG) and Anatolian farmers (ANA). Row one shows zoomed Manhattan plots of the p-
2977 values for each ancestry (significant SNPs sized by their selection coefficients), and row two
2978 shows allele trajectories for the top SNPs across all ancestries.

2976

2977 This peak spanned the region chr4:102507109-103525350, with significant SNPs including:

2977

- rs151411 (*BANK1* - *SLC39A8*; ANA; $p=1.80e-11$; $s=0.0247$), associated with General cognitive ability⁷⁴.

2978

2979

- rs1813006 (*BANK1* - *SLC39A8*; ANA; $p=6.60e-10$; $s=0.0563$), associated with Intelligence (MTAG)⁶².

2980

2981

- rs230489 (*AF213884.2* - *AF213884.3*; ANA; $p=3.63e-08$; $s=0.036$), associated with Brain region volumes; General cognitive ability; Intelligence (MTAG)^{62,74,75}.

2982

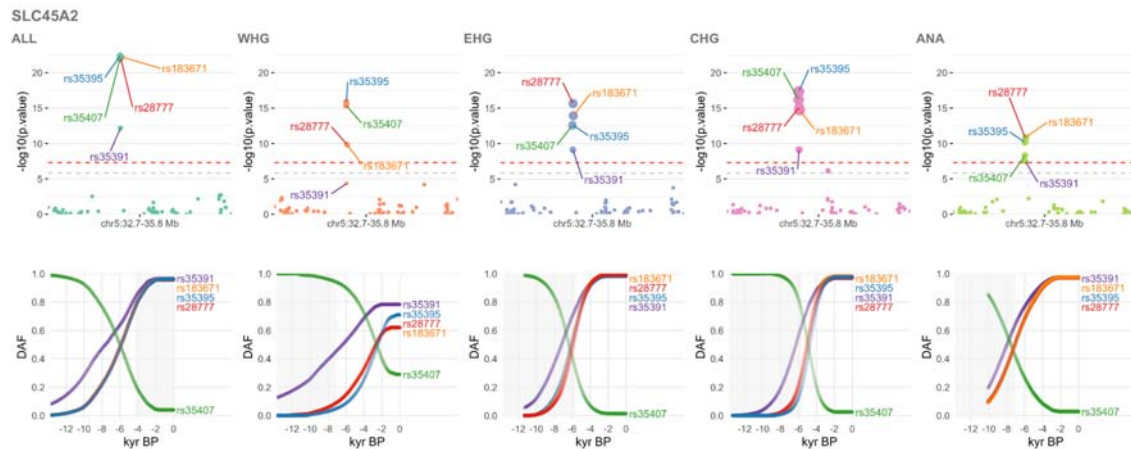
2983

- rs13127398 (*BANK1*, *AP002075.1*; ANA; $p=5.36e-07$; $s=0.0439$), associated with Cardiovascular disease⁵⁸.

2984

2985

2986 Peak 7: SLC45A2



2987

2988 **Figure S4a.25.** Selection at the *SLC45A2* locus, spanning chr5:32710489-35848560.

2989

2990 Results for the pan-ancestry analysis (ALL) plus the four marginal ancestries: Western
 2991 hunter-gatherers (WHG), Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers
 2992 (CHG) and Anatolian farmers (ANA). Row one shows zoomed Manhattan plots of the p-
 2993 values for each ancestry (significant SNPs sized by their selection coefficients), and row two

2994

This peak spanned the region chr5:32710489-35848560, with significant SNPs including:

2995

- rs35395 (*SLC45A2*; ALL; $p=4.13e-23$; $s=0.022$), associated with Skin pigmentation ⁷⁶.

2996

- rs183671 (*SLC45A2*; ALL; $p=5.51e-23$; $s=0.0221$), associated with Hair color; Skin colour saturation; Skin pigmentation traits ^{56,77,78}.

2997

2998

- rs28777 (*SLC45A2*; ALL; $p=8.48e-23$; $s=0.0217$), associated with Black vs. blond hair color; Black vs. red hair color; Skin, hair and eye pigmentation (multivariate analysis) ^{79,80}.

2999

3000

3001

- rs35407 (*SLC45A2*; ALL; $p=1.06e-22$; $s=-0.0221$), associated with Basal cell carcinoma; Cutaneous squamous cell carcinoma; Keratinocyte cancer (MTAG); Melanoma; Squamous cell carcinoma ^{52-54,59}.

3002

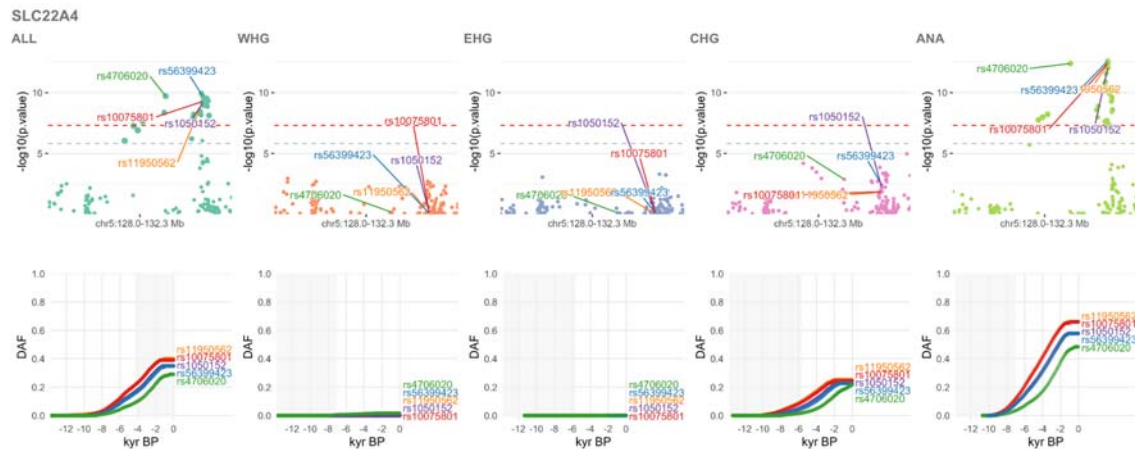
3003

3004

- rs35391 (*SLC45A2*; ALL; $p=6.87e-13$; $s=0.0161$), associated with Tanning ⁸¹.

3005

3006 Peak 8: SLC22A4



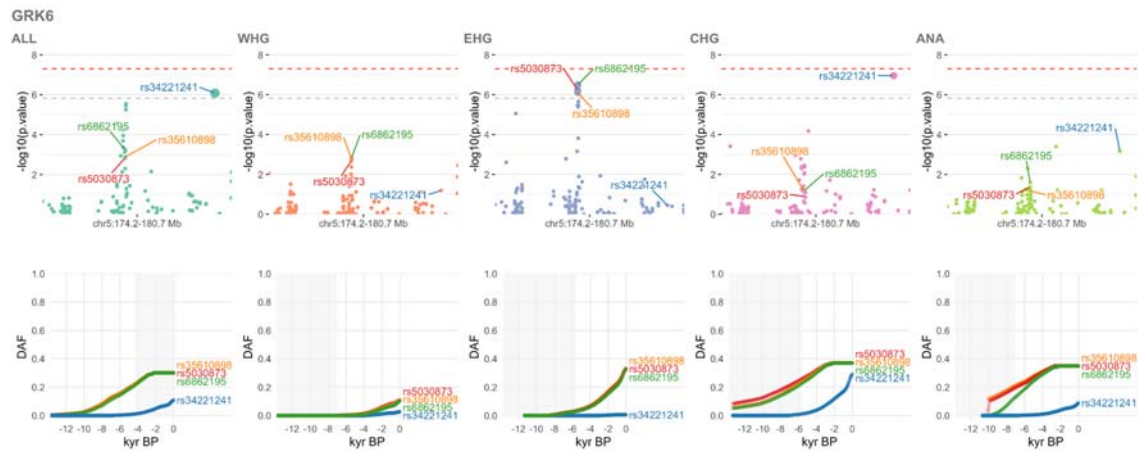
3007 **Figure S4a.26.** Selection at the *SLC22A4* locus, spanning chr5:128016159-132349650.
 3008 Results for the pan-ancestry analysis (ALL) plus the four marginal ancestries: Western
 3009 hunter-gatherers (WHG), Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers
 3010 (CHG) and Anatolian farmers (ANA). Row one shows zoomed Manhattan plots of the p-
 3011 values for each ancestry (significant SNPs sized by their selection coefficients), and row two
 3012 shows allele trajectories for the top SNPs across all ancestries.

3014 This peak spanned the region chr5:128016159-132349650, with significant SNPs including:

- 3015 • rs56399423 (*MIR3936HG,SLC22A4*; ANA; $p=2.48e-13$; $s=0.017$), associated with
 3016 Inflammatory bowel disease ⁸².
- 3017 • rs10075801 (*SLC22A4,MIR3936HG*; ANA; $p=3.75e-13$; $s=0.0164$), associated with
 3018 Granulocyte count; Myeloid white cell count; Neutrophil count; Sum basophil
 3019 neutrophil counts; Sum neutrophil eosinophil counts; White blood cell count ⁸³.
- 3020 • rs4706020 (*CDC42SE2*; ANA; $p=4.03e-13$; $s=0.0184$), associated with Itch intensity
 3021 from mosquito bite; Itch intensity from mosquito bite adjusted by bite size ³⁴.
- 3022 • rs11950562 (*MIR3936HG,SLC22A4*; ANA; $p=5.38e-13$; $s=0.0163$), associated with
 3023 Blood metabolite levels; Mean platelet volume ^{83,84}.
- 3024 • rs1050152 (*SLC22A4,MIR3936HG*; ANA; $p=8.09e-13$; $s=0.0167$), associated with
 3025 Nasal polyps ⁸⁵.

3026

3027 Peak 9: GRK6



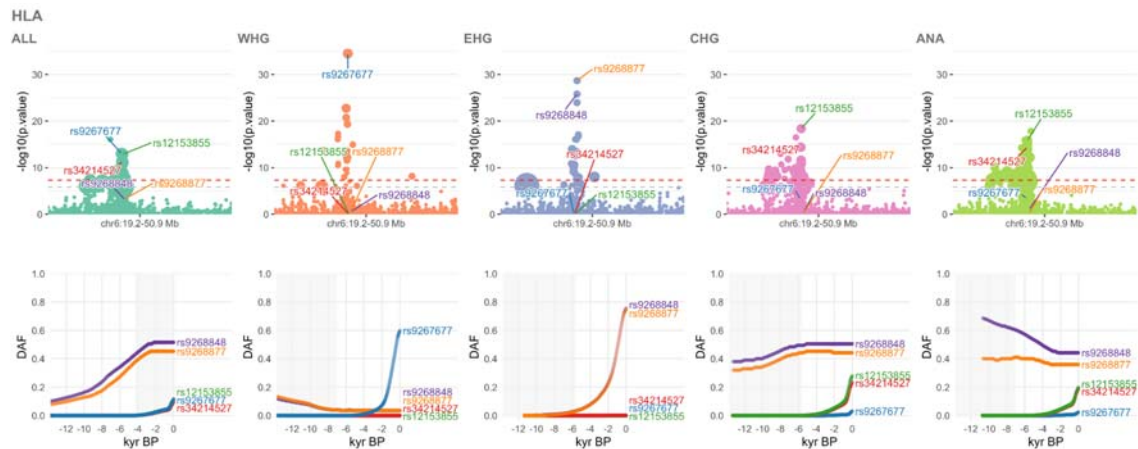
3028 **Figure S4a.27.** Selection at the *GRK6* locus, spanning chr5:174156169-180661980. Results
 3029 for the pan-ancestry analysis (ALL) plus the four marginal ancestries: Western hunter-
 3030 gatherers (WHG), Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers (CHG) and
 3031 Anatolian farmers (ANA). Row one shows zoomed Manhattan plots of the p-values for each
 3032 ancestry (significant SNPs sized by their selection coefficients), and row two shows allele
 3033 trajectories for the top SNPs across all ancestries.
 3034

3035 This peak spanned the region chr5:174156169-180661980, with significant SNPs including:

- 3036 • rs34221241 (*FLT4*; CHG; $p=1.11e-07$; $s=0.0253$), associated with Blood protein
 3037 levels ⁴⁰.
- 3038 • rs6862195 (*SLC34A1*; EHG; $p=3.18e-07$; $s=0.0232$), associated with Estimated
 3039 glomerular filtration rate ⁸⁶.
- 3040 • rs5030873 (*SLC34A1*; EHG; $p=5.85e-07$; $s=0.0228$), associated with Creatinine
 3041 levels ⁸⁷.
- 3042 • rs35610898 (*SLC34A1*; EHG; $p=7.84e-07$; $s=0.0224$), associated with Estimated
 3043 glomerular filtration rate ⁸⁶.

3044

3045 Peak 10: HLA



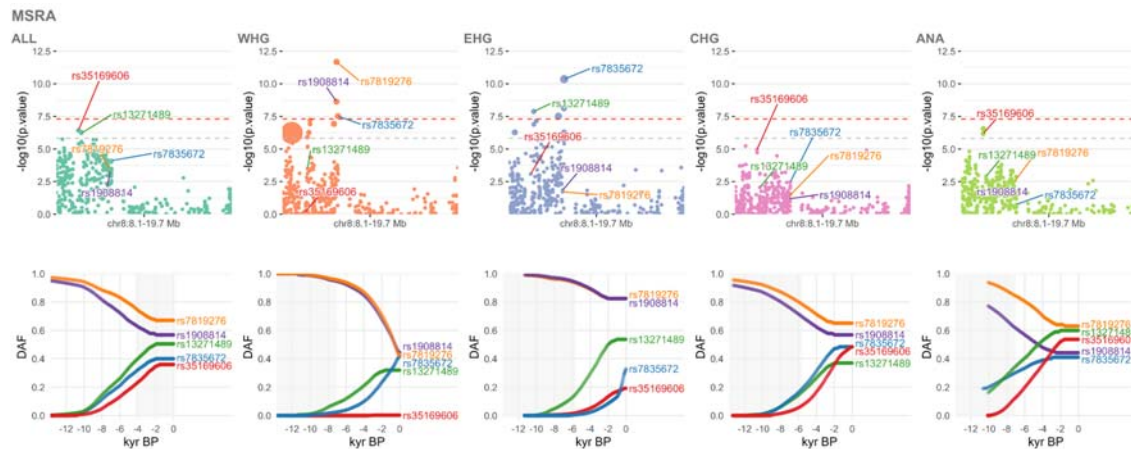
3046 **Figure S4a.28.** Selection at the *HLA* locus, spanning chr6:19191049-50921600. Results for
 3047 the pan-ancestry analysis (ALL) plus the four marginal ancestries: Western hunter-gatherers
 3048 (WHG), Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers (CHG) and Anatolian
 3049 farmers (ANA). Row one shows zoomed Manhattan plots of the p-values for each ancestry
 3050 (significant SNPs sized by their selection coefficients), and row two shows allele trajectories
 3051 for the top SNPs across all ancestries.
 3052

3053 This peak spanned the region chr6:19191049-50921600, with significant SNPs including:

- 3054 • rs9267677 (*C2*; WHG; $p=3.03e-35$; $s=0.0365$), associated with Cognitive
 3055 performance (MTAG); Educational attainment (MTAG); Educational attainment (years
 3056 of education); Highest math class taken (MTAG) ⁶⁶.
- 3057 • rs9268877 (*HLA-DRB9*; EHG; $p=2.07e-29$; $s=0.0249$), associated with Poor
 3058 prognosis in Crohn's disease; Ulcerative colitis ⁸⁸⁻⁹⁰.
- 3059 • rs9268848 (*HLA-DRB9*; EHG; $p=1.56e-26$; $s=0.0247$), associated with Nonatopic
 3060 asthma; Urate levels ^{72,91}.
- 3061 • rs12153855 (*TNXB,AL662884.2*; CHG; $p=4.19e-19$; $s=0.0334$), associated with Age-
 3062 related macular degeneration; Atopic dermatitis ^{92,93}.
- 3063 • rs34214527 (*TNXB*; CHG; $p=4.09e-15$; $s=0.0348$), associated with Highest math
 3064 class taken ⁶⁶.

3065

3066 Peak 11: MSRA



3067

3068

Figure S4a.29. Selection at the *MSRA* locus, spanning chr8:8142579-19746880. Results for the pan-ancestry analysis (ALL) plus the four marginal ancestries: Western hunter-gatherers (WHG), Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers (CHG) and Anatolian farmers (ANA). Row one shows zoomed Manhattan plots of the p-values for each ancestry (significant SNPs sized by their selection coefficients), and row two shows allele trajectories for the top SNPs across all ancestries.

3072

3073

This peak spanned the region chr8:8142579-19746880, with significant SNPs including:

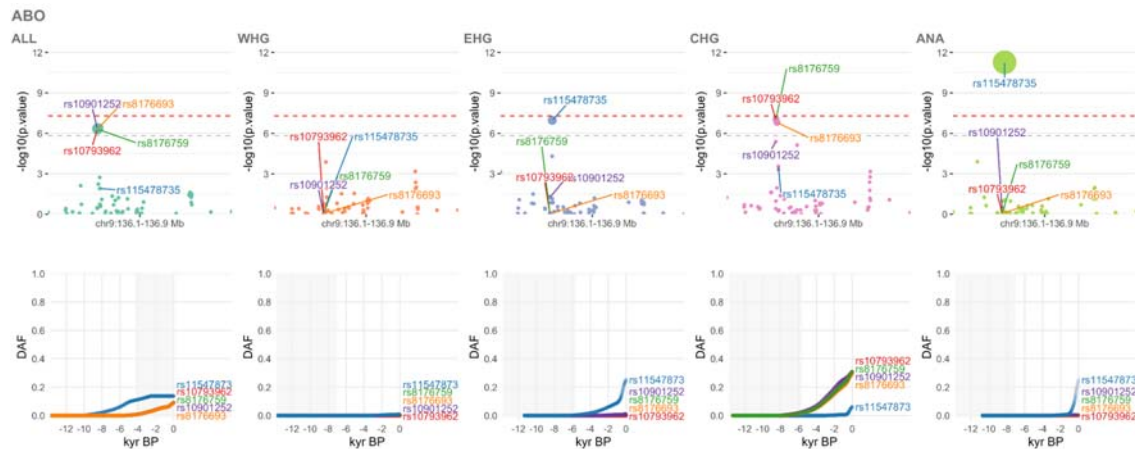
3074

- rs7819276 (*AC025857.1 - OR7E158P*; WHG; $p=2.10e-12$; $s=-0.0207$), associated with General factor of neuroticism ⁹⁴.
- rs7835672 (*AC107918.3*; EHG; $p=4.43e-11$; $s=0.0287$), associated with General factor of neuroticism ⁹⁴.
- rs1908814 (*AC025857.1 - OR7E158P*; WHG; $p=2.33e-09$; $s=-0.0201$), associated with General factor of neuroticism; Neuroticism ^{94,95}.
- rs13271489 (*LINC00599 - AC034111.2*; EHG; $p=1.33e-08$; $s=0.016$), associated with Diastolic blood pressure x smoking status (current vs non-current) interaction (2df test); Diastolic blood pressure x smoking status (ever vs never) interaction (2df test); Systolic blood pressure x smoking status (current vs non-current) interaction (2df test); Systolic blood pressure x smoking status (ever vs never) interaction (2df test) ⁹⁶.
- rs35169606 (*TNKS*; ALL; $p=3.85e-07$; $s=0.0137$), associated with Lifetime smoking index ⁹⁷.

3086

3087

3088 Peak 12: ABO



3089

3090

Figure S4a.30. Selection at the *ABO* locus, spanning chr9:136127999-136925660. Results for the

3091

pan-ancestry analysis (ALL) plus the four marginal ancestries: Western hunter-gatherers (WHG),

3092

Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers (CHG) and Anatolian farmers (ANA).

3093

Row one shows zoomed Manhattan plots of the p-values for each ancestry (significant SNPs sized by

3094

their selection coefficients), and row two shows allele trajectories for the top SNPs across all

3095

ancestries.

3096

This peak spanned the region chr9:136127999-136925660, with significant SNPs including:

3097

- rs115478735 (*ABO*; ANA; $p=5.36e-12$; $s=0.095$), associated with Blood protein levels ⁴⁰.

3098

- rs8176759 (*ABO*; CHG; $p=8.90e-08$; $s=0.0227$), associated with Granulocyte percentage of myeloid white cells; Plateletcrit ⁸³.

3099

- rs10793962 (*ABO*; CHG; $p=9.78e-08$; $s=0.0226$), associated with Blood protein levels; Intraocular pressure ^{40,98,99}.

3100

- rs8176693 (*ABO*; CHG; $p=1.55e-07$; $s=0.0233$), associated with Blood protein levels; Blood protein levels in cardiovascular risk; Endothelial growth factor levels; High serum lipase activity; Serum lipase activity ^{40,100-102}.

3101

- rs10901252 (*ABO*; ALL; $p=4.39e-07$; $s=0.0371$), associated with Blood protein levels; Hematocrit; Hemoglobin concentration; vWF levels ^{43,83,103}.

3102

3103

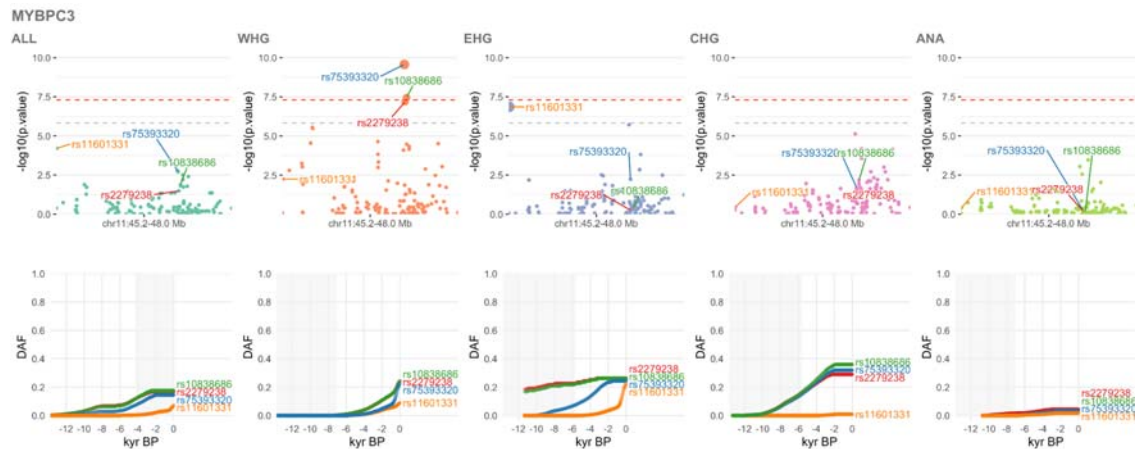
3104

3105

3106

3107

3108 Peak 13: MYBPC3



3109

Figure S4a.31. Selection at the *MYBPC3* locus, spanning chr11:45227569-48018360. Results for the pan-ancestry analysis (ALL) plus the four marginal ancestries: Western hunter-gatherers (WHG), Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers (CHG) and Anatolian farmers (ANA). Row one shows zoomed Manhattan plots of the p-values for each ancestry (significant SNPs sized by their selection coefficients), and row two shows allele trajectories for the top SNPs across all ancestries.

3116

This peak spanned the region chr11:45227569-48018360, with significant SNPs including:

3117

- rs75393320 (*ACP2*; WHG; $p=2.73e-10$; $s=0.0351$), associated with HDL cholesterol ¹⁰⁴.

3118

- rs10838686 (*MADD*; WHG; $p=3.56e-08$; $s=0.0261$), associated with High density lipoprotein cholesterol levels ¹⁰⁵.

3119

- rs2279238 (*NR1H3*; WHG; $p=6.49e-08$; $s=0.0261$), associated with Creatinine levels ¹⁰⁶.

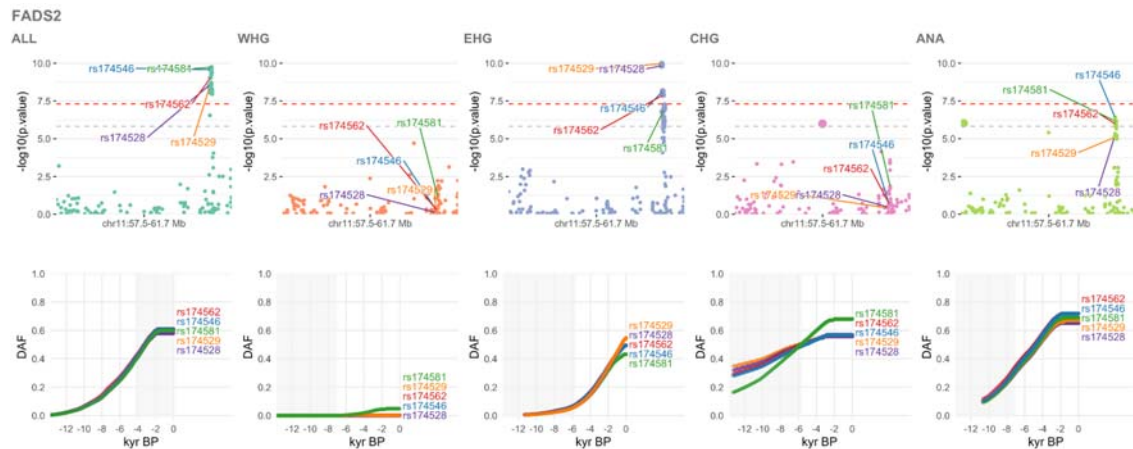
3121

- rs11601331 (*TSPAN18*; EHG; $p=1.38e-07$; $s=0.04$), associated with Cortical brain region measurements (area, volume and thickness) ¹⁰⁷.

3122

3123

3124 Peak 14: FADS2

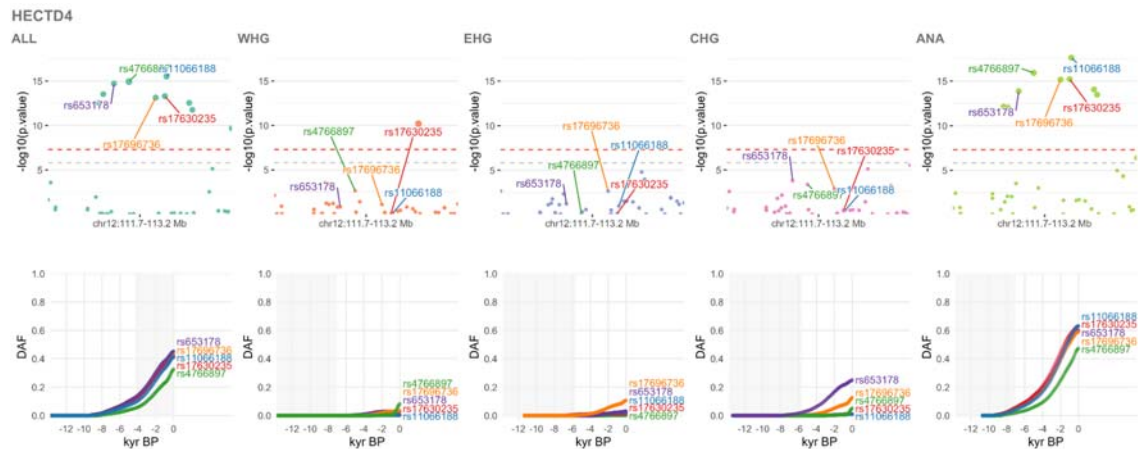


3125 **Figure S4a.32.** Selection at the *FADS2* locus, spanning chr11:57467039-61706010. Results for the
 3126 pan-ancestry analysis (ALL) plus the four marginal ancestries: Western hunter-gatherers (WHG),
 3127 Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers (CHG) and Anatolian farmers (ANA).
 3128 Row one shows zoomed Manhattan plots of the p-values for each ancestry (significant SNPs sized by
 3129 their selection coefficients), and row two shows allele trajectories for the top SNPs across all
 3130 ancestries.
 3131

3132 This peak spanned the region chr11:57467039-61706010, with significant SNPs including:

- 3133 • rs174529 (*TMEM258,MYRF*; EHG; $p=1.12e-10$; $s=0.0198$), associated with HDL
 3134 cholesterol; Heel bone mineral density; High density lipoprotein cholesterol levels; Low
 3135 density lipoprotein cholesterol levels; Total cholesterol levels; Triglyceride levels;
 3136 Triglycerides ^{58,105,108–110}.
- 3137 • rs174528 (*TMEM258,MYRF*; EHG; $p=1.48e-10$; $s=0.0198$), associated with Gondoic acid
 3138 (20:1n-9) levels; Phosphatidylcholine-ether levels; Plasma omega-6 polyunsaturated fatty
 3139 acid levels (arachidonic acid); Serum metabolite ratios in chronic kidney disease; Trans
 3140 fatty acid levels; Vaccenic acid (18:1n-7) levels ^{47,111–114}.
- 3141 • rs174581 (*FADS2*; ALL; $p=1.87e-10$; $s=0.0133$), associated with Male-pattern baldness;
 3142 Serum metabolite ratios in chronic kidney disease ^{46,47}.
- 3143 • rs174546 (*FADS2,FADS1*; ALL; $p=2.65e-10$; $s=0.0131$), associated with C-reactive
 3144 protein levels or HDL-cholesterol levels (pleiotropy); C-reactive protein levels or
 3145 triglyceride levels (pleiotropy); Change in serum metabolite levels; Change in serum
 3146 metabolite levels (CMS); Cholesterol, total; Glycerophospholipid levels; HDL cholesterol;
 3147 HDL cholesterol levels; High density lipoprotein cholesterol levels; LDL cholesterol; LDL
 3148 cholesterol levels; Low density lipoprotein cholesterol levels; Plasma omega-6
 3149 polyunsaturated fatty acid levels (gamma-linolenic acid); QT interval; Serum metabolite
 3150 levels; Serum metabolite levels (CMS); Total cholesterol levels; Trans fatty acid levels;
 3151 Triglyceride levels; Triglycerides ^{105,110–112,115–121}.
- 3152 • rs174562 (*FADS2,FADS1*; ALL; $p=8.63e-10$; $s=0.0127$), associated with Asthma; Serum
 3153 metabolite ratios in chronic kidney disease ^{47,122}.

3154 Peak 15: HECTD4



3155
3156
3157
3158
3159
3160
3161

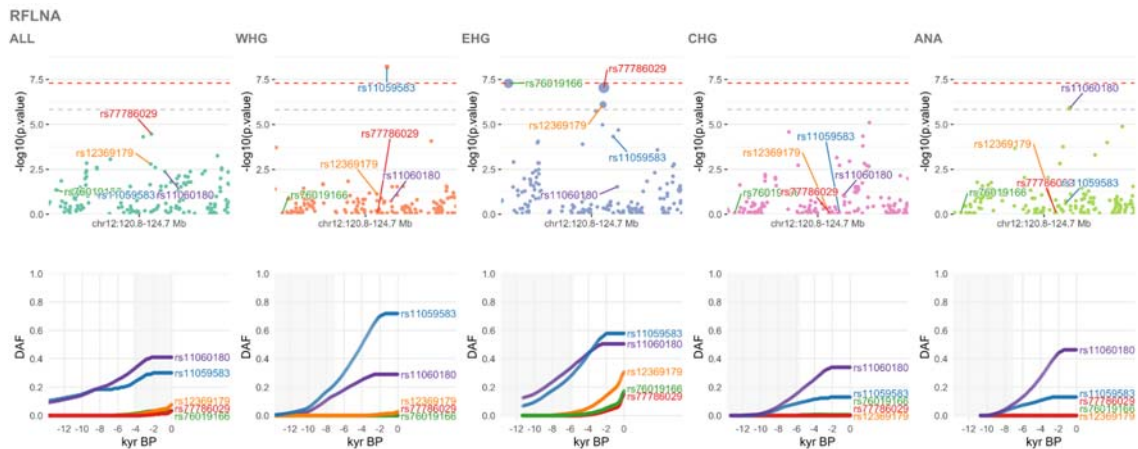
Figure S4a.33. Selection at the *HECTD4* locus, spanning chr12:111706879-113205500. Results for the pan-ancestry analysis (ALL) plus the four marginal ancestries: Western hunter-gatherers (WHG), Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers (CHG) and Anatolian farmers (ANA). Row one shows zoomed Manhattan plots of the p-values for each ancestry (significant SNPs sized by their selection coefficients), and row two shows allele trajectories for the top SNPs across all ancestries.

3162 This peak spanned the region chr12:111706879-113205500, with significant SNPs including:

- 3163 • rs11066188 (*HECTD4*; ANA; $p=2.35e-18$; $s=0.0202$), associated with Celiac disease and
- 3164 Rheumatoid arthritis ⁴⁸.
- 3165 • rs4766897 (*ACAD10*; ANA; $p=1.17e-16$; $s=0.0208$), associated with Fibrinogen levels ¹²³.
- 3166 • rs17630235 (*TRAFD1 - HECTD4*; ANA; $p=5.98e-16$; $s=0.0195$), associated with Body
- 3167 mass index; Diastolic blood pressure; Parental longevity (combined parental attained
- 3168 age, Martingale residuals); Systolic blood pressure; Tonsillectomy ^{39,42,124–126}.
- 3169 • rs17696736 (*NAA25*; ANA; $p=6.77e-16$; $s=0.0196$), associated with Coronary artery
- 3170 disease; Diastolic blood pressure x alcohol consumption interaction (2df test); Estimated
- 3171 glomerular filtration rate; Mean arterial pressure; Mean arterial pressure x alcohol
- 3172 consumption interaction (2df test); Parental longevity (combined parental attained age,
- 3173 Martingale residuals); Systolic blood pressure x alcohol consumption interaction (2df
- 3174 test); Type 1 diabetes; Urate levels ^{91,126–133}.
- 3175 • rs653178 (*ATXN2*; ALL; $p=1.92e-15$; $s=0.0194$), associated with Allergic disease
- 3176 (asthma, hay fever or eczema); Asthma; Blood pressure; Celiac disease; Celiac disease
- 3177 or Rheumatoid arthritis; Chronic kidney disease; Diastolic blood pressure; Eczema;
- 3178 Eosinophil counts; Eosinophil percentage of granulocytes; Eosinophil percentage of white
- 3179 cells; Hay fever and/or eczema; Inflammatory bowel disease; LDL cholesterol; Mean
- 3180 arterial pressure; Monocyte count; Myocardial infarction; Neutrophil percentage of
- 3181 granulocytes; Sarcoidosis; Sum eosinophil basophil counts; Systemic lupus
- 3182 erythematosus; Thyroid peroxidase antibody positivity; Tonsillectomy; Total cholesterol
- 3183 levels; Type 1 diabetes; Urate levels ^{58,82,83,104,125,130,134–148}.

3184

3185 Peak 16: RFLNA



3186
3187
3188
3189
3190
3191
3192

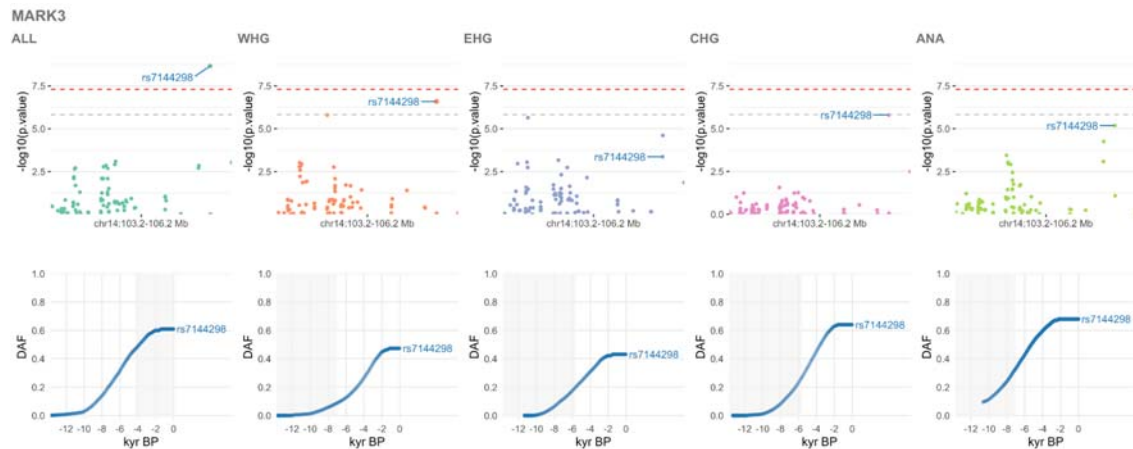
Figure S4a.34. Selection at the *RFLNA* locus, spanning chr12:120813919-124667680. Results for the pan-ancestry analysis (ALL) plus the four marginal ancestries: Western hunter-gatherers (WHG), Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers (CHG) and Anatolian farmers (ANA). Row one shows zoomed Manhattan plots of the p-values for each ancestry (significant SNPs sized by their selection coefficients), and row two shows allele trajectories for the top SNPs across all ancestries.

3193 This peak spanned the region chr12:120813919-124667680, with significant SNPs including:

- 3194 ● rs11059583 (*HCAR3 - HCAR1*; WHG; $p=6.21e-09$; $s=0.0144$), associated with Waist-to-
- 3195 hip ratio adjusted for BMI ⁷².
- 3196 ● rs76019166 (*PXN*; EHG; $p=5.20e-08$; $s=0.0339$), associated with Appendicular lean
- 3197 mass ¹⁴⁹.
- 3198 ● rs77786029 (*ZCCHC8*; EHG; $p=9.02e-08$; $s=0.0399$), associated with Red cell
- 3199 distribution width ⁵⁸.
- 3200 ● rs12369179 (*ZCCHC8*; EHG; $p=7.96e-07$; $s=0.0239$), associated with Body mass index;
- 3201 Waist-to-hip ratio adjusted for BMI ^{72,150}.
- 3202 ● rs11060180 (*CCDC62*; ANA; $p=1.27e-06$; $s=0.0146$), associated with Parkinson's
- 3203 disease ¹⁵¹⁻¹⁵³.

3204

3205 Peak 17: MARK3



3206

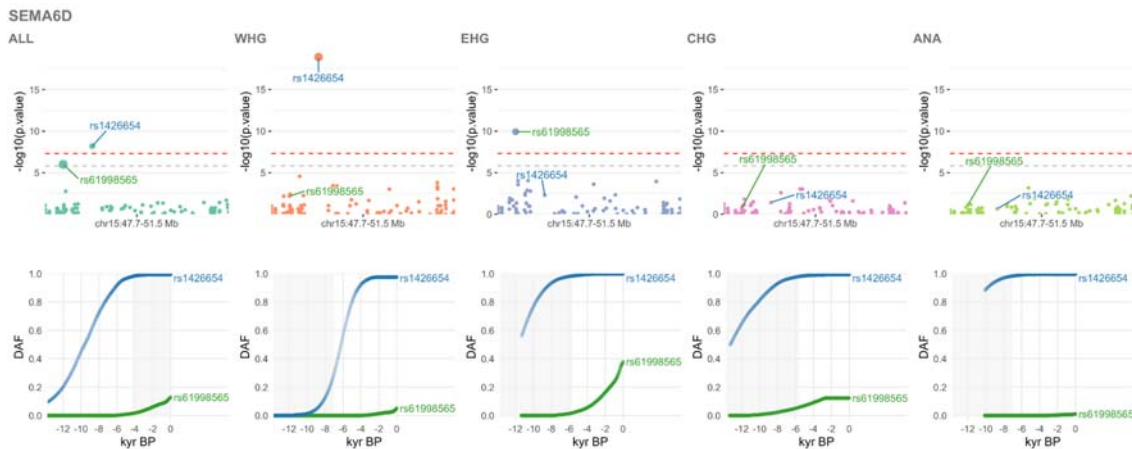
3207 **Figure S4a.35.** Selection at the *MARK3* locus, spanning chr14:103239629-106248400. Results for
 3208 the pan-ancestry analysis (ALL) plus the four marginal ancestries: Western hunter-gatherers (WHG),
 3209 Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers (CHG) and Anatolian farmers (ANA).
 3210 Row one shows zoomed Manhattan plots of the p-values for each ancestry (significant SNPs sized by
 3211 their selection coefficients), and row two shows allele trajectories for the top SNPs across all
 3212 ancestries.

3213 This peak spanned the region chr14:103239629-106248400, with significant SNPs including:

- 3214 • rs7144298 (*IGHG3* - *AL122127.1*; ALL; p=2.19e-09; s=0.0122), associated with Blood
 3215 protein levels ⁴⁰.

3216

3217 Peak 18: SEMA6D



3218

3219

Figure S4a.36. Selection at the *SEMA6D* locus, spanning chr15:47665179-51534060. Results for the

3220

pan-ancestry analysis (ALL) plus the four marginal ancestries: Western hunter-gatherers (WHG),

3221

Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers (CHG) and Anatolian farmers (ANA).

3222

Row one shows zoomed Manhattan plots of the p-values for each ancestry (significant SNPs sized by

3223

their selection coefficients), and row two shows allele trajectories for the top SNPs across all

3224

ancestries.

3225

This peak spanned the region chr15:47665179-51534060, with significant SNPs including:

3226

- rs1426654 (*SLC24A5*; WHG; $p=1.25e-19$; $s=0.0305$), associated with Eye colour; Eye colour (brightness); Eye colour (saturation); Hair colour; Iris colour (b^* coordinate); Skin pigmentation; Skin reflectance (Melanin index); Skin, hair and eye pigmentation (multivariate analysis)^{76,78,80,154,155}.

3227

3228

3229

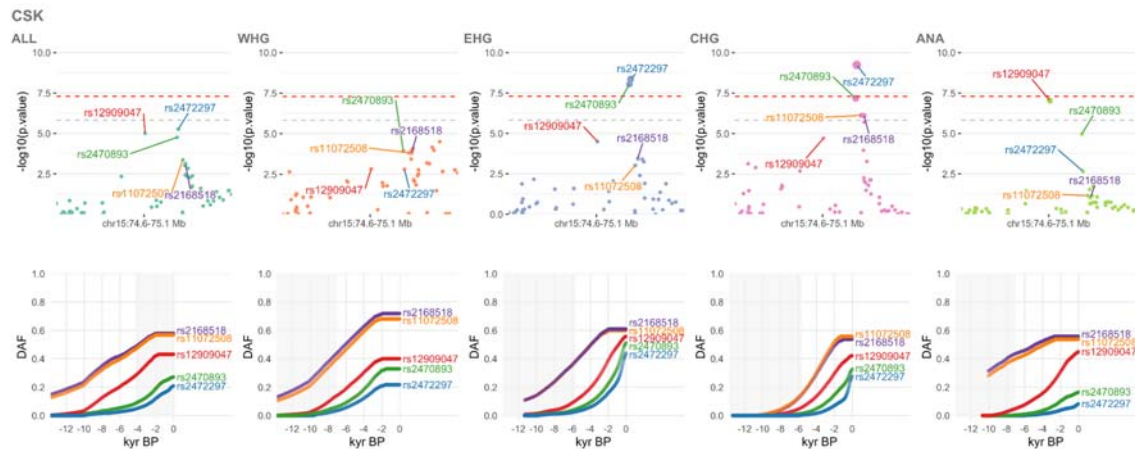
- rs61998565 (*SEMA6D*, *AC023905.1*; EHG; $p=1.19e-10$; $s=0.023$), associated with Heel bone mineral density^{45,58,109}.

3230

3231

3232

3233 Peak 19: CSK



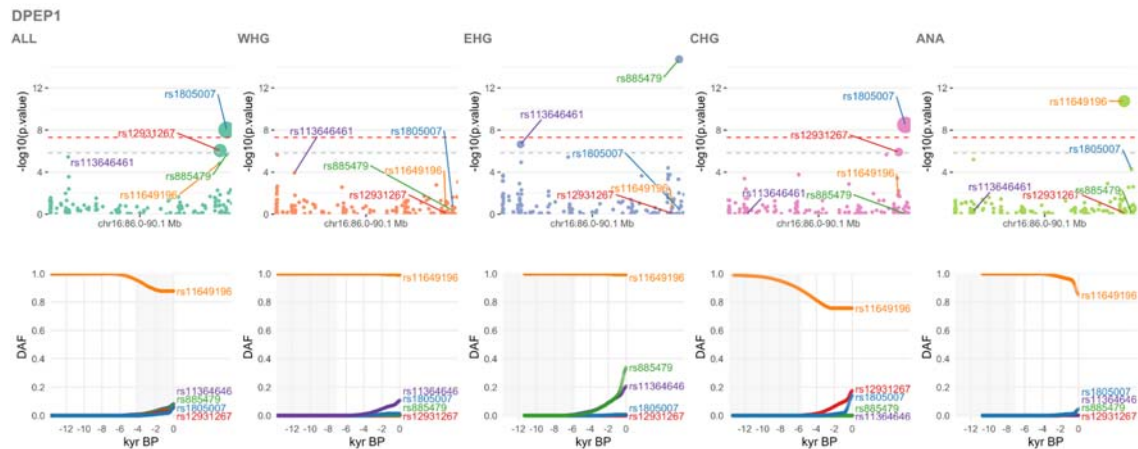
3234 **Figure S4a.37.** Selection at the CSK locus, spanning chr15:74607429-75101530. Results for the
 3235 pan-ancestry analysis (ALL) plus the four marginal ancestries: Western hunter-gatherers (WHG),
 3236 Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers (CHG) and Anatolian farmers (ANA).
 3237 Row one shows zoomed Manhattan plots of the p-values for each ancestry (significant SNPs sized by
 3238 their selection coefficients), and row two shows allele trajectories for the top SNPs across all
 3239 ancestries.
 3240

3241 This peak spanned the region chr15:74607429-75101530, with significant SNPs including:

- 3242 • rs2472297 (*CYP1A1 - CYP1A2*; CHG; $p=5.80e-10$; $s=0.0304$), associated with Alcohol
 3243 consumption (drinks per week); Alcohol consumption (drinks per week) (MTAG); Bitter
 3244 beverage consumption; Bitter non-alcoholic beverage consumption; Caffeine metabolism
 3245 (plasma 1,3,7-trimethylxanthine (caffeine) level); Caffeine metabolism (plasma 1,7-
 3246 dimethylxanthine (paraxanthine) to 1,3,7-trimethylxanthine (caffeine) ratio); Coffee
 3247 consumption; Coffee consumption (cups per day); Estimated glomerular filtration rate;
 3248 Estimated glomerular filtration rate in non-diabetics; Plasma clozapine levels in treatment-
 3249 resistant schizophrenia; Predicted visceral adipose tissue; Tea consumption; Urate
 3250 levels; Urinary albumin excretion; Urinary albumin-to-creatinine ratio; Urinary potassium
 3251 excretion; Urinary sodium excretion ^{86,91,133,156–165}.
- 3252 • rs2470893 (*CYP1A1 - CYP1A2*; EHG; $p=8.93e-09$; $s=0.0217$), associated with Blood
 3253 urea nitrogen levels; Caffeine consumption; Caffeine metabolism (plasma 1,3,7-
 3254 trimethylxanthine (caffeine) level); Caffeine metabolism (plasma 1,7-dimethylxanthine
 3255 (paraxanthine) to 1,3,7-trimethylxanthine (caffeine) ratio); Coffee consumption;
 3256 Microalbuminuria; Platelet distribution width; Urinary albumin excretion (no hypertensive
 3257 medication); Urinary albumin-to-creatinine ratio ^{83,133,157–159,166–168}.
- 3258 • rs12909047 (*AC012435.2,AC012435.3,UBL7-AS1*; ANA; $p=9.39e-08$; $s=0.019$),
 3259 associated with Caffeine metabolism (plasma 1,3,7-trimethylxanthine (caffeine) level);
 3260 Caffeine metabolism (plasma 1,3-dimethylxanthine (theophylline) level) ¹⁵⁸.
- 3261 • rs11072508 (*CYP1A2 - CSK*; CHG; $p=7.45e-07$; $s=0.0163$), associated with
 3262 Cardiovascular disease; Medication use (agents acting on the renin-angiotensin system)
 3263 ^{58,169}.
- 3264 • rs2168518 (*MIR4513,CSK*; CHG; $p=7.73e-07$; $s=0.0163$), associated with Medication use
 3265 (calcium channel blockers) ¹⁶⁹.

3266

3267 Peak 20: DPEP1



3268
3269
3270
3271
3272
3273
3274

Figure S4a.38. Selection at the *DPEP1* locus, spanning chr16:85972609-90084560. Results for the pan-ancestry analysis (ALL) plus the four marginal ancestries: Western hunter-gatherers (WHG), Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers (CHG) and Anatolian farmers (ANA). Row one shows zoomed Manhattan plots of the p-values for each ancestry (significant SNPs sized by their selection coefficients), and row two shows allele trajectories for the top SNPs across all ancestries.

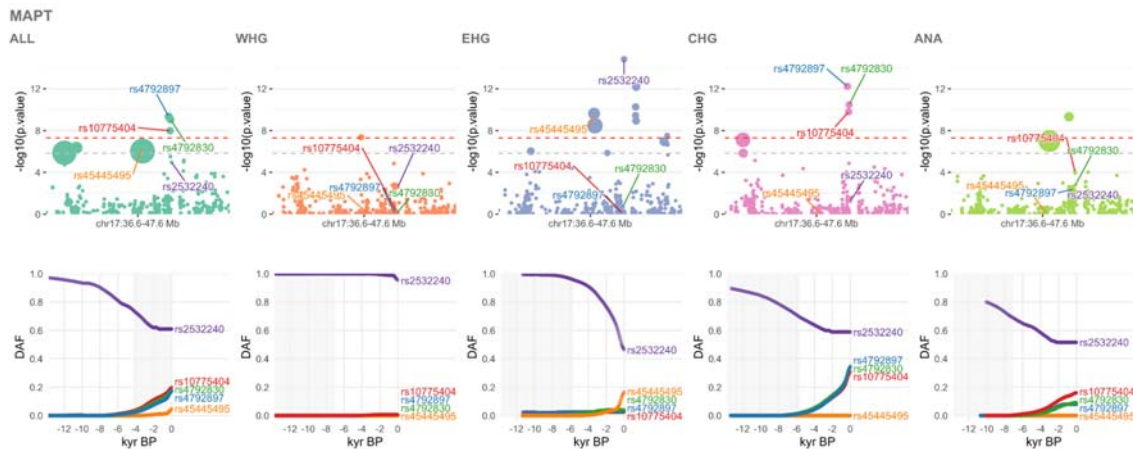
3275
3276
3277
3278
3279
3280
3281
3282
3283
3284
3285
3286
3287
3288
3289
3290

This peak spanned the region chr16:85972609-90084560, with significant SNPs including:

- rs885479 (*AC092143.1,MC1R*; EHG; $p=1.92e-15$; $s=0.0284$), associated with Blond vs. brown/black hair colour ⁵⁷.
- rs11649196 (*ZNF276*; ANA; $p=1.79e-11$; $s=-0.0457$), associated with Keratinocyte cancer (MTAG) ⁵⁹.
- rs1805007 (*AC092143.1,MC1R*; CHG; $p=3.22e-09$; $s=0.0648$), associated with Balding type 1; Basal cell carcinoma; Blond vs. brown hair colour; Blond vs. brown/black hair colour; Brown vs. black hair colour; Cutaneous squamous cell carcinoma; Freckles; Hair colour; Hair morphology traits; Keratinocyte cancer (MTAG); Melanoma; Non-melanoma skin cancer; Red vs non-red hair colour; Red vs. brown/black hair colour; Skin pigmentation traits; Skin sensitivity to sun; Squamous cell carcinoma; Sunburns; Tanning 49–59.
- rs113646461 (*AC092723.5 - AC092723.4*; EHG; $p=2.28e-07$; $s=0.0273$), associated with Monocyte percentage of white cells ⁸³.
- rs12931267 (*FANCA*; ALL; $p=8.87e-07$; $s=0.0506$), associated with Freckling; Hair colour; Hair morphology traits; Skin pigmentation traits; Skin sensitivity to sun ^{56,77,170}.

3291

3292 Peak 21: MAPT



3293 **Figure S4a.39.** Selection at the *MAPT* locus, spanning chr17:36615969-47588760. Results for the
 3294 pan-ancestry analysis (ALL) plus the four marginal ancestries: Western hunter-gatherers (WHG),
 3295 Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers (CHG) and Anatolian farmers (ANA).
 3296 Row one shows zoomed Manhattan plots of the p-values for each ancestry (significant SNPs sized by
 3297 their selection coefficients), and row two shows allele trajectories for the top SNPs across all
 3298 ancestries.
 3299

3300 This peak spanned the region chr17:36615969-47588760, with significant SNPs including:

- 3301 ● rs2532240 (*KANSL1*; EHG; $p=1.51e-15$; $s=-0.0217$), associated with White matter
 3302 microstructure (axial diffusivities) ¹⁷¹.
- 3303 ● rs4792897 (*MAPT*; CHG; $p=5.94e-13$; $s=0.0248$), associated with Snoring ⁶⁰.
- 3304 ● rs10775404 (*KANSL1*; CHG; $p=1.60e-10$; $s=0.0248$), associated with Reaction time ⁷⁴.
- 3305 ● rs45445495 (*SLC4A1*; EHG; $p=2.34e-10$; $s=0.0448$), associated with Appendicular lean
 3306 mass ¹⁴⁹.

3307
 3308
 3309
 3310

3311 Discussion

3312 Using our ancient genomic panel, we sought to identify phenotype-associated variants that
3313 have evidence for directional selection over the last 13,000 years. To estimate allele
3314 frequency trajectories and selection coefficients of trait-associated variants through time, we
3315 used the software *CLUES*¹ which can perform inference of allele frequency trajectories
3316 using marginal trees sampled from a reconstruction of an ancestral recombination graph
3317 (ARG)² for a set of genomic sequences, in combination with genotype likelihoods from
3318 serially sampled ancient DNA (aDNA).

3319

3320 Our results show that the incorporation of ancient DNA considerably improves our power to
3321 detect variants under selection, compared to a method that only uses the ARG inferred from
3322 present-day data alone. Using genomes from the 1,000 Genomes Project project
3323 (populations GBR, FIN and TSI), we inferred allele trajectories and selection coefficients for
3324 35,592 phenotype-associated variants, ascertained from the GWAS Catalog⁴, along with an
3325 equal number of putatively neutral “control” variants. Our analysis identified no genome-wide
3326 significant selective sweeps ($p < 5e-8$) using present-day data alone. However, the trait-
3327 associated variants were significantly enriched for evidence of selection when compared to
3328 the control group (Wilcoxon signed-rank test, $p < 2.2e-16$).

3329 Pan-ancestry selection

3330 In contrast to patterns observed in present-day genomes, our selection analysis based on a
3331 time-series of ancient DNA genotype probabilities identified 11 genome-wide significant ($p <$
3332 $5e-8$) selective sweeps in the GWAS variants, and none in the control group; consistent with
3333 widespread selection acting on trait-associated variants. This analysis confirms many of the
3334 previously reported selection loci in West Eurasians, identified from present-day and ancient
3335 DNA^{16,172–174}, and reveals novel selective sweeps, while refining the temporal dynamics of
3336 the selected alleles.

3337

3338 The strongest overall signal of selection in the pan-ancestry analysis is at the *LCT / MCM6*
3339 locus (rs4988235; $p=9.86e-31$; $s=0.020$), the derived allele of which is casual for lactase
3340 persistence^{35,36}. The inferred trajectory indicates that this allele began rising in frequency c.
3341 6,000 years ago, and has continued to rise in frequency up to the present (Supplementary
3342 Figure S4a.5).

3343

3344 We find a strong signal of selection at the *FADS1* (rs174546; $p=2.65e-10$; $s=0.013$) and
3345 *FADS2* (rs174581; $p=1.87e-10$; $s=0.013$) locus, associated with fatty acid metabolism

3346 ^{116,118,121,175–177}. The trajectories for these variants indicate a rise in frequency, beginning
3347 around 13,000 years ago, and continuing up until c. 2,000 years ago, after which their
3348 frequencies plateaued (Supplementary Figure S4a.11). In contrast to earlier findings ¹⁶, we
3349 do not detect a significant signal of selection at the *DHCR7* and *NADSYN1* locus, associated
3350 with vitamin D levels (most significant SNP rs4423214; $p=5.54e-03$; $s=-0.006$).

3351

3352 We detect an 8 megabase (Mb) wide selection sweep signal in chromosome 6 (chr6:25.2-
3353 33.5 Mb), spanning the human leukocyte antigen (HLA) region. The selection trajectories of
3354 the variants within this locus support multiple independent sweeps, occurring at different
3355 times and with differing intensities. The strongest signal of selection at this locus in the pan-
3356 ancestry analysis is at an intergenic variant, located between *HLA-A* and *HLA-W*
3357 (rs7747253; $p=8.86e-17$; $s=-0.018$), associated with heel bone mineral density ⁴⁵, the derived
3358 allele of which rapidly reduced in frequency, beginning c. 8,000 years ago (Supplementary
3359 Figure S4a.10). In contrast, the signal of selection at *C2* (rs9267677; $p=9.82e-14$; $s=$
3360 0.04463), also found within this sweep, and associated with educational attainment ⁶⁶, shows
3361 a gradual increase in frequency beginning c. 4,000 years ago, before rising more rapidly c.
3362 1,000 years ago; highlighting the complex temporal dynamics of selection at the HLA locus.

3363

3364 We also identify selection signals at the *SLC22A4* (rs35260072; $p=1.15e-10$; $s=0.018$) and
3365 *RAPGEF6* (rs11950815; $p=1.82e-12$; $s=0.021$) loci, associated with asthma ¹⁴⁸ and itch
3366 intensity from mosquito bites ³⁴, and find that these alleles have been steadily rising in
3367 frequency, beginning c. 8,000 years ago (Supplementary Figure S4a.8). However, we find
3368 that the frequency of rs1050152 plateaued c. 1,500 years ago, contrary to previous reports
3369 suggesting a recent rise in frequency ¹⁶. Similarly, we detect selection at the *HECTD4*
3370 (rs11066188; $p=3.02e-16$; $s=0.020$) and *ATXN2* (rs653178; $p=1.92e-15$; $s=0.019$) locus,
3371 associated with celiac disease and rheumatoid arthritis ⁴⁸, which has been rising in frequency
3372 for c. 9,000 years (Supplementary Figure S4a.12), also contrary to previous reports of a
3373 more recent rise in frequency ¹⁶.

3374

3375 We detect strong selection at the *SLC45A2* (rs35395; $p=4.13e-23$; $s=0.022$) locus,
3376 associated with skin pigmentation ^{76,178}, and find that the selected allele began rising in
3377 frequency c. 13,000 years ago, after which it plateaued at high frequency c. 2,000 years ago.
3378 This is similar to the selection trajectory at the independent *GRM5* (rs7119749; $p=8.54e-09$;
3379 $s=0.011$) locus, which plateaued at medium-high frequency at approximately the same time.
3380 We also detect strong selection at the *SLC24A5* (rs1426654; $p=6.45e-09$; $s=0.019$) locus,
3381 and find that the selected allele, also associated with skin pigmentation ^{76,154}, began rising in
3382 frequency even earlier than *SLC45A2*, and reached near fixation c. 3,500 years ago.

3383

3384 We further detect strong selection in an 11 Mb sweep in chromosome 17 (chr17:36.6-47.5),
3385 spanning the 17q21.31 locus, a 900-kb inversion polymorphism^{179,180}. The strongest signal
3386 of selection in this sweep is at *MAPT* (rs4792897; $p=4.65e-10$; $s=0.03$), associated with
3387 snoring⁶⁰, the trajectory for which indicates a steady increase in frequency, beginning c.
3388 7,000 years ago (Supplementary Figure S4a.14).

3389 Ancestry stratified selection trajectories

3390 To account for population structure in our samples, we also applied a novel chromosome
3391 painting technique, based on inference of a sample's nearest neighbours in the marginal
3392 trees of an ARG that contains individuals classified into different ancient populations
3393 (Supplementary Note 3i). This method allows us to accurately assign ancestral population
3394 labels to haplotypes found in both ancient and present-day individuals. By conditioning our
3395 selection analyses on these haplotype backgrounds, we can infer the selection trajectories of
3396 GWAS risk alleles in a manner that is approximately invariant to change in the admixture
3397 proportions through time. These ancestry specific allele trajectories reveal many novel
3398 aspects about the dynamic interplay between selection and admixture in West Eurasia
3399 throughout the Holocene. We find that the allele trajectories of directionally selected sites
3400 become much more apparent once we perform this ancestry partitioning. We often find
3401 variants with strong allele frequency changes in one ancestral population but not another,
3402 and analysing all ancient individuals without accounting for their ancestry composition leads
3403 to a decrease in our ability to identify selected variants, and a blurring of the temporal signal
3404 of allele frequency changes.

3405

3406 When conditioned on one of our four marginal ancestries—Western hunter-gatherers
3407 (WHG), Eastern hunter-gatherers (EHG), Caucasus hunter-gatherers (CHG) and Anatolian
3408 farmers (ANA)—we find 21 genome-wide significant selection peaks (substantially more than
3409 in our pan-ancestry analysis), indicating that admixture between ancestral populations has
3410 masked evidence of selection at many loci. Furthermore, we find that some strong signals of
3411 selection identified in the pan-ancestry analysis are driven by sweeps in the marginal
3412 ancestries which substantially differ in their most significant SNPs, suggesting that multiple
3413 selected alleles may be common within genome-wide significant sweep loci.

3414

3415 For example, in the pan-ancestry analysis, the sweep at *MCM6* is led by rs4988235,
3416 consistent with the widespread interpretation that selection has acted upon the lactase
3417 persistence phenotype. However, in the ancestry stratified analysis, this selection signal is
3418 primarily driven by sweeps in the EHG and CHG ancestral backgrounds, which differ in their

3419 most significant SNPs ([Supplementary Figure S4a.21](#)). The strongest sweep signal in all
3420 marginal ancestries at this locus is in the EHG background; however, conditional on that
3421 background, rs1246580 (*R3HDM1*) is the most significant SNP (associated with blood
3422 protein levels⁴⁰ and mosquito bite size³⁴). In CHG, rs4988235 is the most significant SNP,
3423 but there is no evidence for selection at rs1246580 in this background. Conversely, in WHG,
3424 we find no evidence for selection at rs4988235, and instead find evidence for strong
3425 selection at rs1246580 occurring in the last c. 2,000 years. Despite the highly studied nature
3426 of this locus, a satisfactory explanation for the observed strength of selection has remained
3427 elusive^{181–184}.

3428

3429 In comparison, the sweep at the *SLC45A2* locus shows a much simpler pattern, in which all
3430 marginal ancestries show broad agreement at this locus ([Supplementary Figure S4a.25](#)).
3431 Where they differ lies primarily in the timing of their frequency rises. The ANA ancestry
3432 background shows the earliest evidence for selection at *SLC45A2*, followed by EHG and
3433 WHG, beginning around c. 10,000 years ago, and CHG c. 2,000 years later. In all ancestry
3434 backgrounds except WHG, the selected haplotypes reach near fixation by c. 3,000 years
3435 ago, whilst the WHG haplotype background contains the majority of the ancestral alleles still
3436 segregating in present-day Europeans.

3437

3438 At the *FADS2* locus, the strong signal of selection in the pan-ancestry analysis is driven
3439 primarily by a sweep occurring on the EHG haplotypic background. Interestingly, we find no
3440 evidence for selection at this locus in the WHG background, and most of the frequency rise
3441 in the EHG background occurs after their admixture with CHG, where the selected alleles
3442 were already at close to present-day frequencies.

3443

3444 References

- 3445 1. Stern, A. J., Wilton, P. R. & Nielsen, R. An approximate full-likelihood method for
3446 inferring selection and allele frequency trajectories from DNA sequence data. *PLoS*
3447 *Genet.* **15**, e1008384 (2019).
- 3448 2. Speidel, L., Forest, M., Shi, S. & Myers, S. R. A method for genome-wide genealogy
3449 estimation for thousands of samples. *Nat. Genet.* **51**, 1321–1329 (2019).
- 3450 3. Köster, J. & Rahmann, S. Snakemake--a scalable bioinformatics workflow engine.
3451 *Bioinformatics* **28**, 2520–2522 (2012).

- 3452 4. Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide
3453 association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.*
3454 **47**, D1005–D1012 (2019).
- 3455 5. Jupp, S., Burdett, T., Leroy, C. & Parkinson, H. E. A new Ontology Lookup Service at
3456 EMBL-EBI. in *SWAT4LS* 118–119 (ceur-ws.org, 2015).
- 3457 6. Yates, A. *et al.* The Ensembl REST API: Ensembl Data for Any Language.
3458 *Bioinformatics* **31**, 143–145 (2015).
- 3459 7. Aken, B. L. *et al.* The Ensembl gene annotation system. *Database* **2016**, (2016).
- 3460 8. Kelleher, J., Etheridge, A. M. & McVean, G. Efficient Coalescent Simulation and
3461 Genealogical Analysis for Large Sample Sizes. *PLoS Comput. Biol.* **12**, e1004842
3462 (2016).
- 3463 9. 1000 Genomes Project Consortium *et al.* A global reference for human genetic
3464 variation. *Nature* **526**, 68–74 (2015).
- 3465 10. Vilella, A. J. *et al.* EnsemblCompara GeneTrees: Complete, duplication-aware
3466 phylogenetic trees in vertebrates. *Genome Res.* **19**, 327–335 (2009).
- 3467 11. Rasmussen, M. D., Hubisz, M. J., Gronau, I. & Siepel, A. Genome-wide
3468 inference of ancestral recombination graphs. *PLoS Genet.* **10**, e1004342 (2014).
- 3469 12. Hein, J., Schierup, M. & Wiuf, C. *Gene Genealogies, Variation and Evolution:*
3470 *A primer in coalescent theory.* (Oxford University Press, USA, 2004).
- 3471 13. Stern, A. J., Speidel, L., Zaitlen, N. A. & Nielsen, R. Disentangling selection
3472 on genetically correlated polygenic traits via whole-genome genealogies. *Am. J. Hum.*
3473 *Genet.* (2021) doi:10.1016/j.ajhg.2020.12.005.
- 3474 14. 1000 Genomes Project Consortium *et al.* A map of human genome variation
3475 from population-scale sequencing. *Nature* **467**, 1061–1073 (2010).
- 3476 15. Haller, T., Tasa, T. & Metspalu, A. Manhattan Harvester and Cropper: a
3477 system for GWAS peak detection. *BMC Bioinformatics* **20**, 22 (2019).
- 3478 16. Mathieson, I. *et al.* Genome-wide patterns of selection in 230 ancient
3479 Eurasians. *Nature* **528**, 499–503 (2015).

- 3480 17. Li, H. A statistical framework for SNP calling, mutation discovery, association
3481 mapping and population genetical parameter estimation from sequencing data.
3482 *Bioinformatics* **27**, 2987–2993 (2011).
- 3483 18. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing
3484 genomic features. *Bioinformatics* **26**, 841–842 (2010).
- 3485 19. Cock, P. J. A. *et al.* Biopython: freely available Python tools for computational
3486 molecular biology and bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).
- 3487 20. Anaconda Software Distribution. *Anaconda Documentation* (2020).
- 3488 21. Harris, C. R. *et al.* Array programming with NumPy. *Nature* **585**, 357–362
3489 (2020).
- 3490 22. McKinney, W. & Others. Data structures for statistical computing in python. in
3491 *Proceedings of the 9th Python in Science Conference* vol. 445 51–56 (Austin, TX,
3492 2010).
- 3493 23. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools.
3494 *Bioinformatics* **25**, 2078–2079 (2009).
- 3495 24. Van Rossum, G. & Drake, F. L. *Python 3 Reference Manual: (Python*
3496 *Documentation Manual Part 2)*. (CreateSpace Independent Publishing Platform, 2009).
- 3497 25. R Core Team. R: A Language and Environment for Statistical Computing.
3498 (2019).
- 3499 26. Haider, S. *et al.* A bedr way of genomic interval processing. *Source Code Biol.*
3500 *Med.* **11**, 14 (2016).
- 3501 27. Wickham, H., François, R., Henry, L. & Müller, K. dplyr: A Grammar of Data
3502 Manipulation. (2019).
- 3503 28. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*. (Springer, 2016).
- 3504 29. Petukhov, V., van den Brand, T. & Biederstedt, E. ggrastr: Raster Layers for
3505 ‘ggplot2’. (2020).
- 3506 30. Slowikowski, K. ggrepel: Automatically Position Non-Overlapping Text Labels
3507 with ‘ggplot2’. (2020).

- 3508 31. Wilke, C. O. ggrridges: Ridgeline Plots in 'ggplot2'. (2018).
- 3509 32. Wickham, H. stringr: Simple, Consistent Wrappers for Common String
3510 Operations. (2019).
- 3511 33. Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing
3512 in Python. *Nat. Methods* **17**, 261–272 (2020).
- 3513 34. Jones, A. V. *et al.* GWAS of self-reported mosquito bite size, itch intensity and
3514 attractiveness to mosquitoes implicates immune-related predisposition loci. *Hum. Mol.*
3515 *Genet.* **26**, 1391–1406 (2017).
- 3516 35. Enattah, N. S. *et al.* Identification of a variant associated with adult-type
3517 hypolactasia. *Nat. Genet.* **30**, 233–237 (2002).
- 3518 36. Bersaglieri, T. *et al.* Genetic signatures of strong recent positive selection at
3519 the lactase gene. *Am. J. Hum. Genet.* **74**, 1111–1120 (2004).
- 3520 37. Hill, W. D. *et al.* Genome-wide analysis identifies molecular systems and 149
3521 genetic loci associated with income. *Nat. Commun.* **10**, 5741 (2019).
- 3522 38. Shungin, D. *et al.* New genetic loci link adipose and insulin biology to body fat
3523 distribution. *Nature* **518**, 187–196 (2015).
- 3524 39. Winkler, T. W. *et al.* The Influence of Age and Sex on Genetic Associations
3525 with Adult Body Size and Shape: A Large-Scale Genome-Wide Interaction Study. *PLoS*
3526 *Genet.* **11**, e1005378 (2015).
- 3527 40. Sun, B. B. *et al.* Genomic atlas of the human plasma proteome. *Nature* **558**,
3528 73–79 (2018).
- 3529 41. Emilsson, V. *et al.* Co-regulatory networks of human serum proteins link
3530 genetics to disease. *Science* **361**, 769–773 (2018).
- 3531 42. Hoffmann, T. J. *et al.* A Large Multiethnic Genome-Wide Association Study of
3532 Adult Body Mass Index Identifies Novel Loci. *Genetics* **210**, 499–515 (2018).
- 3533 43. Ahsan, M. *et al.* The relative contribution of DNA methylation and genetic
3534 variants on protein biomarkers for human diseases. *PLoS Genet.* **13**, e1007005 (2017).
- 3535 44. Zhu, Z. *et al.* A genome-wide cross-trait analysis from UK Biobank highlights

3536 the shared genetic architecture of asthma and allergic diseases. *Nat. Genet.* **50**, 857–
3537 864 (2018).

3538 45. Morris, J. A. *et al.* An atlas of genetic influences on osteoporosis in humans
3539 and mice. *Nat. Genet.* **51**, 258–266 (2018).

3540 46. Hageaars, S. P. *et al.* Genetic prediction of male pattern baldness. *PLoS*
3541 *Genet.* **13**, e1006594 (2017).

3542 47. Li, Y. *et al.* Genome-Wide Association Studies of Metabolites in Patients with
3543 CKD Identify Multiple Loci and Illuminate Tubular Transport Mechanisms. *J. Am. Soc.*
3544 *Nephrol.* **29**, 1513–1524 (2018).

3545 48. Gutierrez-Achury, J. *et al.* Functional implications of disease-specific variants
3546 in loci jointly associated with coeliac disease and rheumatoid arthritis. *Hum. Mol. Genet.*
3547 **25**, 180–190 (2016).

3548 49. Sulem, P. *et al.* Genetic determinants of hair, eye and skin pigmentation in
3549 Europeans. *Nat. Genet.* **39**, 1443–1452 (2007).

3550 50. Nan, H. *et al.* Genome-wide association study identifies novel alleles
3551 associated with risk of cutaneous basal cell carcinoma and squamous cell carcinoma.
3552 *Hum. Mol. Genet.* **20**, 3718–3724 (2011).

3553 51. Zhang, M. *et al.* Genome-wide association studies identify several new loci
3554 associated with pigmentation traits and skin cancer risk in European Americans. *Hum.*
3555 *Mol. Genet.* **22**, 2948–2959 (2013).

3556 52. Chahal, H. S. *et al.* Genome-wide association study identifies novel
3557 susceptibility loci for cutaneous squamous cell carcinoma. *Nat. Commun.* **7**, 12048
3558 (2016).

3559 53. Chahal, H. S. *et al.* Genome-wide association study identifies 14 novel risk
3560 alleles associated with basal cell carcinoma. *Nat. Commun.* **7**, 12510 (2016).

3561 54. Ransohoff, K. J. *et al.* Two-stage genome-wide association study identifies a
3562 novel susceptibility locus associated with melanoma. *Oncotarget* **8**, 17586–17592
3563 (2017).

- 3564 55. Hysi, P. G. *et al.* Genome-wide association meta-analysis of individuals of
3565 European ancestry identifies new loci explaining a substantial fraction of hair color
3566 variation and heritability. *Nat. Genet.* **50**, 652–656 (2018).
- 3567 56. Galván-Femenía, I. *et al.* Multitrait genome association analysis identifies new
3568 susceptibility genes for human anthropometric variation in the GCAT cohort. *J. Med.*
3569 *Genet.* **55**, 765–778 (2018).
- 3570 57. Morgan, M. D. *et al.* Genome-wide study of hair colour in UK Biobank explains
3571 most of the SNP heritability. *Nat. Commun.* **9**, 5271 (2018).
- 3572 58. Kichaev, G. *et al.* Leveraging Polygenic Functional Enrichment to Improve
3573 GWAS Power. *Am. J. Hum. Genet.* **104**, 65–75 (2019).
- 3574 59. Liyanage, U. E. *et al.* Combined analysis of keratinocyte cancers identifies
3575 novel genome-wide loci. *Hum. Mol. Genet.* **28**, 3148–3160 (2019).
- 3576 60. Jansen, P. R. *et al.* Genome-wide analysis of insomnia in 1,331,010
3577 individuals identifies new risk loci and functional pathways. *Nat. Genet.* **51**, 394–403
3578 (2019).
- 3579 61. Akinkuolie, A. O. *et al.* Group IIA Secretory Phospholipase A2, Vascular
3580 Inflammation, and Incident Cardiovascular Disease. *Arterioscler. Thromb. Vasc. Biol.*
3581 **39**, 1182–1190 (2019).
- 3582 62. Hill, W. D. *et al.* A combined analysis of genetically correlated traits identifies
3583 187 loci and a role for neurogenesis and myelination in intelligence. *Mol. Psychiatry* **24**,
3584 169–181 (2019).
- 3585 63. Legge, S. E. *et al.* Association of Genetic Liability to Psychotic Experiences
3586 With Neuropsychotic Disorders and Traits. *JAMA Psychiatry* **76**, 1256–1265 (2019).
- 3587 64. Okbay, A. *et al.* Genome-wide association study identifies 74 loci associated
3588 with educational attainment. *Nature* **533**, 539–542 (2016).
- 3589 65. Savage, J. E. *et al.* Genome-wide association meta-analysis in 269,867
3590 individuals identifies new genetic and functional links to intelligence. *Nat. Genet.* **50**,
3591 912–919 (2018).

- 3592 66. Lee, J. J. *et al.* Gene discovery and polygenic prediction from a genome-wide
3593 association study of educational attainment in 1.1 million individuals. *Nat. Genet.* **50**,
3594 1112–1121 (2018).
- 3595 67. Li, M. *et al.* Genome-wide association study of 1,5-anhydroglucitol identifies
3596 novel genetic loci linked to glucose metabolism. *Sci. Rep.* **7**, 2812 (2017).
- 3597 68. Schlosser, P. *et al.* Genetic studies of urinary metabolites illuminate
3598 mechanisms of detoxification and excretion in humans. *Nat. Genet.* **52**, 167–176
3599 (2020).
- 3600 69. Tikkanen, E. *et al.* Biological Insights Into Muscular Strength: Genetic
3601 Findings in the UK Biobank. *Sci. Rep.* **8**, 6451 (2018).
- 3602 70. Jones, S. E. *et al.* Genome-wide association analyses of chronotype in
3603 697,828 individuals provides insights into circadian rhythms. *Nat. Commun.* **10**, 343
3604 (2019).
- 3605 71. Ahola-Olli, A. V. *et al.* Genome-wide Association Study Identifies 27 Loci
3606 Influencing Concentrations of Circulating Cytokines and Growth Factors. *Am. J. Hum.*
3607 *Genet.* **100**, 40–50 (2017).
- 3608 72. Zhu, Z. *et al.* Shared genetic and experimental links between obesity-related
3609 traits and asthma subtypes in UK Biobank. *J. Allergy Clin. Immunol.* **145**, 537–549
3610 (2020).
- 3611 73. Yap, C. X. *et al.* Dissection of genetic variation and evidence for pleiotropy in
3612 male pattern baldness. *Nat. Commun.* **9**, 5407 (2018).
- 3613 74. Davies, G. *et al.* Study of 300,486 individuals identifies 148 independent
3614 genetic loci influencing general cognitive function. *Nat. Commun.* **9**, 2098 (2018).
- 3615 75. Zhao, B. *et al.* Genome-wide association analysis of 19,629 individuals
3616 identifies variants influencing regional brain volumes and refines their genetic co-
3617 architecture with cognitive and mental health traits. *Nat. Genet.* **51**, 1637–1644 (2019).
- 3618 76. Lona-Durazo, F. *et al.* Meta-analysis of GWA studies provides new insights on
3619 the genetic architecture of skin pigmentation in recently admixed populations. *BMC*

3620 *Genet.* **20**, 59 (2019).

3621 77. Liu, F. *et al.* Genetics of skin color variation in Europeans: genome-wide
3622 association studies with functional follow-up. *Hum. Genet.* **134**, 823–835 (2015).

3623 78. Adhikari, K. *et al.* A genome-wide association scan in admixed Latin
3624 Americans identifies loci influencing facial and scalp hair features. *Nat. Commun.* **7**,
3625 10815 (2016).

3626 79. Han, J. *et al.* A genome-wide association study identifies novel alleles
3627 associated with hair color and skin pigmentation. *PLoS Genet.* **4**, e1000074 (2008).

3628 80. Adhikari, K. *et al.* A GWAS in Latin Americans highlights the convergent
3629 evolution of lighter skin pigmentation in Eurasia. *Nat. Commun.* **10**, 358 (2019).

3630 81. Nan, H. *et al.* Genome-wide association study of tanning phenotype in a
3631 population of European ancestry. *J. Invest. Dermatol.* **129**, 2250–2257 (2009).

3632 82. Liu, J. Z. *et al.* Association analyses identify 38 susceptibility loci for
3633 inflammatory bowel disease and highlight shared genetic risk across populations. *Nat.*
3634 *Genet.* **47**, 979–986 (2015).

3635 83. Astle, W. J. *et al.* The Allelic Landscape of Human Blood Cell Trait Variation
3636 and Links to Common Complex Disease. *Cell* **167**, 1415–1429.e19 (2016).

3637 84. Shin, S.-Y. *et al.* An atlas of genetic influences on human blood metabolites.
3638 *Nat. Genet.* **46**, 543–550 (2014).

3639 85. Kristjansson, R. P. *et al.* A loss-of-function variant in ALOX15 protects against
3640 nasal polyps and chronic rhinosinusitis. *Nat. Genet.* **51**, 267–276 (2019).

3641 86. Hellwege, J. N. *et al.* Mapping eGFR loci to the renal transcriptome and
3642 phenome in the VA Million Veteran Program. *Nat. Commun.* **10**, 3842 (2019).

3643 87. Graham, S. E. *et al.* Sex-specific and pleiotropic effects underlying kidney
3644 function identified from GWAS meta-analysis. *Nat. Commun.* **10**, 1847 (2019).

3645 88. Franke, A. *et al.* Sequence variants in IL10, ARPC2 and multiple other loci
3646 contribute to ulcerative colitis susceptibility. *Nat. Genet.* **40**, 1319–1323 (2008).

3647 89. UK IBD Genetics Consortium *et al.* Genome-wide association study of

3648 ulcerative colitis identifies three new susceptibility loci, including the HNF4A region. *Nat.*
3649 *Genet.* **41**, 1330–1334 (2009).

3650 90. Lee, H.-S. *et al.* An Intergenic Variant rs9268877 Between HLA-DRA and
3651 HLA-DRB Contributes to the Clinical Course and Long-term Outcome of Ulcerative
3652 Colitis. *J. Crohns. Colitis* **12**, 1113–1121 (2018).

3653 91. Tin, A. *et al.* Target genes, variants, tissues and transcriptional pathways
3654 influencing human serum urate levels. *Nat. Genet.* **51**, 1459–1474 (2019).

3655 92. Cipriani, V. *et al.* Genome-wide association study of age-related macular
3656 degeneration identifies associated variants in the TNXB-FKBPL-NOTCH4 region of
3657 chromosome 6p21.3. *Hum. Mol. Genet.* **21**, 4138–4150 (2012).

3658 93. Weidinger, S. *et al.* A genome-wide association study of atopic dermatitis
3659 identifies loci with overlapping effects on asthma and psoriasis. *Hum. Mol. Genet.* **22**,
3660 4841–4856 (2013).

3661 94. Hill, W. D. *et al.* Genetic contributions to two special factors of neuroticism are
3662 associated with affluence, higher intelligence, better health, and longer life. *Mol.*
3663 *Psychiatry* **25**, 3034–3052 (2020).

3664 95. Luciano, M. *et al.* Association analysis in over 329,000 individuals identifies
3665 116 independent variants influencing neuroticism. *Nat. Genet.* **50**, 6–11 (2018).

3666 96. Sung, Y. J. *et al.* A Large-Scale Multi-ancestry Genome-wide Study
3667 Accounting for Smoking Behavior Identifies Multiple Significant Loci for Blood Pressure.
3668 *Am. J. Hum. Genet.* **102**, 375–400 (2018).

3669 97. Wootton, R. E. *et al.* Evidence for causal effects of lifetime smoking on risk for
3670 depression and schizophrenia: a Mendelian randomisation study. *Psychol. Med.* **50**,
3671 2435–2443 (2020).

3672 98. Gao, X. R., Huang, H., Nannini, D. R., Fan, F. & Kim, H. Genome-wide
3673 association analyses identify new loci influencing intraocular pressure. *Hum. Mol.*
3674 *Genet.* **27**, 2205–2213 (2018).

3675 99. Khawaja, A. P. *et al.* Genome-wide analyses identify 68 new loci associated

3676 with intraocular pressure and improve risk prediction for primary open-angle glaucoma.
3677 *Nat. Genet.* **50**, 778–782 (2018).

3678 100. Weiss, F. U. *et al.* Fucosyltransferase 2 (FUT2) non-secretor status and blood
3679 group B are associated with elevated serum lipase activity in asymptomatic subjects,
3680 and an increased risk for chronic pancreatitis: a genetic association study. *Gut* **64**, 646–
3681 656 (2015).

3682 101. Lieb, W. *et al.* Genome-wide association study for endothelial growth factors.
3683 *Circ. Cardiovasc. Genet.* **8**, 389–397 (2015).

3684 102. Folkersen, L. *et al.* Mapping of 79 loci for 83 plasma protein biomarkers in
3685 cardiovascular disease. *PLoS Genet.* **13**, e1006706 (2017).

3686 103. Sabater-Lleal, M. *et al.* Genome-Wide Association Transethnic Meta-Analyses
3687 Identifies Novel Associations Regulating Coagulation Factor VIII and von Willebrand
3688 Factor Plasma Levels. *Circulation* **139**, 620–635 (2019).

3689 104. Klarin, D. *et al.* Genetics of blood lipids among ~300,000 multi-ethnic
3690 participants of the Million Veteran Program. *Nat. Genet.* **50**, 1514–1523 (2018).

3691 105. Hoffmann, T. J. *et al.* A large electronic-health-record-based genome-wide
3692 study of serum lipids. *Nat. Genet.* **50**, 401–413 (2018).

3693 106. Kanai, M. *et al.* Genetic analysis of quantitative traits in the Japanese
3694 population links cell types to complex human diseases. *Nat. Genet.* **50**, 390–400
3695 (2018).

3696 107. Bi, X. *et al.* Common genetic variants have associations with human cortical
3697 brain regions and risk of schizophrenia. *Genet. Epidemiol.* **43**, 548–558 (2019).

3698 108. Surakka, I. *et al.* The impact of low-frequency and rare variants on lipid levels.
3699 *Nat. Genet.* **47**, 589–597 (2015).

3700 109. Kim, S. K. Identification of 613 new loci associated with heel bone mineral
3701 density and a polygenic risk score for bone mineral density, osteoporosis and fracture.
3702 *PLoS One* **13**, e0200785 (2018).

3703 110. Wojcik, G. L. *et al.* Genetic analyses of diverse populations improves

3704 discovery for complex traits. *Nature* **570**, 514–518 (2019).

3705 111. Mozaffarian, D. *et al.* Genetic loci associated with circulating phospholipid
3706 trans fatty acids: a meta-analysis of genome-wide association studies from the
3707 CHARGE Consortium. *Am. J. Clin. Nutr.* **101**, 398–406 (2015).

3708 112. Dorajoo, R. *et al.* A genome-wide association study of n-3 and n-6 plasma
3709 fatty acids in a Singaporean Chinese population. *Genes Nutr.* **10**, 53 (2015).

3710 113. Hu, Y. *et al.* Discovery and fine-mapping of loci associated with MUFAs
3711 through trans-ethnic meta-analysis in Chinese and European populations. *J. Lipid Res.*
3712 **58**, 974–981 (2017).

3713 114. Tabassum, R. *et al.* Genetic architecture of human plasma lipidome and its
3714 link to cardiovascular disease. *Nat. Commun.* **10**, 4329 (2019).

3715 115. Teslovich, T. M. *et al.* Biological, clinical and population relevance of 95 loci
3716 for blood lipids. *Nature* **466**, 707–713 (2010).

3717 116. Willer, C. J. *et al.* Discovery and refinement of loci associated with lipid levels.
3718 *Nat. Genet.* **45**, 1274–1283 (2013).

3719 117. Draisma, H. H. M. *et al.* Genome-wide association study identifies novel
3720 genetic variants contributing to variation in blood metabolite levels. *Nat. Commun.* **6**,
3721 7208 (2015).

3722 118. Ligthart, S. *et al.* Bivariate genome-wide association study identifies novel
3723 pleiotropic loci for lipids and inflammation. *BMC Genomics* **17**, 443 (2016).

3724 119. van Setten, J. *et al.* Genome-wide association meta-analysis of 30,000
3725 samples identifies seven novel loci for quantitative ECG traits. *Eur. J. Hum. Genet.* **27**,
3726 952–962 (2019).

3727 120. Bentley, A. R. *et al.* Multi-ancestry genome-wide gene-smoking interaction
3728 study of 387,272 individuals identifies new loci associated with serum lipids. *Nat. Genet.*
3729 **51**, 636–648 (2019).

3730 121. Gallois, A. *et al.* A comprehensive study of metabolite genetics reveals strong
3731 pleiotropy and heterogeneity across time and context. *Nat. Commun.* **10**, 4788 (2019).

- 3732 122. Olafsdottir, T. A. *et al.* Eighty-eight variants highlight the role of T cell
3733 regulation and airway remodeling in asthma pathogenesis. *Nat. Commun.* **11**, 393
3734 (2020).
- 3735 123. de Vries, P. S. *et al.* Comparison of HapMap and 1000 Genomes Reference
3736 Panels in a Large-Scale Genome-Wide Association Study. *PLoS One* **12**, e0167742
3737 (2017).
- 3738 124. Hoffmann, T. J. *et al.* Genome-wide association analyses using electronic
3739 health records identify new loci influencing blood pressure variation. *Nat. Genet.* **49**,
3740 54–64 (2017).
- 3741 125. Tian, C. *et al.* Genome-wide association and HLA region fine-mapping studies
3742 identify susceptibility loci for multiple common infections. *Nat. Commun.* **8**, 599 (2017).
- 3743 126. Pilling, L. C. *et al.* Human longevity: 25 genetic loci associated in 389,166 UK
3744 biobank participants. *Aging* **9**, 2504–2520 (2017).
- 3745 127. Todd, J. A. *et al.* Robust associations of four new chromosome regions from
3746 genome-wide analyses of type 1 diabetes. *Nat. Genet.* **39**, 857–864 (2007).
- 3747 128. Wellcome Trust Case Control Consortium. Genome-wide association study of
3748 14,000 cases of seven common diseases and 3,000 shared controls. *Nature* **447**, 661–
3749 678 (2007).
- 3750 129. Cooper, J. D. *et al.* Meta-analysis of genome-wide association study data
3751 identifies additional type 1 diabetes risk loci. *Nat. Genet.* **40**, 1399–1401 (2008).
- 3752 130. Liu, C. *et al.* Meta-analysis identifies common and rare variants influencing
3753 blood pressure and overlapping with metabolic trait loci. *Nat. Genet.* **48**, 1162–1170
3754 (2016).
- 3755 131. Feitosa, M. F. *et al.* Novel genetic associations for blood pressure identified
3756 via gene-alcohol interaction in up to 570K individuals across multiple ancestries. *PLoS*
3757 *One* **13**, e0198166 (2018).
- 3758 132. Zhou, W. *et al.* Efficiently controlling for case-control imbalance and sample
3759 relatedness in large-scale genetic association studies. *Nat. Genet.* **50**, 1335–1341

3760 (2018).

3761 133. Wuttke, M. *et al.* A catalog of genetic loci associated with kidney function from
3762 analyses of a million individuals. *Nat. Genet.* **51**, 957–972 (2019).

3763 134. Newton-Cheh, C. *et al.* Genome-wide association study identifies eight loci
3764 associated with blood pressure. *Nat. Genet.* **41**, 666–676 (2009).

3765 135. Dubois, P. C. A. *et al.* Multiple common variants for celiac disease influencing
3766 immune gene expression. *Nat. Genet.* **42**, 295–302 (2010).

3767 136. Köttgen, A. *et al.* New loci associated with kidney function and chronic kidney
3768 disease. *Nat. Genet.* **42**, 376–384 (2010).

3769 137. Zhernakova, A. *et al.* Meta-analysis of genome-wide association studies in
3770 celiac disease and rheumatoid arthritis identifies fourteen non-HLA shared loci. *PLoS*
3771 *Genet.* **7**, e1002004 (2011).

3772 138. Wain, L. V. *et al.* Genome-wide association study identifies six new loci
3773 influencing pulse pressure and mean arterial pressure. *Nat. Genet.* **43**, 1005–1011
3774 (2011).

3775 139. Köttgen, A. *et al.* Genome-wide association analyses identify 18 new loci
3776 associated with serum urate concentrations. *Nat. Genet.* **45**, 145–154 (2013).

3777 140. Medici, M. *et al.* Identification of novel genetic Loci associated with thyroid
3778 peroxidase antibodies and clinical thyroid disease. *PLoS Genet.* **10**, e1004123 (2014).

3779 141. Onengut-Gumuscu, S. *et al.* Fine mapping of type 1 diabetes susceptibility
3780 loci and evidence for colocalization of causal variants with lymphoid gene enhancers.
3781 *Nat. Genet.* **47**, 381–386 (2015).

3782 142. Fischer, A. *et al.* Identification of Immune-Relevant Factors Conferring
3783 Sarcoidosis Genetic Risk. *Am. J. Respir. Crit. Care Med.* **192**, 727–736 (2015).

3784 143. Nikpay, M. *et al.* A comprehensive 1,000 Genomes-based genome-wide
3785 association meta-analysis of coronary artery disease. *Nat. Genet.* **47**, 1121–1130
3786 (2015).

3787 144. Kato, N. *et al.* Trans-ancestry genome-wide association study identifies 12

3788 genetic loci influencing blood pressure and implicates a role for DNA methylation. *Nat.*
3789 *Genet.* **47**, 1282–1293 (2015).

3790 145. de Lange, K. M. *et al.* Genome-wide association study implicates immune
3791 activation of multiple integrin genes in inflammatory bowel disease. *Nat. Genet.* **49**,
3792 256–261 (2017).

3793 146. Langefeld, C. D. *et al.* Transancestral mapping and genetic load in systemic
3794 lupus erythematosus. *Nat. Commun.* **8**, 16021 (2017).

3795 147. Johansson, Å., Rask-Andersen, M., Karlsson, T. & Ek, W. E. Genome-wide
3796 association analysis of 350 000 Caucasians from the UK Biobank identifies novel loci
3797 for asthma, hay fever and eczema. *Hum. Mol. Genet.* **28**, 4022–4041 (2019).

3798 148. Han, Y. *et al.* Genome-wide analysis highlights contribution of immune system
3799 pathways to the genetic architecture of asthma. *Nat. Commun.* **11**, 1776 (2020).

3800 149. Hernandez Cordero, A. I. *et al.* Genome-wide Associations Reveal Human-
3801 Mouse Genetic Convergence and Modifiers of Myogenesis, CPNE1 and STC2. *Am. J.*
3802 *Hum. Genet.* **105**, 1222–1236 (2019).

3803 150. Pulit, S. L. *et al.* Meta-analysis of genome-wide association studies for body
3804 fat distribution in 694 649 individuals of European ancestry. *Hum. Mol. Genet.* **28**, 166–
3805 174 (2019).

3806 151. Nalls, M. A. *et al.* Large-scale meta-analysis of genome-wide association data
3807 identifies six new risk loci for Parkinson's disease. *Nat. Genet.* **46**, 989–993 (2014).

3808 152. Pickrell, J. K. *et al.* Detection and interpretation of shared genetic influences
3809 on 42 human traits. *Nat. Genet.* **48**, 709–717 (2016).

3810 153. Chang, D. *et al.* A meta-analysis of genome-wide association studies
3811 identifies 17 new Parkinson's disease risk loci. *Nat. Genet.* **49**, 1511–1516 (2017).

3812 154. Martin, A. R. *et al.* An Unexpectedly Complex Architecture for Skin
3813 Pigmentation in Africans. *Cell* **171**, 1340–1353.e14 (2017).

3814 155. Jonnalagadda, M. *et al.* A Genome-Wide Association Study of Skin and Iris
3815 Pigmentation among Individuals of South Asian Ancestry. *Genome Biol. Evol.* **11**,

3816 1066–1076 (2019).

3817 156. Sulem, P. *et al.* Sequence variants at CYP1A1-CYP1A2 and AHR associate
3818 with coffee consumption. *Hum. Mol. Genet.* **20**, 2071–2077 (2011).

3819 157. Coffee and Caffeine Genetics Consortium *et al.* Genome-wide meta-analysis
3820 identifies six novel loci associated with habitual coffee consumption. *Mol. Psychiatry* **20**,
3821 647–656 (2015).

3822 158. Cornelis, M. C. *et al.* Genome-wide association study of caffeine metabolites
3823 provides new insights to caffeine metabolism and dietary caffeine-consumption
3824 behavior. *Hum. Mol. Genet.* **25**, 5472–5482 (2016).

3825 159. Haas, M. E. *et al.* Genetic Association of Albuminuria with Cardiometabolic
3826 Disease and Blood Pressure. *Am. J. Hum. Genet.* **103**, 461–473 (2018).

3827 160. Liu, M. *et al.* Association studies of up to 1.2 million individuals yield new
3828 insights into the genetic etiology of tobacco and alcohol use. *Nat. Genet.* **51**, 237–244
3829 (2019).

3830 161. Zanetti, D. *et al.* Identification of 22 novel loci associated with urinary
3831 biomarkers of albumin, sodium, and potassium excretion. *Kidney Int.* **95**, 1197–1208
3832 (2019).

3833 162. Pardiñas, A. F. *et al.* Pharmacogenomic Variants and Drug Interactions
3834 Identified Through the Genetic Analysis of Clozapine Metabolism. *Am. J. Psychiatry*
3835 **176**, 477–486 (2019).

3836 163. Zhong, V. W. *et al.* A genome-wide association study of bitter and sweet
3837 beverage consumption. *Hum. Mol. Genet.* **28**, 2449–2457 (2019).

3838 164. Pazoki, R. *et al.* GWAS for urinary sodium and potassium excretion highlights
3839 pathways shared with cardiovascular traits. *Nat. Commun.* **10**, 3653 (2019).

3840 165. Karlsson, T. *et al.* Contribution of genetics to visceral adiposity and its relation
3841 to cardiovascular and metabolic disease. *Nat. Med.* **25**, 1390–1395 (2019).

3842 166. Cornelis, M. C. *et al.* Genome-wide meta-analysis identifies regions on 7p21
3843 (AHR) and 15q24 (CYP1A2) as determinants of habitual caffeine consumption. *PLoS*

3844 *Genet.* **7**, e1002033 (2011).

3845 167. Amin, N. *et al.* Genome-wide association analysis of coffee drinking suggests
3846 association with CYP1A1/CYP1A2 and NRCAM. *Mol. Psychiatry* **17**, 1116–1129 (2012).

3847 168. Teumer, A. *et al.* Genome-wide association meta-analyses and fine-mapping
3848 elucidate pathways influencing albuminuria. *Nat. Commun.* **10**, 4130 (2019).

3849 169. Wu, Y. *et al.* Genome-wide association study of medication-use and
3850 associated disease in the UK Biobank. *Nat. Commun.* **10**, 1891 (2019).

3851 170. Eriksson, N. *et al.* Web-based, participant-driven studies yield novel genetic
3852 associations for common traits. *PLoS Genet.* **6**, e1000993 (2010).

3853 171. Zhao, B. *et al.* Large-scale GWAS reveals genetic architecture of brain white
3854 matter microstructure and genetic overlap with cognitive and mental health traits (n =
3855 17,706). *Mol. Psychiatry* (2019) doi:10.1038/s41380-019-0569-z.

3856 172. Voight, B. F., Kudravalli, S., Wen, X. & Pritchard, J. K. A map of recent
3857 positive selection in the human genome. *PLoS Biol.* **4**, e72 (2006).

3858 173. Sabeti, P. C. *et al.* Genome-wide detection and characterization of positive
3859 selection in human populations. *Nature* **449**, 913–918 (2007).

3860 174. Pickrell, J. K. *et al.* Signals of recent positive selection in a worldwide sample
3861 of human populations. *Genome Res.* **19**, 826–837 (2009).

3862 175. Buckley, M. T. *et al.* Selection in Europeans on Fatty Acid Desaturases
3863 Associated with Dietary Changes. *Mol. Biol. Evol.* **34**, 1307–1318 (2017).

3864 176. Ye, K., Gao, F., Wang, D., Bar-Yosef, O. & Keinan, A. Dietary adaptation of
3865 FADS genes in Europe varied across time and geography. *Nat Ecol Evol* **1**, 167 (2017).

3866 177. Mathieson, S. & Mathieson, I. FADS1 and the Timing of Human Adaptation to
3867 Agriculture. *Mol. Biol. Evol.* **35**, 2957–2970 (2018).

3868 178. Ju, D. & Mathieson, I. The evolution of skin pigmentation-associated variation
3869 in West Eurasia. *Proc. Natl. Acad. Sci. U. S. A.* **118**, (2021).

3870 179. Stefansson, H. *et al.* A common inversion under selection in Europeans. *Nat.*
3871 *Genet.* **37**, 129–137 (2005).

- 3872 180. Boettger, L. M., Handsaker, R. E., Zody, M. C. & McCarroll, S. A. Structural
3873 haplotypes and recent evolution of the human 17q21.31 region. *Nat. Genet.* **44**, 881–
3874 885 (2012).
- 3875 181. Itan, Y., Powell, A., Beaumont, M. A., Burger, J. & Thomas, M. G. The origins
3876 of lactase persistence in Europe. *PLoS Comput. Biol.* **5**, e1000491 (2009).
- 3877 182. Antelope, C. X., Marnetto, D., Casey, F. & Huerta-Sanchez, E. Leveraging
3878 Multiple Populations across Time Helps Define Accurate Models of Human Evolution: A
3879 Reanalysis of the Lactase Persistence Adaptation. *Hum. Biol.* **89**, 81–97 (2017).
- 3880 183. Segurel, L. *et al.* Why and when was lactase persistence selected for?
3881 Insights from Central Asian herders and ancient DNA. *PLoS Biol.* **18**, e3000742 (2020).
- 3882 184. Burger, J. *et al.* Low Prevalence of Lactase Persistence in Bronze Age
3883 Europe Indicates Ongoing Strong Selection over the Last 3,000 Years. *Curr. Biol.* **30**,
3884 4307–4315.e13 (2020).
- 3885

4b) Detangling Direct and Indirect impacts of sample age from the Mesolithic-Neolithic data on genotype imputation

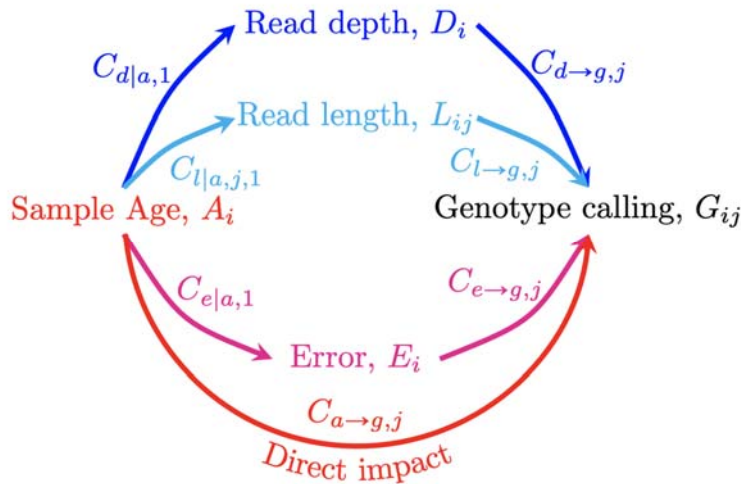
Rasmus Amund Henriksen¹, Rasmus Nielsen^{1,2}, Lei Zhao¹, Thorfinn Sand Korneliusen¹

¹Lundbeck Foundation GeoGenetics Centre, GLOBE Institute, University of Copenhagen, Copenhagen, Denmark

²Department of Integrative Biology, University of California, Berkeley

Many factors can influence the genotype imputation/calling at specific sites, especially when present-day reference panels are used to impute missing genotypes for ancient samples. The unique features of aDNA relative to modern DNA include 1) relatively shallower read depths ¹, 2) shorter read lengths ², and 3) different error profiles ³. Mapping biases that depend on read length and error rates are of particular concern for aDNA and may cause spurious signals of selection. Naturally, read depth may also affect genotype calling. In order to filter SNPs that might be affected by such biases, and to generally correct for and quantify the biases, we develop a causal inference method for distinguishing direct effects of sample age on allele frequency from other (indirect) effects mediated by age-dependent errors, read depth, and read length.

As shown in Figure 4b.1, to distinguish the true selection signal of age on allele frequency, from the signal caused by mapping biases, and other biases, in the ancient samples, we create a model and a workflow which can decompose the influence of sample age (A_i , age of individual i) on genotype (G_{ij} , Genotype of individual i at position j), into its indirect and direct effects. The imputed genotype for each individual is converted into allele frequencies, representing the homozygous for the reference allele, heterozygous and homozygous for the alternative allele as 0, 0.5 and 1 respectively. The indirect effects are mediated through the three unique features of aDNA as previously described, while the direct impact reflects the true change in allele frequency over the time range from the ancient samples to the modern ones. The three factors are the mean read depth across all sites per individual, the mean read length for each individual at each site, and the third is an overall error estimate for each individual. These three factors are all identified using ANGSD ⁴ from the aligned BAM files for all individuals used in the study. The depth and length are calculated based solely on the imputed positions (j) extracted from the imputed vcf files, whereas the overall error is calculated as described in Orlando et al. 2013 from the entire BAM file containing all the invariant sites.



3925

3926

3927

3928

3929

3930

3931

3932

3933

3934

3935

3936

3937

3938

3939

3940

3941

3942

3943

3944

3945

3946

3947

3948

3949

3950

Figure S4b.1. Illustrative figure of the factors that influence genotype calling procedure. The denotations with symbol “|” can be obtained by regression, while the denotations with symbol “→” are derived after PCA.

The initial step is to conduct three linear regressions, one for each of the three factors, mean depth (equation 1), mean length (equation 2) and error estimate (equation 3) with respect to one explanatory variable, i.e., sample age (i.e., A_i). This step is to investigate the influences of sample age on each of the three factors.

$$D_i \triangleq \overline{D_{ij}} \sim C_{d|a,1} A_i + C_{d|a,0} + \epsilon_{d|a,i} \quad (1)$$

$$L_{ij} \triangleq \overline{L_{ijk}} \sim C_{l|a,j,1} A_i + C_{l|a,j,0} + \epsilon_{l|a,ij} \quad (2)$$

$$E_i \sim C_{e|a,1} A_i + C_{e|a,0} + \epsilon_{e|a,i}. \quad (3)$$

Where \triangleq means definition, $\overline{D_{ij}}$ means the average depth with regard to j for a specific individual i , and $\overline{L_{ijk}}$ means the average read length, for all reads k stretching over site j of individual i . E_i means the overall error rate of individual i .

We assume that the sample age influences these factors, as such it is necessary to eliminate this influence as done in equation 4-6.

$$\Delta_i \triangleq D_i - C_{d|a,1} A_i \quad (4)$$

$$\Lambda_{ij} \triangleq L_{ij} - C_{l|a,j,1} A_i \quad (5)$$

$$\Sigma_i \triangleq E_i - C_{e|a,1} A_i \quad (6)$$

A second round of regression is necessary to detect the direct impact of the three factors as well as the direct impact of age on the genotype. The explanatory variables can be the

3951 sample age (A_i) and the three remainders (or the linear combinations of the three
 3952 remainders). To avoid potential correlations between the three remainders, we perform a
 3953 Principal Component Analysis. The relationship between the PCA scores and the remainders
 3954 can be represented as in equation 7 and 8.

3955

$$\left(\frac{\Lambda_{ij} - \Lambda_j}{sd(\Lambda_{.j})}, \frac{\Sigma_i - \Sigma}{sd(\Sigma_{.})}, \frac{\Delta_i - \Delta}{sd(\Delta_{.})} \right)_{n \times 3} = (\alpha_{ij} \beta_{ij}, \gamma_{ij})_{n \times 3} \begin{pmatrix} \vec{\omega}'_{j,1} \\ \vec{\omega}'_{j,2} \\ \vec{\omega}'_{j,3} \end{pmatrix} \quad (7)$$

3956

3957

$$(\alpha_{ij} \beta_{ij}, \gamma_{ij})_{n \times 3} = \left(\frac{\Lambda_{ij} - \Lambda_j}{sd(\Lambda_{.j})}, \frac{\Sigma_i - \Sigma}{sd(\Sigma_{.})}, \frac{\Delta_i - \Delta}{sd(\Delta_{.})} \right)_{n \times 3} (\vec{\omega}_{j,1}, \vec{\omega}_{j,2}, \vec{\omega}_{j,3}) \quad (8)$$

3958

3959 Where Λ_j and $sd(\Lambda_{.j})$ are the mean and the standard deviation of Λ across individuals i at
 3960 fixed position j . Σ and $sd(\Sigma_{.})$ are the mean and the standard deviation of Σ across individuals
 3961 i . Δ and $sd(\Delta_{.})$ are the mean and the standard deviation of Δ across individuals i . $\vec{\omega}_{j,2}, \vec{\omega}_{j,1},$
 3962 $\vec{\omega}_{j,3}$ are the three principle directions (column eigenvectors) and the $\vec{\omega}'_{j,2}, \vec{\omega}'_{j,1}, \vec{\omega}'_{j,3}$ are the
 3963 corresponding transposed vectors. $\alpha_{ij}, \beta_{ij}, \gamma_{ij}$ are the principal component scores of

3964 $\left(\frac{\Lambda_{ij} - \Lambda_j}{sd(\Lambda_{.j})}, \frac{\Sigma_i - \Sigma}{sd(\Sigma_{.})}, \frac{\Delta_i - \Delta}{sd(\Delta_{.})} \right)$, with n being the number of individuals carrying site j .

3965

3966 Once all the explanatory variables are independent, the coefficient $C_{g|a,j,1}$ represent
 3967 the total impact of age which can be obtained by conducting a final round of regression
 3968 (equation 9).

$$G_{ij} \sim C_{g|\alpha,j,1} \alpha_{ij} + C_{g|\beta,j,1} \beta_{ij} + C_{g|\gamma,j,1} \gamma_{ij} + C_{g|a,j,1} A_i + C_{0,j} + \epsilon_{ij} \quad (9)$$

3969

3970

3971 The sum of the first three terms of equation 9 can be obtained by multiplying the
 3972 coefficients $(C_{g|\alpha,j,1}, C_{g|\beta,j,1}, C_{g|\gamma,j,1})'$ with equation 8 (shown in equation 10). And the effective
 3973 regression coefficients for $(\Delta_{ij}, \Lambda_i, \Sigma_i)$ will be observed when calculating the linear slopes of
 3974 the remainders in Equation 10. Such effective coefficients can be viewed as measurements
 3975 of the direct impact of the corresponding factors, i.e. length ($C_{l \rightarrow g,j}$), error ($C_{e \rightarrow g,j}$) and depth
 3976 ($C_{d \rightarrow g,j}$), on the genotype calling (equation 11 - 13).

3977

$$(\alpha_{ij} \beta_{ij}, \gamma_{ij})_{n \times 3} \begin{pmatrix} C_{g|\alpha,j,1} \\ C_{g|\beta,j,1} \\ C_{g|\gamma,j,1} \end{pmatrix} = \left(\frac{\Lambda_{ij} - \Lambda_j}{sd(\Lambda_{.j})}, \frac{\Sigma_i - \Sigma}{sd(\Sigma_{.})}, \frac{\Delta_i - \Delta}{sd(\Delta_{.})} \right)_{n \times 3} (\vec{\omega}_{j,1}, \vec{\omega}_{j,2}, \vec{\omega}_{j,3})_{3 \times 3} \begin{pmatrix} C_{g|\alpha,j,1} \\ C_{g|\beta,j,1} \\ C_{g|\gamma,j,1} \end{pmatrix} \quad (10)$$

3978

3979

$$\vec{\omega}_{j,1} = \begin{pmatrix} \omega_{j,1,1} \\ \omega_{j,2,1} \\ \omega_{j,3,1} \end{pmatrix}, \vec{\omega}_{j,2} = \begin{pmatrix} \omega_{j,1,2} \\ \omega_{j,2,2} \\ \omega_{j,3,2} \end{pmatrix}, \vec{\omega}_{j,3} = \begin{pmatrix} \omega_{j,1,3} \\ \omega_{j,2,3} \\ \omega_{j,3,3} \end{pmatrix}$$

3980

3981

3982

$$C_{l \rightarrow g,j} = \frac{1}{\text{sd}(\Lambda_j)} (\omega_{j,1,1} C_{g|\alpha,j,1} + \omega_{j,1,2} C_{g|\beta,j,1} + \omega_{j,1,3} C_{g|\gamma,j,1}) \quad (11)$$

3983

$$C_{e \rightarrow g,j} = \frac{1}{\text{sd}(\Sigma_j)} (\omega_{j,2,1} C_{g|\alpha,j,1} + \omega_{j,2,2} C_{g|\beta,j,1} + \omega_{j,2,3} C_{g|\gamma,j,1}) \quad (12)$$

3984

$$C_{d \rightarrow g,j} = \frac{1}{\text{sd}(\Delta_j)} (\omega_{j,3,1} C_{g|\alpha,j,1} + \omega_{j,3,2} C_{g|\beta,j,1} + \omega_{j,3,3} C_{g|\gamma,j,1}) \quad (13)$$

3985

3986

3987

3988

3989

3990

3991

Any influence of age on the genotype calling imposed through one of the three factors (depth, length, error) is an indirect impact. The measurements of such indirect impacts are calculated by multiplying equations 11-13 with each of the factors corresponding coefficients obtained from our first round of regression (equation 4-6) for each site j .

3992

$$C_{a \rightarrow l \rightarrow g,j} = C_{l|a,j,1} \cdot C_{l \rightarrow g,j} \quad (14)$$

3993

$$C_{a \rightarrow e \rightarrow g,j} = C_{e|a,1} \cdot C_{e \rightarrow g,j} \quad (15)$$

3994

$$C_{a \rightarrow d \rightarrow g,j} = C_{d|a,1} \cdot C_{d \rightarrow g,j} \quad (16)$$

3995

3996

3997

3998

3999

Finally, the direct impact of age on the genotype calling at site j (equation 17) can be obtained by subtracting all indirect impacts of age (equation 14-16) from the total impact of age $C_{g|a,j,1}$, obtained from the second round of regression (equation 9).

4000

$$C_{a \rightarrow g,j} = C_{g|a,j,1} - C_{a \rightarrow l \rightarrow g,j} - C_{a \rightarrow e \rightarrow g,j} - C_{a \rightarrow d \rightarrow g,j} \quad (17)$$

4001

4002

4003

4004

4005

4006

4007

4008

4009

4010

$$R_j = \frac{C_{a \rightarrow l \rightarrow g,j} + C_{a \rightarrow e \rightarrow g,j} + C_{a \rightarrow d \rightarrow g,j}}{C_{a \rightarrow g,j}} \quad (18)$$

4011

4012

4013 R_j can also be converted into a fraction representing the proportion of indirect effects on age

4014 relative to the total effect (equation 19), which can equivalently be used for filtering:

4015

$$F_j = \frac{R_j}{1 + R_j} = \frac{1}{1 + \frac{1}{R_j}} \quad (19)$$

4016

4017

4018 To filter out sites in selection analyses that may be affected by biases, we use a fixed

4019 threshold of $0.5 < F_j \leq 1$.

4020

4021 References

4022 1. Shapiro, B. & Hofreiter, M. A paleogenomic perspective on evolution and gene

4023 function: new insights from ancient DNA. *Science* **343**, 1236573 (2014).

4024 2. Paabo, S. Ancient DNA: extraction, characterization, molecular cloning, and

4025 enzymatic amplification. *Proceedings of the National Academy of Sciences* vol. 86

4026 1939–1943 (1989).

4027 3. Orlando, L. *et al.* Recalibrating Equus evolution using the genome sequence of an

4028 early Middle Pleistocene horse. *Nature* **499**, 74–78 (2013).

4029 4. Korneliussen, T. S., Albrechtsen, A. & Nielsen, R. ANGSD: Analysis of Next

4030 Generation Sequencing Data. *BMC Bioinformatics* **15**, 356 (2014).

4031

4032

4033 4c) Over-dispersion in polygenic scores across ancient 4034 populations

4035

Alba Refoyo Martínez¹, Fernando Racimo¹

4036

4037 ¹Lundbeck Foundation GeoGenetics Centre, GLOBE Institute, University of Copenhagen,

4038

Copenhagen, Denmark

4039 Introduction

4040 We aimed to test whether there was evidence for over-dispersion in polygenic scores across
4041 different ancient populations. This could indicate if there was evidence for strong
4042 differentiation in the genetics of traits across ancient populations, beyond what can be
4043 explained by genetic drift alone.

4044 Methods

4045 We used summary statistics from GWASs performed on the UK Biobank cohort ¹ by the
4046 Neale lab (Round 2: <http://www.nealelab.is/uk-biobank/>), which include both quantitative and
4047 case-control traits. From the case-control case traits, we filtered out those where the N-
4048 case/N ratio is lower than 1/100 and those that did not have any associated variants at the
4049 standard genome-wide significance threshold ($P = 5e-8$). We then filtered out variants with
4050 minor allele frequency lower than 5%, variants with INFO score $< 50\%$, variants with a
4051 genotype probability lower than 0.8 in more than 10% of the individuals, triallelic variants and
4052 variants flagged as “low-confident” after imputation (Supplementary Note S2).

4053 We divided the genome into 1,703 non-overlapping and approximately independent linkage
4054 disequilibrium (LD) blocks ². From each block, we retrieved the variant with the lowest
4055 association P-value. We then selected only the variants with a P-value lower than the
4056 genome-wide significance threshold, $5e-8$, for downstream analyses.

4057 Polygenic scores were calculated by summing over the allele frequency of the filtered trait-
4058 associated variants and weighting them by the effect size obtained from the summary
4059 statistics from the UK Biobank (round 2). The allele frequency was retrieved from the
4060 imputed ancient West-Eurasian individuals. Individuals were clustered into groups of closely-
4061 related individuals using MDS (see Supplementary Note S3). In our overdispersion test, we
4062 did not include genomes that were not successfully classified as belonging to any of the
4063 clusters, and removed clusters that contained less than four individuals. This resulted in
4064 1,119 individuals clustered into 41 groups.

4065 We evaluated our choice of the summary statistics by comparing the Neale Lab scores to
4066 scores obtained using a different GWAS on the same phenotype and on the same cohort -
4067 the GWAS ATLAS ^{3,4} (<https://atlas.ctglab.nl/>). Despite both association tests being performed
4068 on the same cohort, the two studies applied different filters on the data, which resulted in
4069 differences in the effect size estimates and in the significance of the SNP associations
4070 (Figure **S4c.1**, panel C), which in turn lead to discrepancies on the polygenic scores,
4071 particularly for pre-Neolithic populations, like the Western hunter-gatherers, and for Eastern

4072 Eurasian individuals, like the ancient Siberians, which are more distantly related to the
4073 present-day British individuals included in the GWAS cohorts (Figure **S4c.1**, panel B). We
4074 followed the Sohail et al. ⁵ procedure to detect residual population stratification along the
4075 axes of population variation by looking for strong correlations between axes of population
4076 structure and the magnitude of effect size estimates. We found that the GWAS ATLAS had
4077 substantially more uncorrected population stratification than the Neale Lab GWAS (Figure
4078 **S4c.1**, panel A). Therefore, all results presented below were based on the scores obtained
4079 using effect sizes from the Neale Lab (round 2) estimates. We also observed that polygenic
4080 scores for pre-Neolithic individuals were particularly sensitive to the choice of cohort and to
4081 the SNP filtering scheme, so we urge caution in the interpretation of those values.

4082 The Q_x statistic was introduced by ⁶ to look for overdispersion in polygenic scores across
4083 populations that cannot be explained simply by genetic drift, assuming the polygenic scores
4084 were not biased by population stratification in the GWAS cohort from which effect size
4085 estimates were obtained. Polygenic scores are also assumed to follow a multivariate normal
4086 distribution under a null model of genetic drift, and this statistic serves to look for departures
4087 from this model. To compute the Q_x statistic, an empirical genome-wide covariance matrix is
4088 needed, which we constructed using a subset of SNPs with a trait-association p-value larger
4089 than $5e-8$, and then we sampled every 20th “non-associated” SNPs across the genome.

4090 To further test the significance of the overdispersion and account for possible deviations from
4091 the assumptions the Q_x statistic makes, we also computed P-values using two randomization
4092 schemes simulating a neutral scenario. The first method was based on the randomization of
4093 the signs of the effect size estimates of trait-associated SNPs while the other was based on
4094 sampling variants across the genome with frequencies matching those of trait-associated
4095 variants in the GBR panel from the 1000 Genomes Project.

4096

4097 To address mapping biases we performed a one-tailed wilcoxon rank-sum test for each trait.
4098 We evaluated if the candidate associated SNPs had higher values of the artefactual effect
4099 estimates than the non-associated SNPs used as our neutral baseline (Supplementary Note
4100 S4b). None of the tests were significant (min-P value = 0.19; max-P Value = 0.99;
4101 Mean=0.56).

4102 Results

4103 Polygenic scores across ancient populations

4104 We computed polygenic scores of trait-associated SNPs across the 41 ancient population
4105 groups. We used the effect size estimates from the UK Biobank Neale lab GWAS ¹ and the
4106 allele frequencies from the 41 population clusters previously described (Supplementary Note
4107 S3). We filtered out those traits that had less than 10 genome-wide significantly associated
4108 SNPs, restricting our analysis to a total of 320 polygenic traits. We then used the Q_x statistic
4109 on these traits to test for overdispersion across clusters.

4110 We applied Q_x statistic to each of the 320 traits. Figure **S4c.2** shows p-values for Q_x statistic
4111 for the standard genome-wide set of trait-associated SNPs ($P < 5e-8$). From these, 119
4112 resulted in a nominally significant Q_x statistic and only 39 remained significantly over-
4113 dispersed after controlling for multiple testing via a Bonferroni correction (Figure **S4c.3**). We
4114 grouped these 39 traits into ten broad categories in agreement with different sub-categories
4115 levels from the Data Showcase (<https://biobank.ctsu.ox.ac.uk/showcase/>): “body
4116 measurements”, “impedance measures”, “spirometry”, “sun exposure: hair and skin
4117 pigmentation”, “assay results”, “diabetes, cholesterol or blood pressure”, “diet”, “medical
4118 conditions”, “mental health” and “sex-specific factors” (Figure **S4c.4**).

4119 We also computed p-values using two randomization schemes to account for possible
4120 violations of the assumption that the Q_x statistic is chi-squared distributed: one based on
4121 randomising of the effect sizes before calculating polygenic scores and the other based on
4122 frequency-matched variants that were not significantly associated with the trait under study
4123 ^{6,7}.

4124 Individual scores across time and space

4125 We also computed polygenic scores individually for each ancient genome. Figure **S4c.5**
4126 shows the ancient genomes projected into the first two principal components, colouring each
4127 genome with its corresponding polygenic score for height. Maps for each of the clusters are
4128 also available to better interpret how scores change across space and time (Figures **S4c.6 -**
4129 **S4c.8**).

4130 Differentiation among ancient populations and GBR

4131 Finally, we were interested in determining how differentiated our ancient populations were to
4132 each other, and to a present-day British panel: the 1000 Genomes GBR panel ⁸. This would

4133 enable us to better understand how portable the polygenic scores created using the UK
4134 Biobank were to each ancient population, by drawing analogies with pairwise comparisons
4135 with known estimates of score portability between present-day populations^{9,10}. In Figure
4136 **S4c.9.B**, we show genome-wide pairwise F_{st} estimates¹¹ computed between GBR and all
4137 the other present-day population panels from the 1000 Genomes Project using vcfTools¹². In
4138 turn, in Figure **S4c.9.A**, we show genome-wide pairwise F_{st} estimates computed between
4139 each of our ancient population clusters, as well as between each ancient population cluster
4140 and GBR. GBR shows lower F_{st} values when compared to IBD Global Cluster
4141 “Eurasia_5000BP_200BP” which represents modern Eurasian diversity after the major
4142 Bronze Age migrations events. Within the “Eurasia_5000BP_200BP” cluster, GBR is more
4143 genetically distant to the Steppe populations. GBR shows low values in the order of 0.01
4144 when compared to “EuropeWCAAsia_25000BP_300BP”. The
4145 “EuropeWCAAsia_25000BP_300BP” cluster shows west-asian populations that are thought to
4146 have brought Neolithic farming practises into Europe, as well as individuals with ancestry
4147 associated with Caucasus hunter-gatherers. This cluster shows strong differentiation when
4148 compared to “Europe_15000_4000BP”. The cluster “Europe_15000_4000BP” (which
4149 includes both eastern and western hunter-gatherer populations that share more ancestry
4150 with Siberian samples) clustered together with “Asia_45000BP_2000BP”. The
4151 “Asia_45000BP_2000BP” cluster represents Siberian hunter-gatherers sampled during the
4152 Neolithic and Bronze Age.

4153 Discussion

4154 After testing for genetic polygenic score differentiation among ancient populations, we found
4155 39 traits with significant overdispersion in scores after controlling for multiple testing (Table
4156 **S4c.10**). Most significant traits are related to pigmentation, body size differences, disorders
4157 related to diet and sugar levels and mental health conditions. A few of these differences
4158 among ancient populations have been reported before^{13–15} but our analysis provides a
4159 much more fine-scale account of these differences across time and space, due to the
4160 relatively higher availability of population genomic data, and across different phenotypes.

4161 The top ten most significantly overdispersed traits are related with either hair or skin
4162 pigmentation. These phenotypes (Data-field: 1737, 2267, 1717, 1727, C44, C_SKIN and
4163 C3_SKIN) are very closely linked to sun exposure, from which the last three represent
4164 associations with malignant neoplasms of skin. Polygenic scores for these three traits tend to
4165 be higher in Bronze Age Europeans and Neolithic Farmers compared to northeastern asian
4166 and hunter-gathers groups. Skin colour polygenic scores are higher in hunter-gatherers
4167 groups (Panel C of Figure **S4c.5**). Similarly, Europeans tend to have higher polygenic scores

4168 for blonde and light brown hair colour, while East Asian have higher scores for dark brown
4169 and black hair (Panel A and B of Figure **S4c.5**).

4170 Intriguingly, we find that Eastern hunter-gatherers had much higher polygenic scores for
4171 height than Western hunter-gatherers, indicating that even in Mesolithic Europe, before the
4172 arrival of Neolithic farmers, there was already strong genetic differentiation at variants
4173 associated with height between local populations. We find that the Yamnaya steppe people
4174 and the Caucasus hunter-gatherers / Iranian Neolithic people generally have high genetic
4175 scores for height, while Levant Neolithic peoples, early European farmers and western
4176 hunter-gatherers have low genetic height scores (Figure **S4c.6**).

4177 Lung capacity is very influenced by height, BMI and ethnicity. One of the most common
4178 clinical practices to measure lung capacity is forced expiratory volume in the first second
4179 (FEV-1), which happens to be one of the significant traits (Data-field: 20153) ¹⁶ (Figure
4180 **S4c.4**).

4181 Impedance measures are used for estimating body composition, especially body fat and
4182 muscle mass ¹⁷. The difference in height between the samples could also partially explain
4183 the distribution of polygenic scores in the impedance traits category. For instance, taller
4184 populations require a higher energy intake which is proportional to their basal metabolic rate
4185 (BMR) ¹⁸. As we can see in Panel D of Figure **S4c.5**, farmers, who have a very similar
4186 distribution of height polygenic scores as Eastern hunter-gatherers, have lower BMR
4187 polygenic scores and also a more sedentary life ¹⁹. The body mass index (BMI) is also
4188 strongly influenced by the percentage of body fat and muscle mass. Therefore, ethnic
4189 differences in BMI-associated health risks may be caused by differences in impedance
4190 measures ¹⁷. Height could be a confounding factor when interpreting the significance of
4191 some closely related traits.

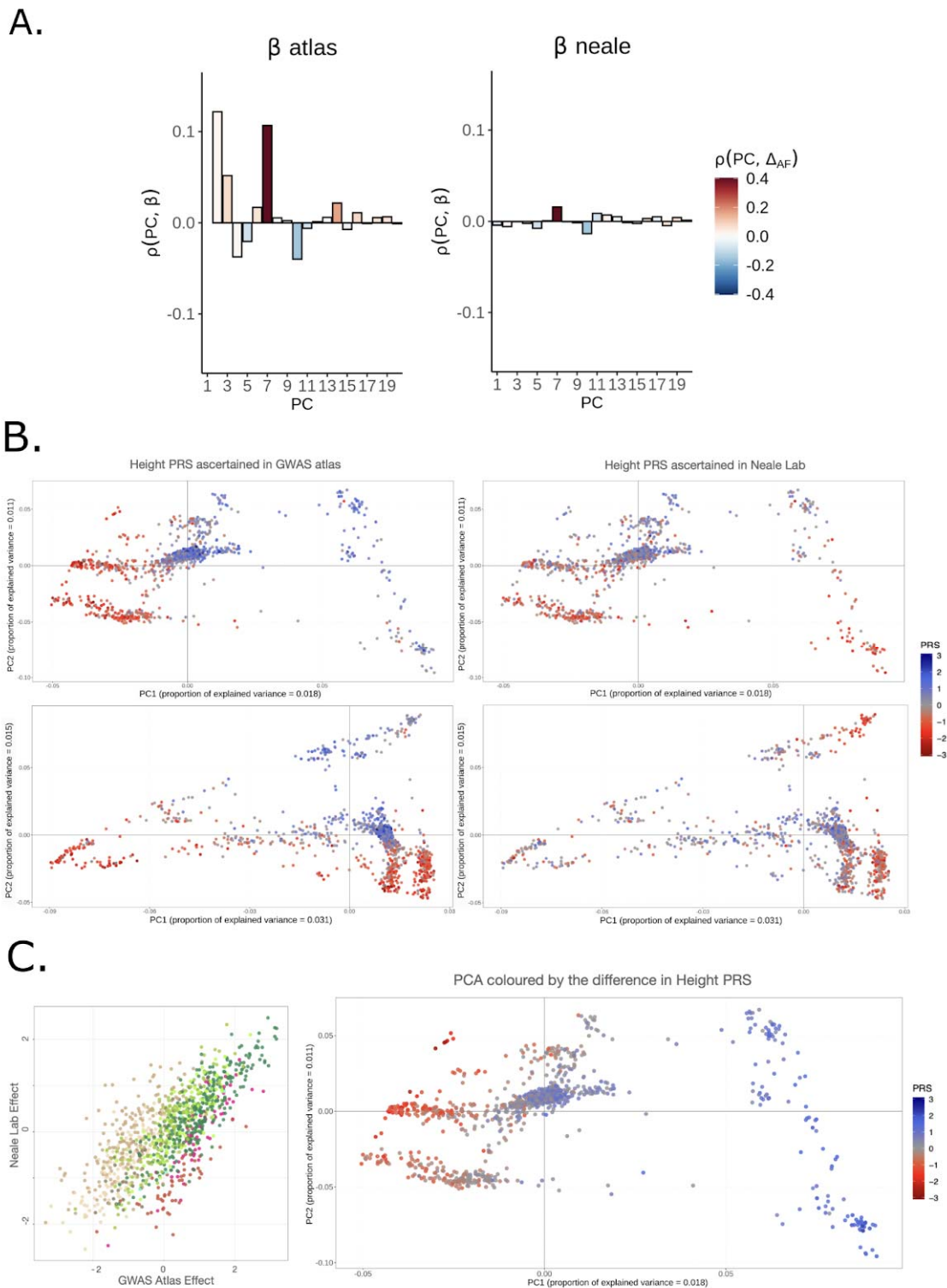
4192 There are a few traits directly or indirectly connected with human diets and local environment
4193 adaptation. Some of these traits are linked to sugar levels, cholesterol and blood pressure,
4194 which can increase the risk of developing a cardiovascular disease or diabetes. For instance,
4195 low levels of potassium in urine is a risk factor for cardiovascular disease and diabetes ²⁰⁻²².
4196 Abnormal low levels of mean corpuscular haemoglobin concentration have been associated
4197 with local adaptations to high altitude environments ²³.

4198 Three of the significant traits (UK Biobank codes "1940", "2030" and "1980") correspond to
4199 human feelings and mood instability which can have an effect on mental health. Both worrier

4200 and nervous feelings are associated with anxiety disorders while irritability is one of the
4201 critical systems to clinically assess depression and/or bipolar affective disorders ^{24,25}.

4202 It is important to keep in mind that all these polygenic scores rely on effect size estimates
4203 obtained from the UK Biobank and, in particular, those individuals identifying as “white
4204 British” within that panel. This might lead to biases in the portability of these scores to
4205 distantly related ancient populations. For example, ancient Western Eurasian hunter-
4206 gatherers are as differentiated from present-day British individuals as present-day British
4207 individuals are to some of the South East Asian populations (ITU, STU, BEB, GIH, PJJ) and
4208 American (MXL). This means that the expected portability of these scores in Western hunter-
4209 gatherers should be comparable to that observed in South Asians when computing polygenic
4210 scores using European effect size estimates ^{9,10}.

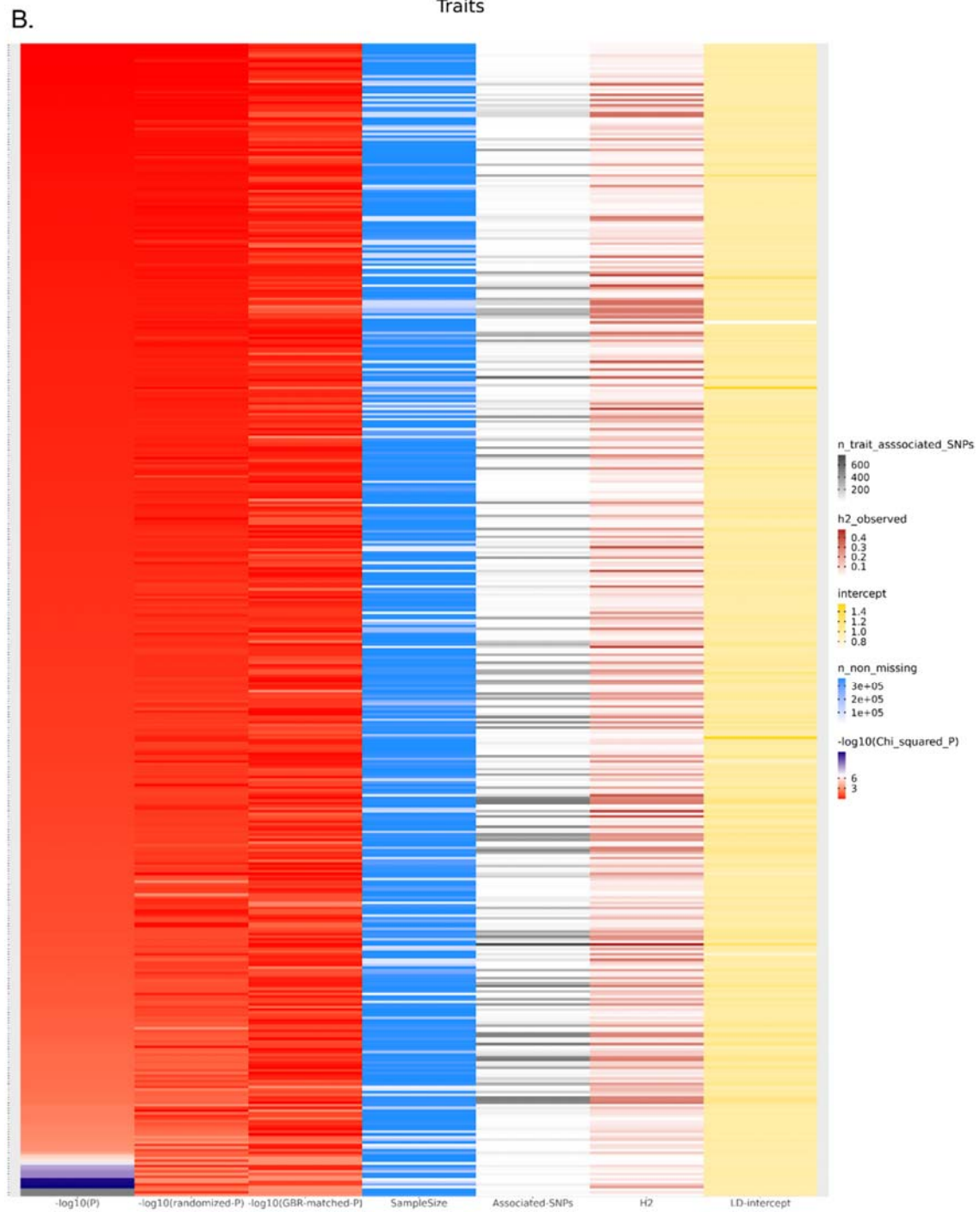
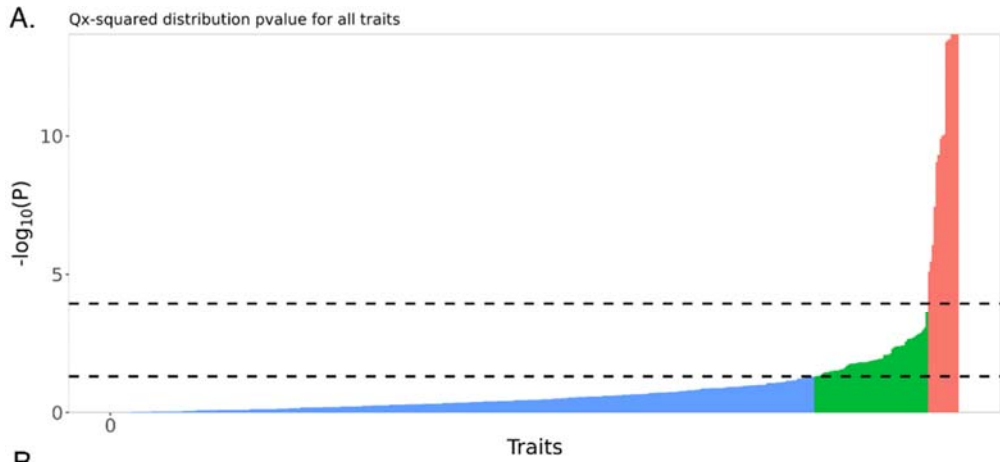
4211 We also note that the value of polygenic scores and the magnitude of the Q_x statistic may
4212 both be affected by population stratification in the GWAS panel from which effect size
4213 estimates were obtained. This seems to be less of a problem in UK Biobank GWAS than in
4214 other GWAS based on meta-analyses ^{5,26,27} but we should nevertheless be cautious about
4215 conclusions drawn purely from these effect size estimates.



4217
4218
4219
4220
4221

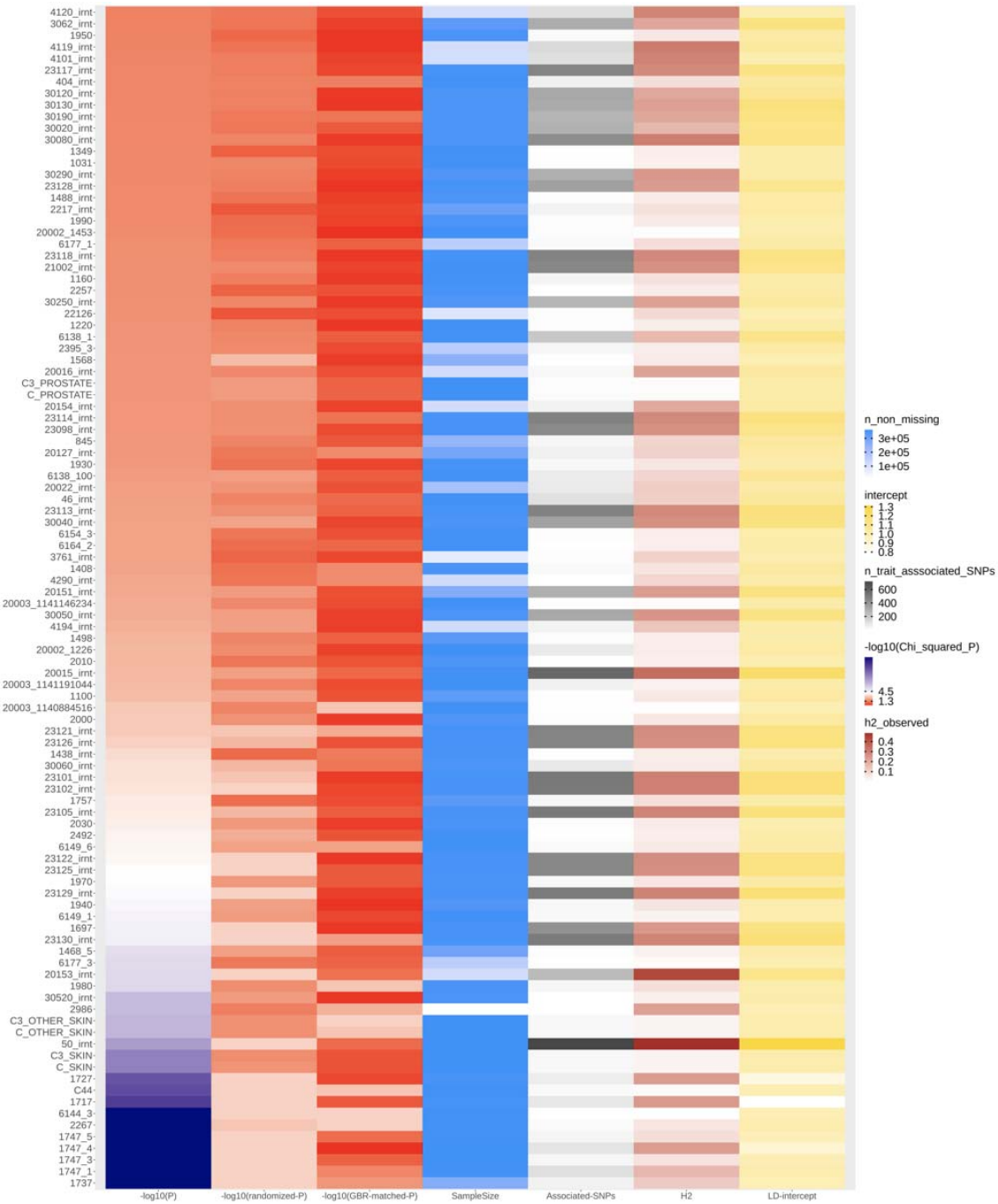
Figure S4c.1. A. Pearson correlations between 20 PC loadings and height effect size estimates from the GWAS atlas, compared to the same correlation using effect size estimates from the Neale Lab GWAS, both summary statistics performed on the UKBiobank. The correlations were computed using SNPs that are present in both the GWAS atlas and Neale

4222 Lab GWAS summary statistics, and in the 1000 Genomes Project. The barplots are coloured
4223 by the correlation between each loading and the allele frequency difference between GBR and
4224 TSI. **B.** Height PRS scores ascertained in GWAS atlas (left panel) and in Neale Lab (right
4225 panel). **C.** Right panel: SNP-associated effect size in GWAS atlas against their effect size in
4226 the Neale Lab. Left panel: Height PRS scores are coloured by the difference between the two
4227 ascertainment (GWAS atlas scores versus Neale Lab scores.)
4228
4229
4230



4232
 4233
 4234
 4235
 4236
 4237
 4238
 4239

Figure S4c.2. Qx p-value for 320 traits. A. Trait-associated SNPs are selected using the 5e-8 cutoff. 119 traits are significant (p -value < 0.05) and only 39 are significant after Bonferroni correction (p -value $< 0.05/320$). B. Heatmap includes the log values for QX p-value, randomised p-value and GBR-matched p-value, sample size, number of associated SNPs used to compute Qx statistic, heritability coefficient and LD-intercept value.



4240
 4241
 4242
 4243
 4244
 4245

Figure S4c.3. Qx p-value for 119 significant traits. Traits are labelled with the corresponding Uk Biobank data coding system. Heatmap includes the log values for QX p-value, randomised p-value and GBR-matched p-value, sample size, number of associated SNPs used to compute Qx statistic, heritability coefficient and LD-intercept value.

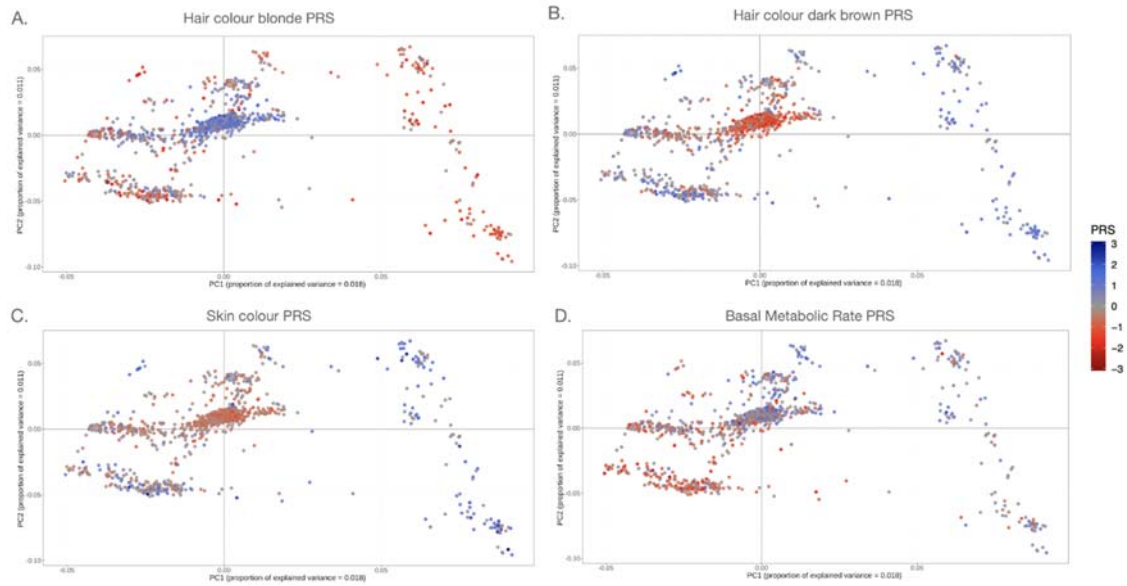
4246
4247
4248
4249
4250
4251
4252



4253

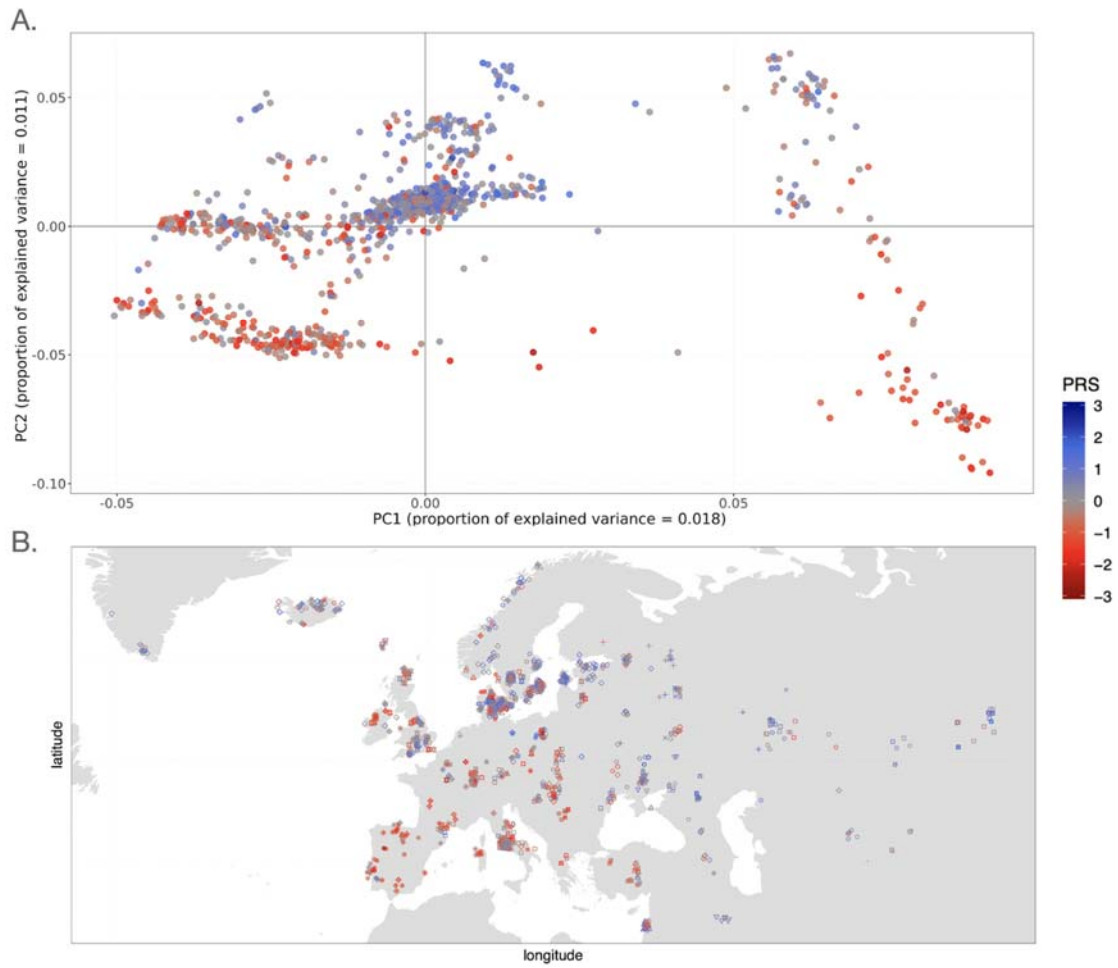
4254
4255
4256
4257
4258
4259
4260

Figure S4c.4. Polygenic scores for 39 significant traits after Bonferroni correction grouped by traits category (trait-associated SNPs with a p-value < 5e-08). Polygenic scores for each of the different populations are shown. They are coloured by broader population groups.



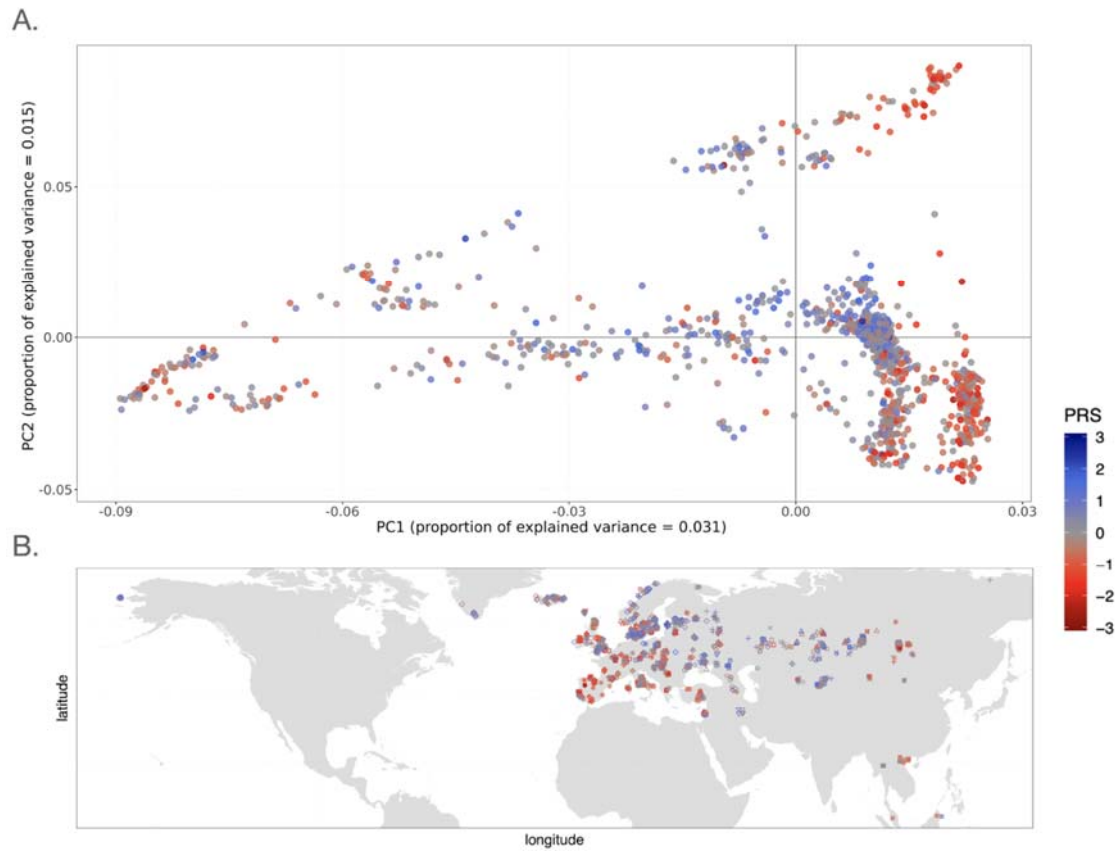
4261
4262
4263
4264
4265
4266

Figure S4c.5. Principal component analysis on West Eurasian samples coloured by individual polygenic scores. A-C. Polygenic scores of sun exposure traits. D. Polygenic scores of Impedance measures trait.



4267
4268
4269
4270

Figure S4c.6. Principal component analysis on West Eurasian samples coloured by individual polygenic scores.



4271

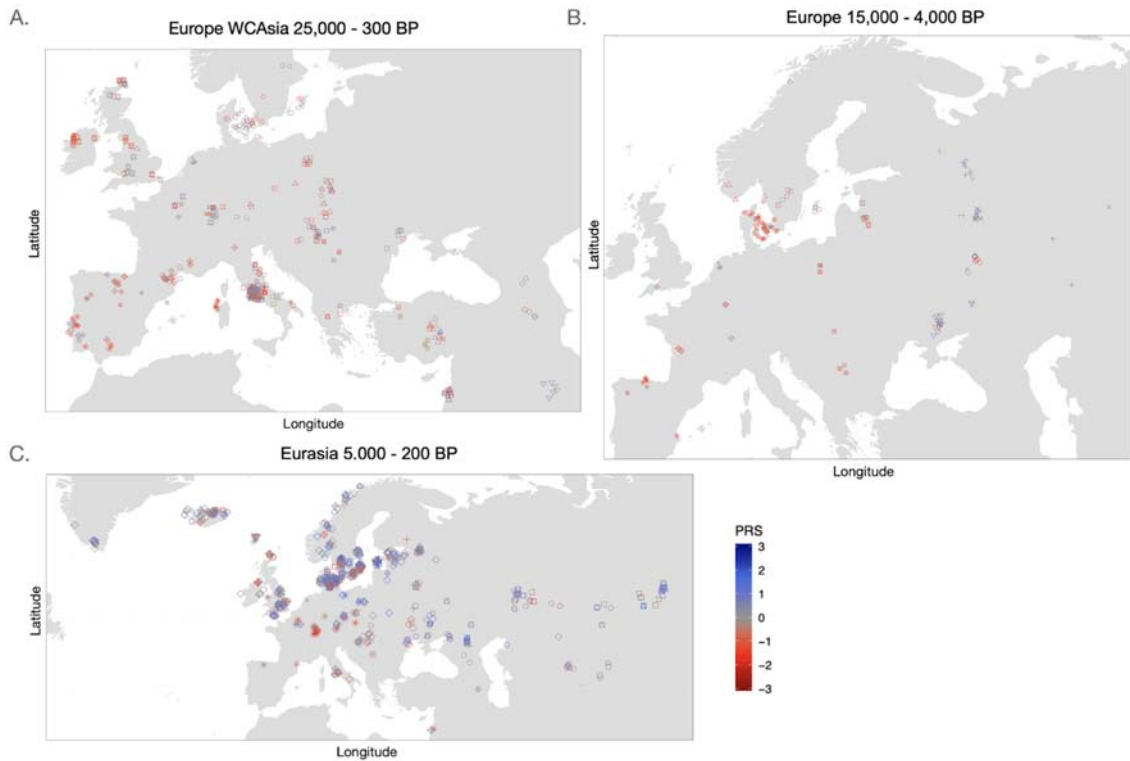
4272

4273

4274

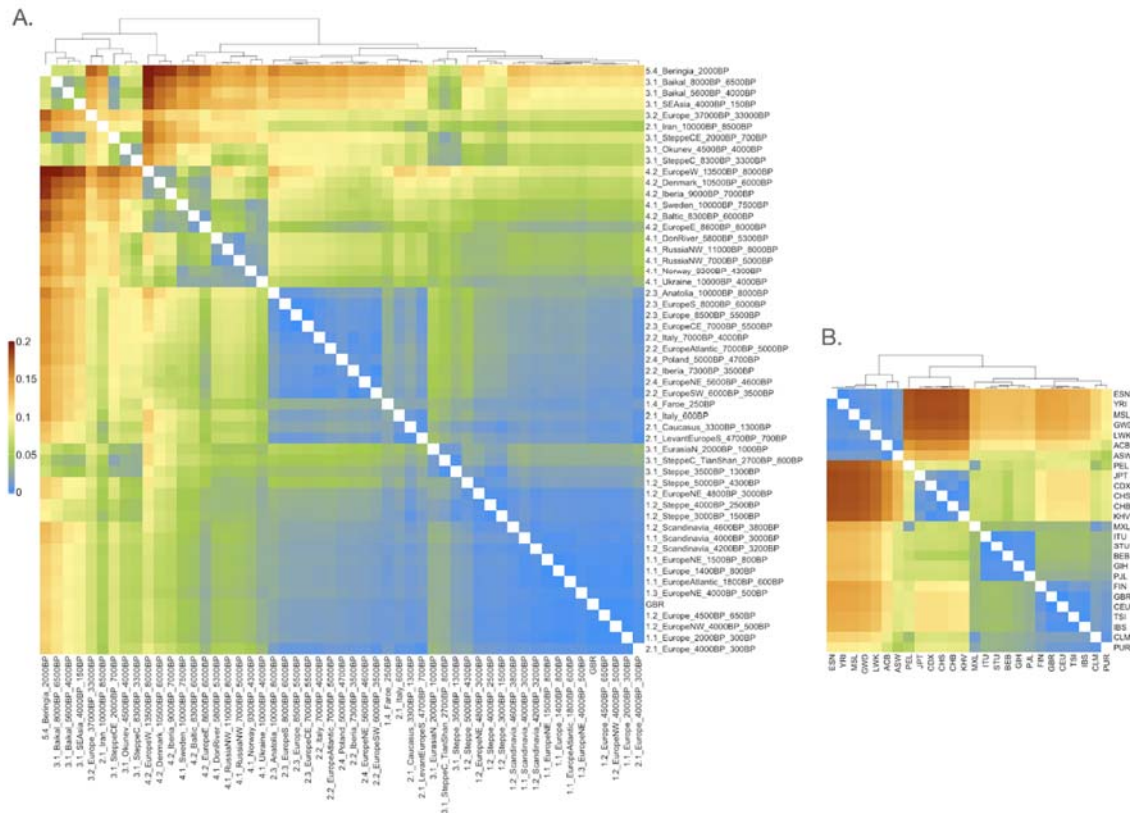
4275

Figure S4c.7. A. Principal component analysis on 1,332 Eurasian samples coloured by individual polygenic scores. B. World map of all Eurasian samples coloured by polygenic scores for height.



4276
4277
4278

Figure S4c.8. Map of Western Eurasian polygenic scores for height.



4279

4280
4281
4282
4283
4284
4285

Figure S4c.9. Genetic differentiation between A. ancient clusters and B. modern population in 1000 Genomes Project data set. Pairwise Fst

Population Code	Population Description	Super Population Code
CHB	Han Chinese in Beijing, China	EAS
JPT	Japanese in Tokyo, Japan	EAS
CHS	Southern Han Chinese	EAS
CDX	Chinese Dai in Xishuangbanna, China	EAS
KHV	Kinh in Ho Chi Minh City, Vietnam	EAS
CEU	Utah Residents (CEPH) with Northern and Western European Ancestry	EUR
TSI	Toscani in Italia	EUR
FIN	Finnish in Finland	EUR
GBR	British in England and Scotland	EUR
IBS	Iberian Population in Spain	EUR
YRI	Yoruba in Ibadan, Nigeria	AFR
LWK	Luhya in Webuye, Kenya	AFR
GWD	Gambian in Western Divisions in the Gambia	AFR
MSL	Mende in Sierra Leone	AFR
ESN	Esan in Nigeria	AFR
ASW	Americans of African Ancestry in SW USA	AFR
ACB	African Caribbeans in Barbados	AFR
MXL	Mexican Ancestry from Los Angeles USA	AMR
PUR	Puerto Ricans from Puerto Rico	AMR
CLM	Colombians from Medellin, Colombia	AMR
PEL	Peruvians from Lima, Peru	AMR
GIH	Gujarati Indian from Houston, Texas	SAS
PJL	Punjabi from Lahore, Pakistan	SAS
BEB	Bengali from Bangladesh	SAS
STU	Sri Lankan Tamil from the UK	SAS
ITU	Indian Telugu from the UK	SAS

4286
4287
4288
4289

Figure S4c.10. Full population descriptions of 1000 Genomes Project panels used in our analysis.

4290 References

- 4291 1. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.
4292 *Nature* 562, 203–209 (2018).
- 4293 2. Berisa, T. & Pickrell, J. K. Approximately independent linkage disequilibrium blocks in
4294 human populations. *Bioinformatics* 32, 283–285 (2016).
- 4295 3. Watanabe, K. *et al.* A global overview of pleiotropy and genetic architecture in
4296 complex traits. doi:10.1101/500090.
- 4297 4. Canela-Xandri, O., Rawlik, K. & Tenesa, A. An atlas of genetic associations in UK
4298 Biobank. *Nat. Genet.* 50, 1593–1599 (2018).
- 4299 5. Sohail, M. *et al.* Polygenic adaptation on height is overestimated due to uncorrected
4300 stratification in genome-wide association studies. *Elife* 8, (2019).
- 4301 6. Berg, J. J. & Coop, G. A population genetic signal of polygenic adaptation. *PLoS*
4302 *Genet.* 10, e1004412 (2014).
- 4303 7. Racimo, F., Berg, J. J. & Pickrell, J. K. Detecting Polygenic Adaptation in Admixture
4304 Graphs. *Genetics* 208, 1565–1584 (2018).
- 4305 8. Consortium, T. 1000 G. P. & The 1000 Genomes Project Consortium. A global
4306 reference for human genetic variation. *Nature* vol. 526 68–74 (2015).
- 4307 9. Martin, A. R. *et al.* Clinical use of current polygenic risk scores may exacerbate health
4308 disparities. *Nat. Genet.* 51, 584–591 (2019).
- 4309 10. Martin, A. R. *et al.* Human Demographic History Impacts Genetic Risk Prediction
4310 across Diverse Populations. *Am. J. Hum. Genet.* 100, 635–649 (2017).
- 4311 11. Weir, B. S. & Cockerham, C. C. ESTIMATING F-STATISTICS FOR THE ANALYSIS
4312 OF POPULATION STRUCTURE. *Evolution* 38, 1358–1370 (1984).
- 4313 12. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* 27, 2156–
4314 2158 (2011).
- 4315 13. Mathieson, I. *et al.* Genome-wide patterns of selection in 230 ancient Eurasians.
4316 *Nature* 528, 499–503 (2015).

- 4317 14. Cox, S. L., Ruff, C. B., Maier, R. M. & Mathieson, I. Genetic contributions to variation
4318 in human stature in prehistoric Europe. *Proc. Natl. Acad. Sci. U. S. A.* 116, 21484–21492
4319 (2019).
- 4320 15. Ju, D. & Mathieson, I. The evolution of skin pigmentation associated variation in West
4321 Eurasia. doi:10.1101/2020.05.08.085274.
- 4322 16. Gao, C. *et al.* Reference values for lung function screening in 10- to 81-year-old,
4323 healthy, never-smoking residents of Southeast China. *Medicine* 97, (2018).
- 4324 17. Jensen, B. *et al.* Ethnic differences in fat and muscle mass and their implication for
4325 interpretation of bioelectrical impedance vector analysis. *Appl. Physiol. Nutr. Metab.* 44,
4326 619–626 (2019).
- 4327 18. Bosy-Westphal, A., Plachta-Danielzik, S., Dörhöfer, R.-P. & Müller, M. J. Short
4328 stature and obesity: positive association in adults but inverse association in children and
4329 adolescents. *Br. J. Nutr.* 102, 453–461 (2009).
- 4330 19. Sjödin, A. M. *et al.* The influence of physical activity on BMR. *Med. Sci. Sports*
4331 *Exercise* 28, 85 (1996).
- 4332 20. [No title]. <https://www.ahajournals.org/doi/10.1161/HYPERTENSIONAHA.119.14028>.
- 4333 21. Zanetti, D. *et al.* Urinary Albumin, Sodium, and Potassium and Cardiovascular
4334 Outcomes in the UK Biobank: Observational and Mendelian Randomization Analyses.
4335 *Hypertension* 75, 714–722 (2020).
- 4336 22. O'Donnell, M. *et al.* Urinary sodium and potassium excretion, mortality, and
4337 cardiovascular events. *N. Engl. J. Med.* 371, 612–623 (2014).
- 4338 23. Marciniak, S. & Perry, G. H. Harnessing ancient genomes to study the history of
4339 human adaptation. *Nature Reviews Genetics* vol. 18 659–674 (2017).
- 4340 24. Balbuena, L., Bowen, R., Baetz, M. & Marwaha, S. Mood Instability and Irritability as
4341 Core Symptoms of Major Depression: An Exploration Using Rasch Analysis. *Front.*
4342 *Psychiatry* 7, (2016).
- 4343 25. Ward, J. *et al.* Genome-wide analysis in UK Biobank identifies four loci associated
4344 with mood instability and genetic correlation with major depressive disorder, anxiety

4345 disorder and schizophrenia. *Transl. Psychiatry* 7, 1–9 (2017).

4346 26. Berg, J. J. *et al.* Reduced signal for polygenic adaptation of height in UK Biobank.

4347 doi:10.1101/354951.

4348 27. Refoyo-Martínez, A. *et al.* How robust are cross-population signatures of polygenic

4349 adaptation in humans? doi:10.1101/2020.07.13.200030.

4350

4351 4d) Identifying candidates for positive selection using patterns 4352 of ancient population differentiation

4353 Alba Refoyo Martínez¹, Fernando Racimo¹

4354

4355 ¹Lundbeck Foundation GeoGenetics Centre, GLOBE Institute, University of Copenhagen,
4356 Copenhagen, Denmark

4357

4358 Introduction

4359 We aimed to detect whether there is evidence of positive selection in the past 15,000 years
4360 by searching for loci with strong differentiation in allele frequencies both between populations
4361 and across time.

4362 Methods

4363 We worked with the imputed dataset (Supplementary Note S2). We applied several variant-
4364 level filters: 1) we only used SNPs with MAF > 5% , 2) genotype missingness rate < 50% and
4365 3) variants where <10% individuals had post-imputation genotype probability (GP) <= 0.8.

4366 We used *pcadapt*¹ - a method for detection of allele frequency outliers based on principal
4367 component analysis - to search for loci that might have evidence for positive selection in the
4368 past. After performing a PCA on the ancient genomic data, we used a scree plot to visualise
4369 the percentage of variance explained by each PC.

4370

4371 We performed two scans to look for candidates under positive selection. First, we performed
4372 a “eurasian” scan, in which we used the first three principal components of a PCA of all
4373 ancient genomes (higher components explained less than 1% of the total variance in allele
4374 frequencies) (Figures **S4d.1**, **S4d.2** and **S4d.3**). We also performed a second scan of
4375 selection restricting only to ancient West Eurasian individuals (excluding ancient Siberian
4376 populations) and we also used the first three components of the PCA (Figures **S4d.4**, **S4d.5**

4377 **and S4d.6**). We call this the “west-eur” scan. We also performed a third scan, which we call
4378 “hg-neo” scan, in which we only tested for significant loadings corresponding to the
4379 component that separates hunter-gatherers (WHG+EHG) from Neolithic farming peoples in
4380 the aforementioned Western Eurasian PCA (first principal component). The latter scan
4381 should serve to find loci with particularly strong allele frequencies between these groups with
4382 two distinct modes of subsistence. Figure **S4d.2** and **S4d.5** shows Manhattan plots for each
4383 of the three scans.

4384

4385 We selected the top-scoring 300 SNPs with the lowest p-values and merged them into
4386 candidate regions if they were within 100kb of each other. Each region was labelled with the
4387 P-value of its highest scoring SNP, which was then used to rank regions. HGNC protein-
4388 coding genes within each region were retrieved using *biomaRt*². Table **S4d.1** lists the 25
4389 top-scoring candidate regions from the “eurasian” scan, while Tables **S4d.2** and list the top
4390 candidate regions from the “west-eur” and “hg-neo” scan, respectively.

4391

4392 None of our candidate regions seem to be affected by mapping biases. The F_j value for
4393 highest scoring SNP in each region is lower than 0.5 (Supplementary Note S4b).

4394 Results

4395 Eurasian scan

4396 In the Eurasian-wide scan, the strongest peak contains the gene SLC24A5, involved in skin
4397 and eye pigmentation, and previously reported to be under positive selection in Western
4398 Eurasia^{3–6} (Figure **S4d.2**). We also recover the region encompassing the *EDAR* gene which
4399 has been implicated in numerous studies of positive selection involving East Asian
4400 populations^{7–9}. Variants in this gene are associated with numerous ectoderm-related traits,
4401 including hair thickness and ectodermal dysplasia¹⁰ and tooth morphology¹¹.

4402

4403 We also found a candidate peak (chr2:25874547-26568094), containing several genes
4404 (*ASXL2*, *RAB10*, *HADHA*, *GPR113*) involved in glucose homeostasis mainly in response to a
4405 high fat diet. The *ASXL2* gene regulates skeletal, glucose, and adipocyte homeostasis. It
4406 promotes adipogenesis, the formation of osteoclasts and insulin resistance^{12–14}. A similar
4407 activity is carried out by *RAB10*, it also regulates glucose homeostasis by improving
4408 hyperglycemia, regulating at the skeletal muscle level^{15–17}. *HADHA* participates in long fatty
4409 acid oxidation for energy production in different tissues^{18–20}. Finally, *GPR113* participates in
4410 energy expenditure in fat tissue and glucose homeostasis. It improves insulin sensitivity and

4411 prevents obesity when it binds to bile acids ²¹. High allele frequencies in East Asian
4412 populations (Panel B of Figure **S4d.7**).
4413
4414 We found four peaks containing genes associated with cardiovascular disorders and obesity.
4415 One of these (chr20:18991679-19454079) overlaps with the SLC24A3 gene. This is a salt
4416 sensitivity gene which is significantly expressed in obese individuals ²²⁻²⁴ and it is associated
4417 with hypertension ²⁵. We also find a peak in an intergenic region (chr3:123433220-
4418 123889576) containing ROPN1 and KALRN, two genes involved in vascular disorders ²⁶⁻²⁸.
4419 The alternative allele at the top SNP in this region is at particularly high frequency in ancient
4420 Steppe populations. Another candidate region (chr1:234142067-234549596) contains
4421 SLC35F3, which codes for a thiamine transport and has been associated with hypertension
4422 in a Han Chinese cohort ^{29,30}. In the same region, we also find COA6, which has high
4423 expression in cardiac pathologies ³¹. The alternative allele frequency at the top SNP in the
4424 region is high in East Asian populations (Panel B of Figure **S4d.7**). Finally, the region
4425 (chr10:90592757 - 91009553) contains several genes (CH25H, FAS) associated with obesity
4426 and lipid metabolism ³²⁻³⁴, and immune responses ³⁵.

4427
4428 One of the top candidate regions (chr11:131077365-131516733) contains a gene - NTM -
4429 involved in neuropsychiatric disorders ^{25,36}, while another region (chr11:44634764 -
4430 45073989) contains a gene associated with schizophrenia, TSPAN18. It has previously
4431 shown that the schizophrenia-risk SNPs within this region are highly diverged between
4432 Europeans and East Asians ³⁷. In chromosome 7 (chr7:95947959-96347959), SLC25A13,
4433 which is highly associated with citrin deficiency in East Asian populations ^{38,39}. High
4434 alternative allele frequencies in East Asians (Panel A of Figure **S4d.7**) ^{40,41}.

4435

4436 West Eurasian scan

4437 In the scan restricting to ancient populations in Western Eurasia ("west-eur"), we recover
4438 three regions which are involved in skin, hair and eye pigmentation, and have been
4439 previously implicated in differences in these traits across present-day Eurasians: two in
4440 chromosome 15 containing the genes *SLC24A5/MYEF2/CTXN2* and *OCA2/HERC2*
4441 respectively and one in chromosome 5, containing gene *SLC45A2* ^{3,42-45}. Alternative allele
4442 frequencies shown in Figure **S4d.8**.

4443

4444 The region containing *LCT/MCM6* - responsible for lactase persistence in Europe - is also a
4445 candidate region in the selection scan ⁴⁶⁻⁴⁸. We also recover the TLR-1-6-10 gene cluster,

4446 which is known to be a target of selection in Europe and is associated with the immune
4447 response ⁴⁶⁻⁴⁸ (Panel A, Figure **S4d.8**).

4448

4449 Additionally, we find some new potentially important candidate regions for positive selection.
4450 The region showing the strongest evidence of selection is located in chromosome 6
4451 (chr6:134192815-134628278), around the SLC2A12 gene, which codes for a glucose
4452 transporter that participates in glucose homeostasis ^{49,50}. Variants in this gene are associated
4453 with mean corpuscular levels, heart diseases and height. The alternative allele of the
4454 highest-scoring SNP was at high frequency in hunter-gatherers but at much lower
4455 frequencies in Neolithic farmer populations and other, more recent, populations (Panel B,
4456 Figure **S4d.8**).

4457

4458 Another novel candidate region overlaps with the VAMP 5-8 gene cluster in chromosome
4459 2:85369379-85885211 a region associated with cardiovascular diseases ^{51,52}. In
4460 chr9:27009422-27434948, we found a peak overlapping the TEK gene, which codes for a
4461 tyrosine kinase receptor expressed in endothelial cells. This gene has an important role in
4462 angiogenesis and cardiovascular development and stability, and it is involved in several
4463 vascular disorders ⁵³⁻⁵⁵. Recent studies have also investigated its role in asthma and allergic
4464 conjunctivitis ⁵³⁻⁵⁶. Intermediate allele frequencies in Eastern hunter-gatherers in both
4465 regions (Panel B, Figure **S4d.8**). Region chr1:227020437 - 227877723 contains CDC42BPA,
4466 also known as MRCK α , an important gene involved in iron utilisation ⁵⁷ is involved in the
4467 erythropoiesis regulation ⁵⁸. The alternative allele at the top-scoring SNP in this region is at
4468 high frequencies found in hunter-gatherer groups, predominantly in eastern hunter-gatherers
4469 (Panel B, Figure **S4d.8**).

4470

4471 In chr15:38464638-38992430, we recovered RASGRP1, associated with immunity and
4472 related to systemic lupus erythematosus ⁵⁹, rheumatoid arthritis ⁶⁰ or Epstein-Barr virus ⁶¹
4473 among other disorders. The alternative allele at this SNP is at high frequencies in eastern
4474 hunter gatherers and other ancient Baltic populations (Panel B, Figure **S4d.8**). We also
4475 found a wide candidate region containing several high-scoring SNPs in chromosome
4476 16:66852047-67871804. The gene that falls in the highest peak of the region is ATP6V0D1,
4477 and it plays a very important role in the replication of influenza virus ^{62,63}. In this region, there
4478 are also several genes - *TPPP3/ZDHHC1* and *HSD11B2* - associated with obesity,
4479 cardiovascular diseases and hypertension ⁶⁴⁻⁷⁰. The alternative allele at the top-scoring SNP
4480 has intermediate allele frequencies in hunter-gatherer populations, and higher frequencies in
4481 eastern hunter-gatherers.

4482

4483 Neolithic vs. hunter-gatherer scan

4484 The “neo-hg” scan specifically recovers patterns of allele differentiation along the axis
4485 separating hunter-gatherer and farmer populations in West Eurasia. In this scan, we find a
4486 large number of high-scoring regions associated with lipid and sugar metabolism, and
4487 various metabolic disorders.

4488

4489 For example, we recover the *FADS* gene cluster, involved in lipid metabolism. This region is
4490 presumed to be important in the transition to a diet rich in grains, as a consequence of the
4491 expansion of agriculture in Europe and/or the demographic transitions subsequent to it ^{47,71–}
4492 ⁷³. This region is also found in the “west-eurasia” scan. The alternative allele frequency at the
4493 top-scoring SNP is very high in hunter-gatherer populations.

4494

4495 Another region is located in chromosome 22:31353354-31759255 and also contains genes
4496 involved in lipid metabolism: *PATZ1*, *LIMK2*, *MORC2* and *PLA2G3*. *PATZ1* down-regulates
4497 *FADS1* ⁷⁴, *LIMK2* shows particularly elevated expression in metabolic syndrome ⁷⁵, *MORC2*
4498 plays an important role in cellular lipid metabolism ^{76–78} and *PLA2G3* contributes to
4499 atherogenesis ^{79–81}. The top-scoring SNP in this region has high alternative allele frequencies
4500 in Neolithic farmer populations (Panel A, Figure **S4d.9**).

4501

4502 At chromosome 12:6875213-7366672, we find a region with several genes involved in lipid
4503 metabolism: *PTPN6*, *EMG1*, *PHB2*, *LPCAT3*, *C1S*. *LPCAT3* is essential in high fat diets and
4504 associated with oleic acid levels and linoleic acid ⁸² and its deficiency alters cholesterol
4505 promoting atherosclerosis ⁸³ and plays an important role in hyperuricemia ^{83,84}. *C1S* also
4506 plays a crucial role in innate immunity ⁸⁵ and has been recently associated with coronary
4507 syndrome ⁸⁶.

4508

4509 In chromosome 17:8655348-9223981, we found the *PIK3* family (*PIK3R5*, *PIK3R6*) which is
4510 involved in glucose homeostasis and plays an important role in obesity and insulin resistance
4511 in type 2 diabetes ^{87–89}.

4512

4513 In chromosome 11:27432440-27832440 we find a region containing *BDNF* (brain-derived
4514 neurotrophic factor), expression of which is associated with obesity ^{90,91}. It participates in the
4515 reduction of free fatty acids, cholesterol and glucose levels and enhances energy
4516 expenditure ⁹². Its expression has been shown to be suppressed with a high-fat sucrose diet
4517 ^{90,9394, 90,93}. High alternative allele frequency in HG (Panel A, Figure **S4d.9**).

4518

4519 One of the strongest peaks in this scan corresponds to SLC2A12, coding for a glucose
4520 transporter. The next peak is located in chromosome 17:79055998-79469799. SLC38A10,
4521 which falls in the tip of the peak, is involved in absorption of amino acids from the GI tract ⁹⁵
4522 and has been suggested to play a role in pathways involved in neurotransmission ^{96,97}.
4523 BAIAP2's ^{98,99} and AATK's expression ^{100,101}, which are also found in the same region, are
4524 related to high-fat and omega3 fatty acids. Higher alternative allele frequencies are found in
4525 hunter-gatherers, in particular in Eastern-HG (Panel A, Figure **S4d.9**).

4526

4527 The highest peak in chromosome 4 contains several genes involved in alcohol metabolism,
4528 ADH1B, ADH1C and ADH7 ^{102,103}. In Panel B, Figure **S4d.9**, we can see the highest
4529 alternative allele frequencies for the oldest samples, which correspond to hunter-gatherers.

4530

4531 Two of the top peaks are related to innate immune response in humans. In chromosome
4532 14:73102086-73502086, ZFYVE1 is involved in TLR3-mediated immune response and
4533 regulates antiviral response ^{104,105}. In chromosome 3:98028061-98476960, GPR15 is related
4534 to immune tolerance and it regulates the homeostasis in the intestine mucosa ^{104,106,107}. In
4535 chromosome 18:46132358-46558884, SMAD7 is associated with inflammatory bowel
4536 diseases such as crohn's disease ¹⁰⁸⁻¹¹⁰.

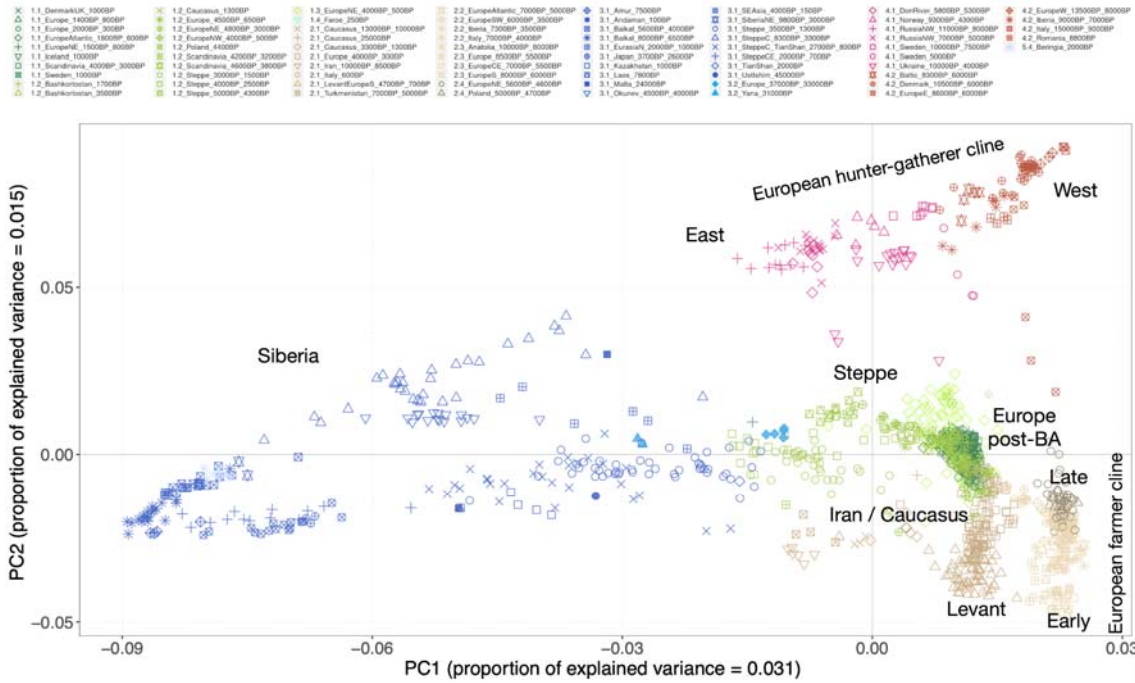
4537

4538 Two regions are related to brain disorders. In chromosome 1:110588519-111127548,
4539 several genes regulate neuronal ion channels: KCNC4, SLC6A17, and STRIP1. Mutations in
4540 SLC6A17 cause intellectual disabilities associated with speech impairment and behavioural
4541 problems ¹¹¹, while the KCNC gene family has also been associated with intellectual
4542 disability ¹¹².

4543

4544 Figures

4545

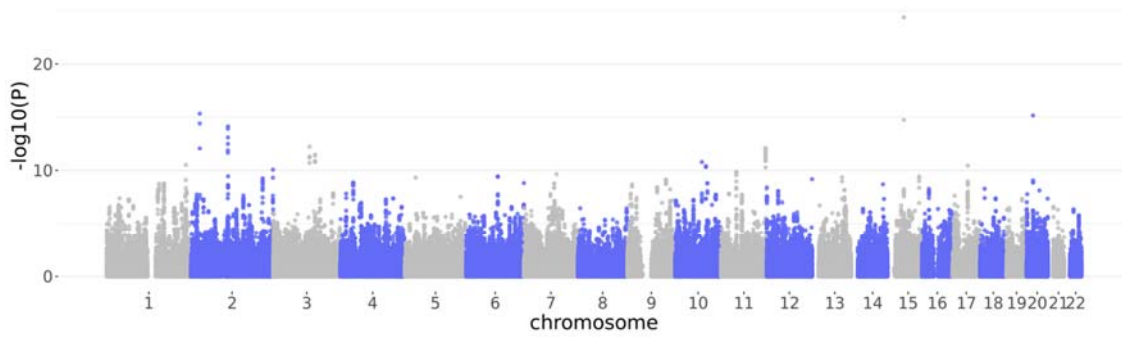


4546

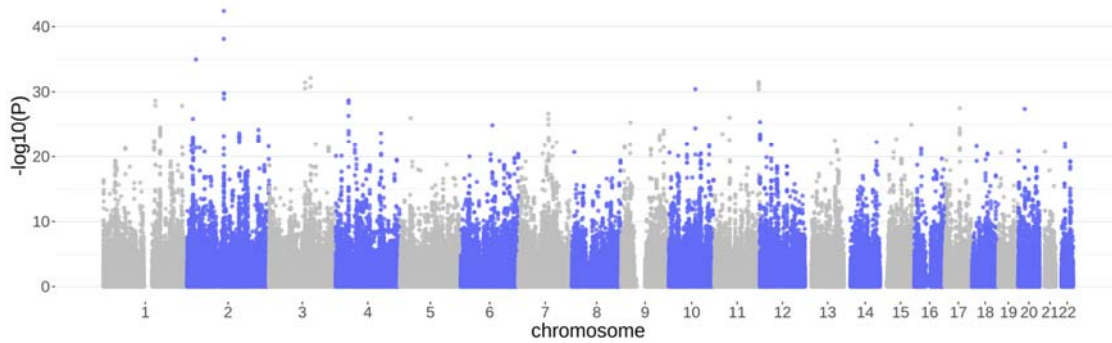
4547 **Figure S4d.1. Principal component analysis on 1402 Eurasian samples.** The first
 4548 component explains 3.1% of the variance and separates East Asian, Steppe and European
 4549 samples. The second component separates farmers and Hunter-gatherers (1.5%).

4550 Eurasian scan

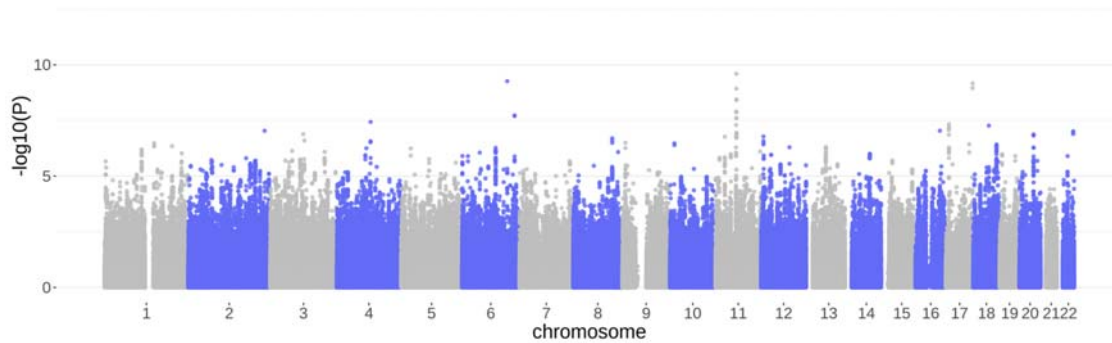
Manhattan plot using k=3



Manhattan plot using k=3 :component-wise pc=1

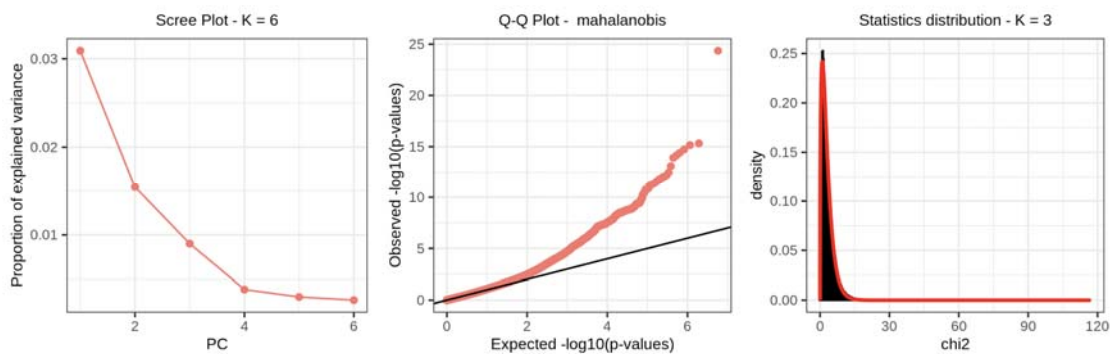


Manhattan plot using k=3 :component-wise pc=2



4551
4552
4553
4554
4555

Figure **S4d.2**. Genome scan using pcadapt k=3. A. Method: mahalanobis distances. Manhattan plot scanning the whole genome. B and C. Component-wise genome scans for component PC1 and component PC2, respectively.

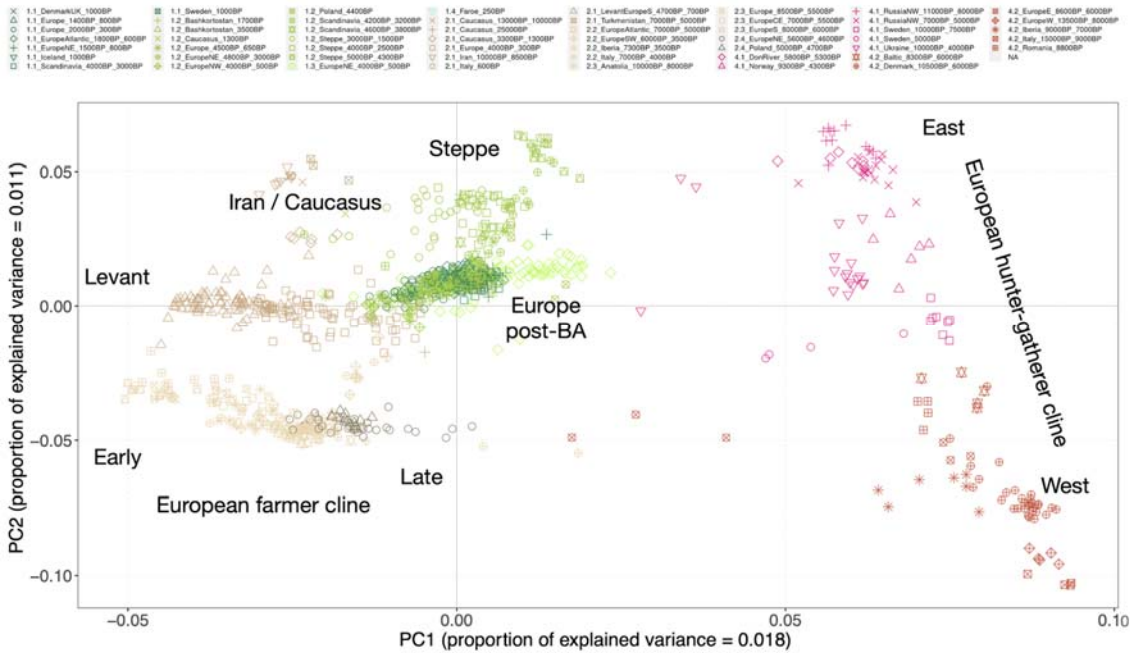


4556

4557 Figure S4d.3. A. Scree plot showing proportion of explained variance of the first PCs in the
 4558 pcadapt analysis. B. Q-Q plot using mahalanobis method for K=3. Distribution of pcadapt
 4559 scores (k=3) compared to chi-squared distribution with one degree of freedom (red line).
 4560

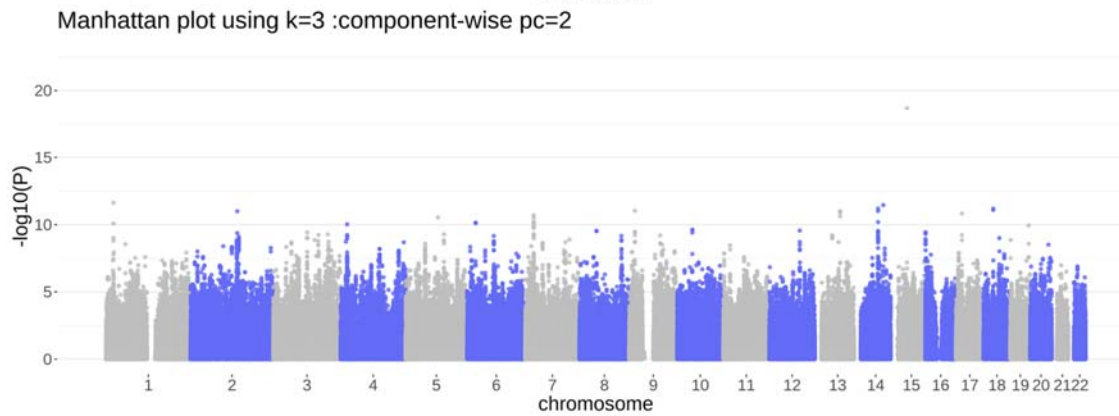
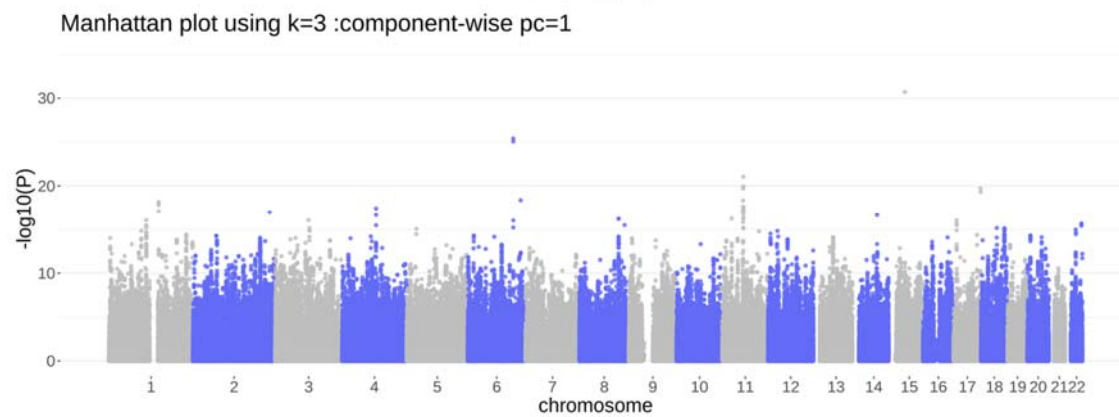
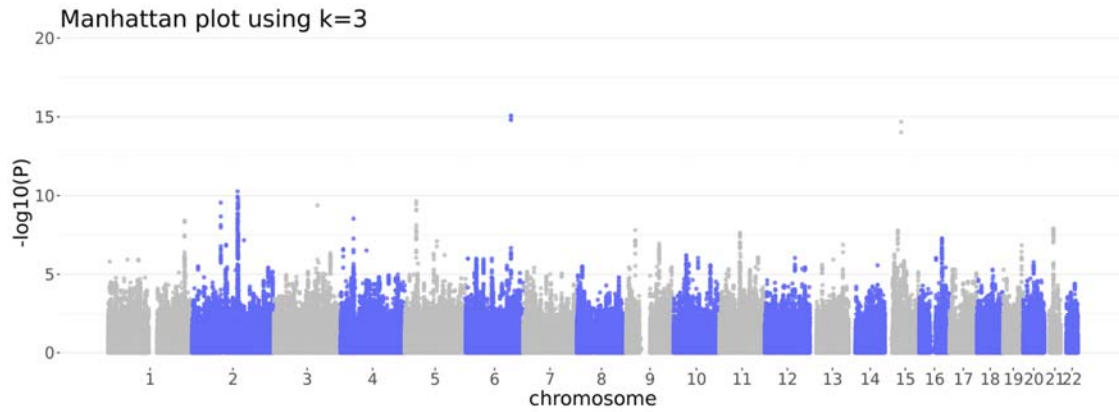
4561 West-Eurasian

4562



4564

4565 **Figure S4d.4. Principal component analysis on 1165 Eurasian samples.** The first
 4566 component explains 1.8% of the variance and separates East Asian, Steppe and European
 4567 samples. The second component separates farmers and Hunter-gatherers (1.1%).



4568

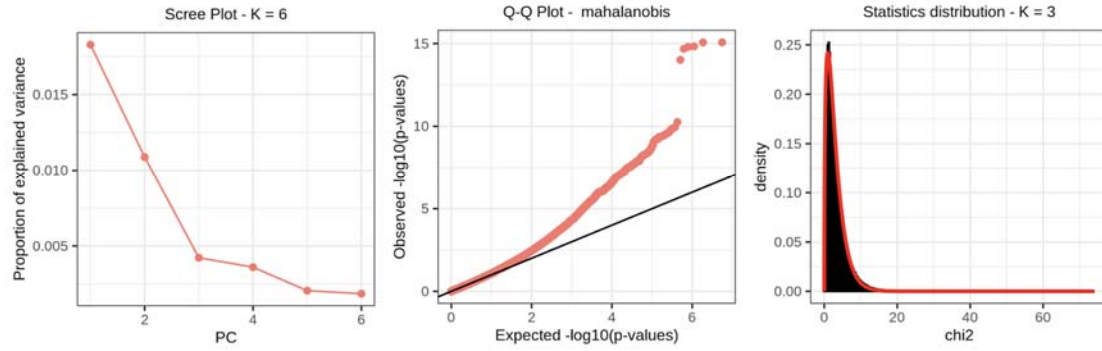
4569

4570

4571

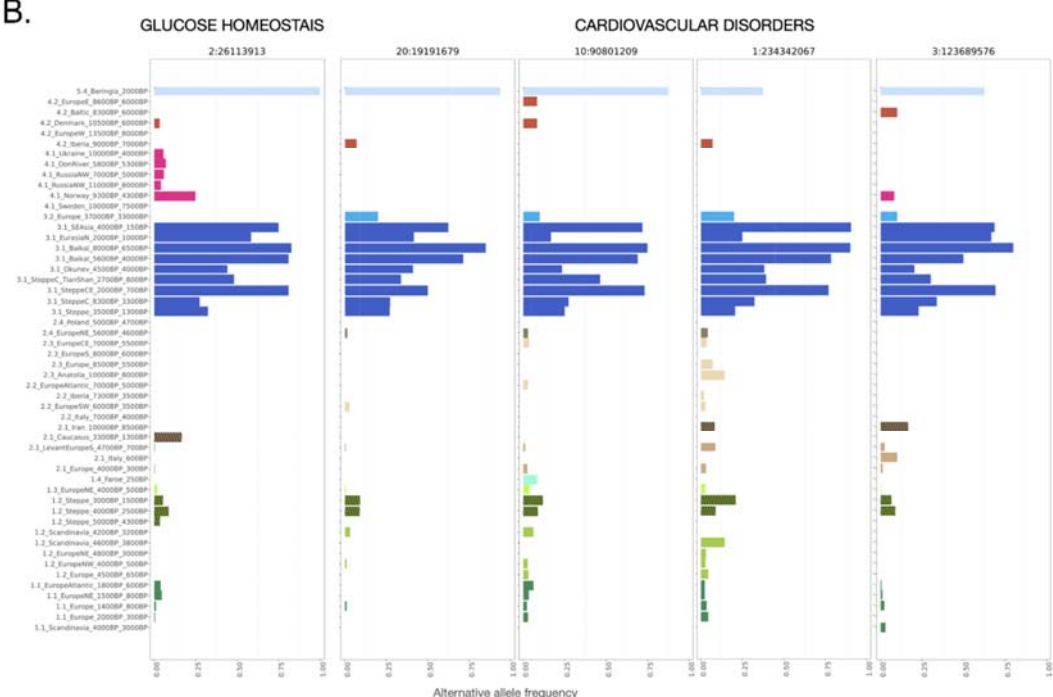
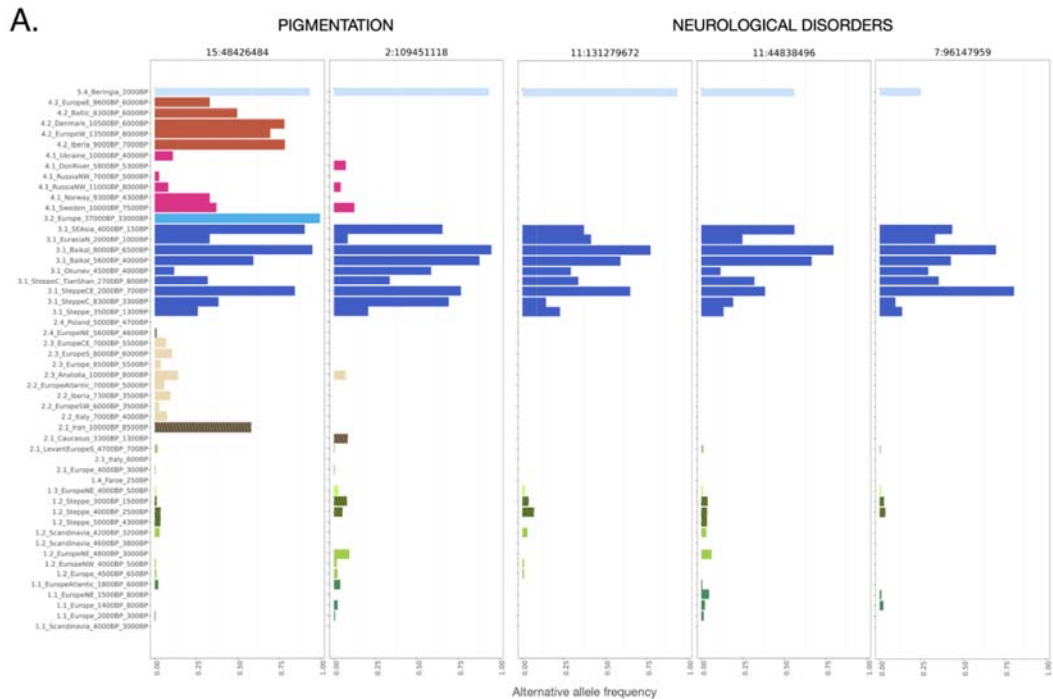
4572

Figure **S4d.5**. Genome scan using pcadapt k=3. A. Method: mahalanobis distances. Manhattan plot scanning the whole genome. B and C. Component-wise genome scans for component PC1 and component PC2, respectively.



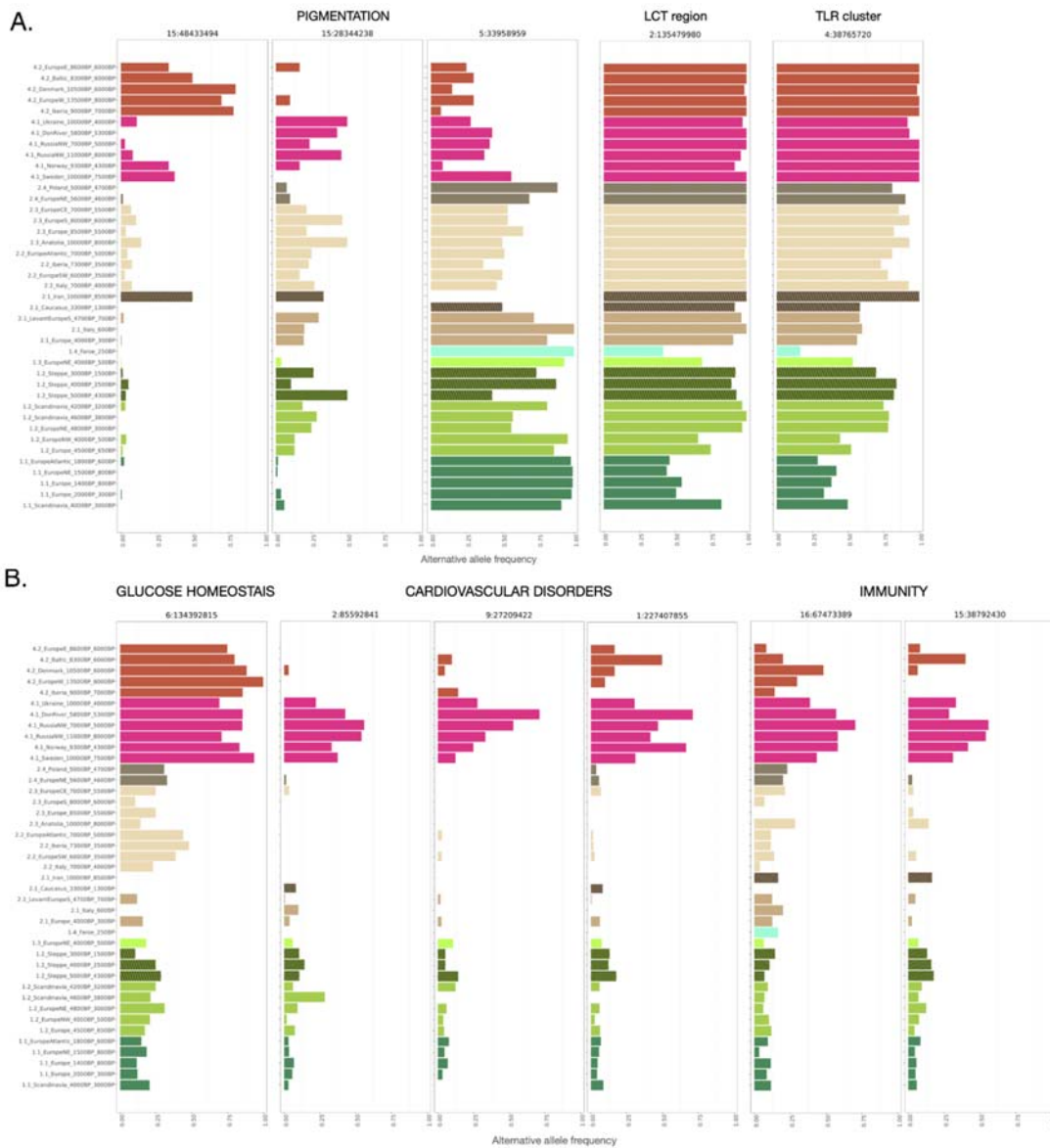
4573
 4574
 4575
 4576
 4577
 4578

Figure **S4d.6**. A. Scree plot showing proportion of explained variance of the first PCs in the pcadapt analysis. B. Q-Q plot using mahalanobis method for K=3. Distribution of pcadapt scores (k=3) compared to chi-squared distribution with one degree of freedom (red line).



4579
4580
4581
4582
4583
4584

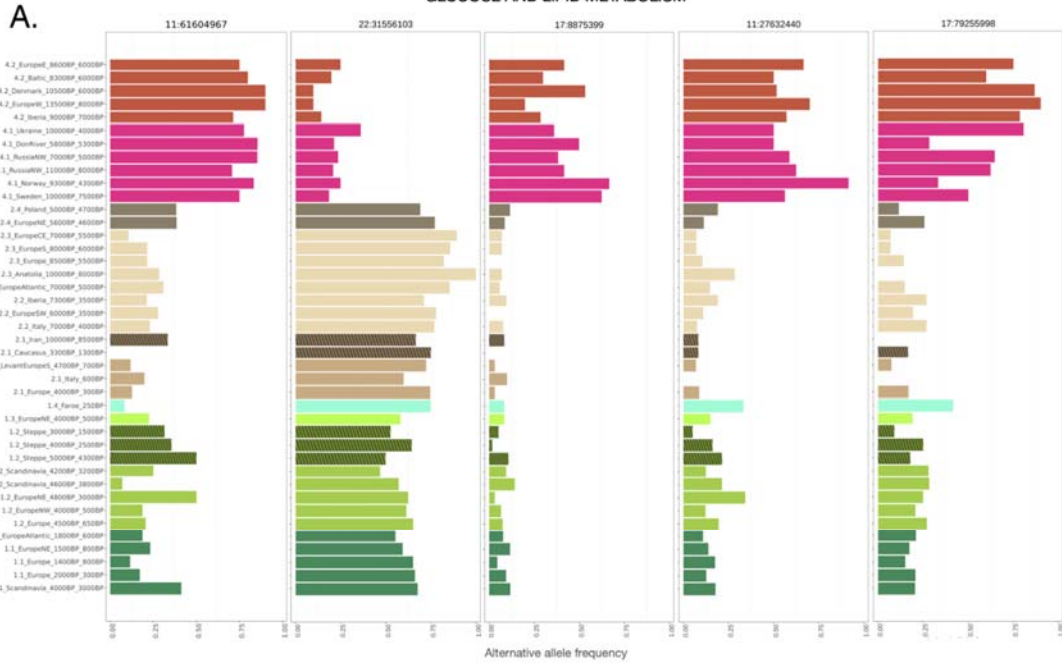
Figure **S4d.7**. Alternative allele frequencies of the position with the lowest p-value in the top regions for the eurasian scan.



4585
4586
4587
4588
4589
4590

Figure **S4d.8**. Alternative allele frequencies of the position with the lowest p-value in the top regions for the West- Eurasian scan.

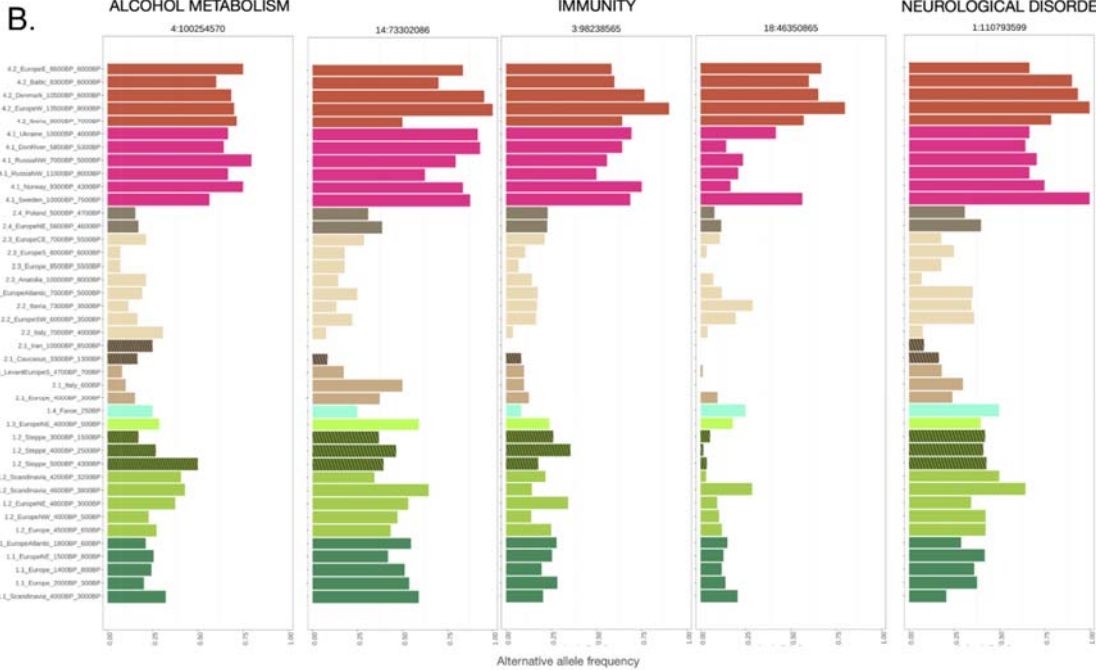
GLUCOSE AND LIPID METABOLISM



ALCOHOL METABOLISM

IMMUNITY

NEUROLOGICAL DISORDERS



4591
4592
4593
4594
4595
4596
4597

Figure **S4d.9**. Alternative allele frequencies of the position with the lowest p-value in the top regions for the HG West- Eurasian scan.

chr	start	end	bestpos	minpvalue	ensembl	hgnc	disease_trait_best	disease_trait
15	48211821	48714309	48426484	4.40E-25	ENSG00000188467, ENSG00000104177, ENSG00000233932, ENSG00000074803, ENSG00000128951, ENSG00000166147	SLC24A5, MYEF2, CTXN2, SLC12A1, DUT, FBN1	ENSG00000188467: Hair color, Skin pigmentation, Body mass index, Skin, hair and eye pigmentation (multivariate analysis), Eye color, Eye color (brightness), Eye color (saturation), Iris color (a* coordinate), Eye color (brightness), Iris color (saturation), Iris color (a* coordinate), Skin reflectance (Melanin index), Iris color (b* coordinate); ENSG00000104177: Skin pigmentation, Hair color; ENSG00000233932: ; ENSG00000074803: Systemic lupus erythematosus, Longevity, Skin reflectance (Melanin index), Lymphocyte counts; ENSG00000128951: Protein quantitative trait loci; ENSG00000166147: Breast cancer, Spherical equivalent (joint analysis main effects and education interaction), Height, Pulse pressure, Refractive error, Spherical equivalent, Systolic blood pressure, Thoracic aortic aneurysms and dissections, Colorectal cancer, Spherical equivalent or myopia (age of diagnosis), Central corneal thickness, Intracranial, abdominal aortic or thoracic aortic aneurysm (pleiotropy), Systolic blood pressure x alcohol consumption interaction (2df test), Pulse pressure x alcohol consumption interaction (2df test), Skin reflectance (Melanin index), Spontaneous coronary artery dissection, Waist circumference adjusted for body mass index, Intraocular pressure, Macular thickness, Heel bone mineral density, Lung function (FEV1/FVC)	ENSG00000188467: Hair color, Skin pigmentation, Body mass index, Skin, hair and eye pigmentation (multivariate analysis), Eye color, Eye color (brightness), Eye color (saturation), Iris color (a* coordinate), Skin reflectance (Melanin index), Iris color (b* coordinate); ENSG00000104177: Skin pigmentation, Hair color; ENSG00000233932: ; ENSG00000074803: Systemic lupus erythematosus, Longevity, Skin reflectance (Melanin index), Lymphocyte counts; ENSG00000128951: Protein quantitative trait loci; ENSG00000166147: Breast cancer, Spherical equivalent (joint analysis main effects and education interaction), Height, Pulse pressure, Refractive error, Spherical equivalent, Systolic blood pressure, Thoracic aortic aneurysms and dissections, Colorectal cancer, Spherical equivalent or myopia (age of diagnosis), Central corneal thickness, Intracranial, abdominal aortic or thoracic aortic aneurysm (pleiotropy), Systolic blood pressure x alcohol consumption interaction (2df test), Pulse pressure x alcohol consumption interaction (2df test), Skin reflectance (Melanin index), Spontaneous coronary artery dissection, Waist circumference adjusted for body mass index, Intraocular pressure, Macular thickness, Heel bone mineral density, Lung function (FEV1/FVC)
2	25874547	26568094	26113913	4.68E-16	ENSG00000138101, ENSG00000138101	DTNB, ASXL2	ENSG00000084731: Atrial fibrillation, Type 2 diabetes, Height, Red blood cell count	ENSG00000138101: Multiple myeloma, LDL cholesterol, Multiple myeloma (IgH translocation), Type 2 diabetes, Cutaneous

			(rs78404020)		43970, ENSG00000084731, ENSG00000084733, ENSG00000157833, ENSG00000084754, ENSG00000138029, ENSG00000173567, ENSG00000138018	KIF3C, RAB10, , GARE ML, HADH A, HADH B, GPR1 13, EPT1		malignant melanoma, Nevus count or cutaneous melanoma, B-cell malignancies (chronic lymphocytic leukemia, Hodgkin lymphoma or multiple myeloma) (pleiotropy), C-reactive protein levels, Response to platinum-based chemotherapy (cisplatin), Waist circumference adjusted for body mass index, Body mass index, Height; ENSG00000143970: CTACK levels, Hemoglobin levels, Red cell distribution width, Hemoglobin concentration, Thyroid stimulating hormone levels, Alcohol consumption (drinks per week) (MTAG), Hair colour, Mean corpuscular hemoglobin; ENSG00000084731: Atrial fibrillation, Type 2 diabetes, Height, Red blood cell count; ENSG00000084733: Immune response to smallpox vaccine (IL-6), Mean corpuscular volume, Apolipoprotein B levels, Mean corpuscular hemoglobin; ENSG00000157833: Red cell distribution width; ENSG00000084754: ; ENSG00000138029: ; ENSG00000173567: ; ENSG00000138018: Coronary artery disease, Apolipoprotein A1 levels
20	18991679	19454079	19191679 (rs4141981)	6.92E-16	ENSG00000185052	SLC24A3	ENSG00000185052: Matrix metalloproteinase levels, Pulmonary function decline, Age at smoking initiation in chronic obstructive pulmonary disease, Metabolite levels, QT interval (sulfonylurea treatment interaction), Migraine, Pulse pressure, Smoking status (ever vs never smokers), Cataracts (operation), Mean platelet volume, Diastolic blood pressure, Lifetime smoking index, Psychosis (atypical), Daytime sleepiness, Nicotine dependence and major depression (severity of comorbidity), Nicotine dependence symptom count, Platelet	ENSG00000185052: Matrix metalloproteinase levels, Pulmonary function decline, Age at smoking initiation in chronic obstructive pulmonary disease, Metabolite levels, QT interval (sulfonylurea treatment interaction), Migraine, Pulse pressure, Smoking status (ever vs never smokers), Cataracts (operation), Mean platelet volume, Diastolic blood pressure, Lifetime smoking index, Psychosis (atypical), Daytime sleepiness, Nicotine dependence and major depression (severity of comorbidity), Nicotine dependence symptom count, Platelet

							smoking index, Psychosis (atypical), Daytime sleepiness, Nicotine dependence and major depression (severity of comorbidity), Nicotine dependence symptom count, Platelet count, Multisite chronic pain, Breast cancer and/or colorectal cancer, Educational attainment (years of education), Educational attainment (MTAG), Highest math class taken (MTAG), Height, Menarche (age at onset), Lung function (FEV1/FVC), Smoking status	count, Multisite chronic pain, Breast cancer and/or colorectal cancer, Educational attainment (years of education), Educational attainment (MTAG), Highest math class taken (MTAG), Height, Menarche (age at onset), Lung function (FEV1/FVC), Smoking status
2	108527043	109781319	109451118 (rs72627476)	7.39E-15	ENSG00000115665, ENSG00000196228, ENSG00000198203, ENSG00000198075, ENSG00000135968, ENSG00000169756, ENSG00000153201, ENSG00000163006, ENSG00000135960, ENSG00000172985	SLC5A7, SULT1C3, SULT1C2, SULT1C4, GCC2, LIMS1, RANBP2, CCDC138, EDAR, SH3RF3	ENSG00000163006: Lobe attachment (rater scored), Birth weight	ENSG00000115665: ; ENSG00000196228: Systolic blood pressure, Diastolic blood pressure; ENSG00000198203: Lobe size; ENSG00000198075: Ear protrusion, Lobe attachment, Helix rolling, Lobe attachment (rater-scored or self-reported), Lobe attachment (rater scored); ENSG00000135968: Low density lipoprotein cholesterol levels; ENSG00000169756: LDL cholesterol levels, Apolipoprotein B levels, Low density lipoprotein cholesterol levels, Excessive hairiness, Monocyte chemoattractant protein-1 levels, Total cholesterol levels, LDL cholesterol, Eyebrow thickness, Birth weight, Corneal endothelial cell density; ENSG00000153201: Attention deficit hyperactivity disorder; ENSG00000163006: Lobe attachment (rater scored), Birth weight; ENSG00000135960: Monobrow thickness, Scalp hair shape, Beard thickness, Eyebrow thickness, Ear

								<p>protrusion, Tragus size, Ear morphology, Lobe size, Lobe attachment, Helix rolling, lower facial morphology traits (quantitative measurement), Lobe attachment (rater scored), Straight vs curly hair, Blood protein levels, Gamma glutamyl transferase levels in excessive alcohol consumption, Male-pattern baldness, Monobrow, Thick vs thin eyebrows, Lung function (FEV1), Lung function (FVC), Balding type 1, Hair color; ENSG00000172985: Neurocognitive impairment in HIV-1 infection (dichotomous), Height, Cerebrospinal fluid sTREM-2 levels, Heel bone mineral density, Pediatric bone mineral density (spine), Bitter taste perception (6-n-propylthiouracil) in obesity with metabolic syndrome, Interleukin-5 levels, General cognitive ability, Rapid automatised naming of digits, Corneal endothelial cell density, Highest math class taken (MTAG)</p>
3	106950936	107415936	107150936	6.02E-13	ENSG00000138483, ENSG00000114439	CCDC54, BBX	ENSG00000138483:	<p>ENSG00000138483 ; ENSG00000114439: Smoking behaviour, Metabolite levels, Neuroticism, Irritable mood, Experiencing mood swings, Feeling worry, Male-pattern baldness, Depressed affect, Neuroticism, Waist-to-hip ratio adjusted for BMI (adjusted for smoking behaviour), Waist-to-hip ratio adjusted for BMI, Measles, Small cell lung carcinoma, Body mass index, Mood instability, Waist-hip ratio, White blood cell count, Neutrophil count, Platelet count, General cognitive ability, Multisite chronic pain, Monocyte count, Mean platelet volume, Cognitive ability, years of educational attainment or schizophrenia (pleiotropy), Educational attainment (MTAG), Cognitive performance (MTAG), Height, Balding type</p>

								1, Educational attainment (years of education), Highest math class taken (MTAG)
1 1	1310773 65	13151 6733	13127 9672	7.74E- 13	ENSG000001 82667	NTM	ENSG00000182667: Bipolar disorder and schizophrenia, Sunburns, Obesity-related traits, Male fertility, Asperger disorder, Educational attainment (years of education), Aggressiveness in attention deficit hyperactivity disorder, Itch intensity from mosquito bite adjusted by bite size, Food addiction, Feeling worry, Refractive error, Smoking status (ever vs never smokers), Leisure sedentary behaviour (television watching), Leisure sedentary behaviour (computer use), Spherical equivalent, Myopia, Cardiac Troponin-T levels, Chronic obstructive pulmonary disease-related biomarkers, Metabolite levels, Serum polyunsaturated fatty acid concentration x sex interaction in metabolic syndrome, Myopia (age of diagnosis), Spherical equivalent or myopia (age of diagnosis), Plasma anti-thyroglobulin levels, C-reactive protein levels, Schizophrenia, Peripheral arterial disease (traffic-related air pollution interaction), Body mass index, Intake of total sugars,	ENSG00000182667: Bipolar disorder and schizophrenia, Sunburns, Obesity-related traits, Male fertility, Asperger disorder, Educational attainment (years of education), Aggressiveness in attention deficit hyperactivity disorder, Itch intensity from mosquito bite adjusted by bite size, Food addiction, Feeling worry, Refractive error, Smoking status (ever vs never smokers), Leisure sedentary behaviour (television watching), Leisure sedentary behaviour (computer use), Spherical equivalent, Myopia, Cardiac Troponin-T levels, Chronic obstructive pulmonary disease-related biomarkers, Metabolite levels, Serum polyunsaturated fatty acid concentration x sex interaction in metabolic syndrome, Myopia (age of diagnosis), Spherical equivalent or myopia (age of diagnosis), Plasma anti-thyroglobulin levels, C-reactive protein levels, Schizophrenia, Peripheral arterial disease (traffic-related air pollution interaction), Body mass index, Intake of total sugars, Familial lung adenocarcinoma, Mental health study participation (completed survey), Reaction time, Smoking initiation (ever regular vs never regular), Age of smoking initiation (MTAG), Cognitive ability, years of educational attainment or schizophrenia (pleiotropy), Smoking initiation (ever regular vs never regular) (MTAG), Pre-treatment viral load in HIV-1 infection, Educational attainment (MTAG), Highest math class taken (MTAG), Smoking status, White blood cell count, Lung function (FEV1/FVC)

							<p>Familial lung adenocarcinoma, Mental health study participation (completed survey), Reaction time, Smoking initiation (ever regular vs never regular), Age of smoking initiation (MTAG), Cognitive ability, years of educational attainment or schizophrenia (pleiotropy), Smoking initiation (ever regular vs never regular) (MTAG), Pre-treatment viral load in HIV-1 infection, Educational attainment (MTAG), Highest math class taken (MTAG), Smoking status, White blood cell count, Lung function (FEV1/FVC)</p>	
3	123433220	123889576	123689576 (rs55935332)	3.47E-12	ENSG00000065534, ENSG00000175455, ENSG00000065371, ENSG00000160145	MYLK, CCDC14, ROPN1 , KALRN	<p>ENSG00000065371: ENSG00000065534: Adolescent idiopathic scoliosis, General cognitive ability, Cognitive empathy, Educational attainment (MTAG), Educational attainment (years of education); ENSG00000175455: Intelligence (MTAG), Drug-induced liver injury (fluoroquinolones), General cognitive ability; ENSG00000065371: ; ENSG00000160145: Mean platelet volume, Platelet count, Post bronchodilator FEV1/FVC ratio, Plateletcrit, Schizophrenia, Systolic blood pressure x dichotomous lifestyle risk score interaction (1df test), Systolic blood pressure x dichotomous lifestyle risk score interaction (2df test), Intraocular pressure, Moderate-to-late spontaneous preterm birth, Heschl's gyrus morphology, PR interval in Tripanosoma cruzi seropositivity, Thyroid peroxidase antibody levels, Aspartate</p>	

								aminotransferase levels in low alcohol consumption, Hematocrit, Platelet distribution width, Amyotrophic lateral sclerosis, Acute graft versus host disease in bone marrow transplantation (donor effect), Childhood steroid-sensitive nephrotic syndrome, Diffuse large B-cell lymphoma or systemic lupus erythematosus, Marginal zone lymphoma or systemic lupus erythematosus, Longevity (age >99th survival percentile), Red blood cell count, Benign childhood epilepsy with centro-temporal spikes, Lymphocyte counts, Childhood ALL/LBL (acute lymphoblastic leukemia/lymphoblastic lymphoma) treatment-related venous thromboembolism, Mean corpuscular hemoglobin concentration, General cognitive ability, Monocyte count, Urine pH measurement, Low urine pH, Educational attainment (MTAG), Cognitive performance (MTAG), Educational attainment (years of education), Red cell distribution width
10	78689487	79110042	78889487	1.62E-11	ENSG00000156113	KCNMA1	ENSG00000156113: Glucocorticoid-induced osteonecrosis, Lean body mass, Angioedema in response to angiotensin-converting enzyme inhibitor and/or angiotensin receptor blocker, Refractive error, Spherical equivalent, Myopia, Hypospadias, Obesity, Myopia (age of diagnosis), Non-melanoma skin cancer, Spherical equivalent or myopia (age of diagnosis), Male-pattern baldness, Initial pursuit acceleration, Educational attainment, Blood pressure, Mumps, Estimated glomerular filtration rate, Body mass index, Response to ranibizumab in age-related macular degeneration	ENSG00000156113: Glucocorticoid-induced osteonecrosis, Lean body mass, Angioedema in response to angiotensin-converting enzyme inhibitor and/or angiotensin receptor blocker, Refractive error, Spherical equivalent, Myopia, Hypospadias, Obesity, Myopia (age of diagnosis), Non-melanoma skin cancer, Spherical equivalent or myopia (age of diagnosis), Male-pattern baldness, Initial pursuit acceleration, Educational attainment, Blood pressure, Mumps, Estimated glomerular filtration rate, Body mass index, Response to ranibizumab in age-related macular degeneration

							<p>pursuit acceleration, Educational attainment, Blood pressure, Mumps, Estimated glomerular filtration rate, Body mass index, Response to ranibizumab in age-related macular degeneration (exudative), Smoking cessation in chronic obstructive pulmonary disease, DNA methylation variation (age effect), Heart rate in heart failure with reduced ejection fraction, Balding type 1, Height</p>	<p>(exudative), Smoking cessation in chronic obstructive pulmonary disease, DNA methylation variation (age effect), Heart rate in heart failure with reduced ejection fraction, Balding type 1, Height</p>
1	234142067	234549596	234342067	3.07E-11	ENSG00000183780, ENSG00000168275, ENSG00000059588	SLC35F3, COA6, TARBP1	<p>ENSG00000183780: Post bronchodilator FEV1/FVC ratio, Epstein-Barr virus copy number in lymphoblastoid cell lines, Trunk fat mass, Creatinine levels, Metabolite levels, Intracranial aneurysm, Pediatric bone mineral content (spine), Interleukin-6 levels, Chronic obstructive pulmonary disease or high blood pressure (pleiotropy), Adolescent idiopathic scoliosis, Diverticular disease</p>	<p>ENSG00000183780: Post bronchodilator FEV1/FVC ratio, Epstein-Barr virus copy number in lymphoblastoid cell lines, Trunk fat mass, Creatinine levels, Metabolite levels, Intracranial aneurysm, Pediatric bone mineral content (spine), Interleukin-6 levels, Chronic obstructive pulmonary disease or high blood pressure (pleiotropy), Adolescent idiopathic scoliosis, Diverticular disease; ENSG00000168275: ; ENSG00000059588: Cognitive test performance</p>
17	44595190	45055683	44802774	3.52E-11	ENSG00000238083, ENSG00000185829, ENSG00000073969, ENSG00000108379, ENSG00000158955,	LRR37A2, ARL17A, NSF, WNT3, WNT9B, GOSR	<p>ENSG00000073969: Ovarian cancer in BRCA1 mutation carriers, Parkinson's disease, Sense of smell, Intelligence</p>	<p>ENSG00000238083: ; ENSG00000185829: Hemoglobin levels, Reaction time, Red cell distribution width, Handedness (Left-handed vs. non-left-handed), Handedness (Right-handed vs. non-right-handed), Mean corpuscular hemoglobin concentration, General cognitive ability, Monocyte count; ENSG00000073969: Ovarian cancer in BRCA1 mutation carriers, Parkinson's disease, Sense of smell, Intelligence</p>

					ENSG00000108433, ENSG00000179673	2, RPRM L	<p>Cortical surface area (global PC1), Parkinson's disease or first degree relation to individual with Parkinson's disease, Brain region volumes, Hemoglobin levels, Intelligence, Male-pattern baldness, Worry, White matter microstructure (axial diuivities), Reaction time, White matter microstructure (mean diuivities), Epithelial ovarian cancer, Thyroid stimulating hormone levels, Cortical surface area, General factor of neuroticism, Smoking initiation (ever regular vs never regular), White matter microstructure (radial diuivities), White matter microstructure (fractional anisotropy), General cognitive ability, Macular thickness, Balding type 1</p>	<p>(MTAG), Neuroticism, Neurociticism, Feeling miserable, Experiencing mood swings, Feeling hurt, Feeling fed-up, Feeling nervous, Feeling worry, Cortical surface area (global PC1), Parkinson's disease or first degree relation to individual with Parkinson's disease, Brain region volumes, Hemoglobin levels, Intelligence, Male-pattern baldness, Worry, White matter microstructure (axial diuivities), Reaction time, White matter microstructure (mean diuivities), Epithelial ovarian cancer, Thyroid stimulating hormone levels, Cortical surface area, General factor of neuroticism, Smoking initiation (ever regular vs never regular), White matter microstructure (radial diuivities), White matter microstructure (fractional anisotropy), General cognitive ability, Macular thickness, Balding type 1; ENSG00000108379: Parkinson's disease, Celiac disease, Post bronchodilator FEV1, Hematocrit, Intelligence (MTAG), Coronary artery disease, Irritable mood, Neuroticism, Feeling guilty, Experiencing mood swings, Feeling hurt, Cortical surface area (global PC1), Parkinson's disease or first degree relation to individual with Parkinson's disease, Red blood cell count, Multiple system atrophy, Cognitive function, Alzheimer's disease in APOE e4-carriers, Hemoglobin levels, Depressed affect, Male-pattern baldness, White matter microstructure (axial diuivities), Hemoglobin concentration, Breast cancer, Itch intensity from mosquito bite adjusted by bite size, Handedness (non-right-handed vs right-handed), Handedness (Left-handed vs. non-left-handed), Handedness (left-handed vs. right-handed),</p>
--	--	--	--	--	----------------------------------	-----------	--	--

							Alcohol consumption (drinks per week), Reaction time, General factor of neuroticism, Atrial fibrillation, Intracranial volume, Lung function (FEV1), White matter microstructure (radial diffusivities), Lung function (FVC), Waist-to-hip ratio adjusted for BMI, General cognitive ability, Snoring, White matter microstructure (fractional anisotropy), Smoking initiation (ever regular vs never regular) (MTAG), Cognitive performance (MTAG); ENSG00000158955: Antineutrophil cytoplasmic antibody-associated vasculitis, Intraocular pressure, Mean corpuscular hemoglobin; ENSG00000108433: Blood pressure, Systolic blood pressure, Nonsyndromic cleft lip with cleft palate, Mean arterial pressure, Coronary artery disease, Blood urea nitrogen levels, QRS duration, QRS complex (Cornell), Pulse pressure, Hemoglobin levels, Myocardial infarction, Aortic root size, Orofacial clefts, Cleft lip with or without cleft palate, Idiopathic pulmonary fibrosis, Atrial fibrillation, Myocardial fractal dimension (slice 2), Myocardial fractal dimension (slice 3), Myocardial fractal dimension (slice 4), Medication use (agents acting on the renin-angiotensin system), Medication use (calcium channel blockers), Height, Cardiovascular disease; ENSG00000179673:	
10	90592757	91009553	90801209	4.00E-11	ENSG00000152766, ENSG00000138134, ENSG00000107796, ENSG00000026103,	ANKR D22, STAM BPL1, ACTA2 , FAS, CH25	ENSG00000026103: Immunoglobulin A, Chronic lymphocytic leukemia, Mosquito bite size, Blood protein levels, Blood protein levels in cardiovascular risk, Juvenile idiopathic arthritis (oligoarticular or rheumatoid	ENSG00000152766: Mucinous adenocarcinoma in colorectal cancer, Late-onset Alzheimer's disease, Accelerometer-based physical activity measurement (average acceleration); ENSG00000138134: Lung cancer, Pulse pressure, Brain region volumes; ENSG00000107796: Lung cancer, Pulse pressure, Chronic inflammatory

					ENSG00000138135, ENSG00000107798	H, LIPA	factor-negative polyarticular), Ankylosing spondylitis, Mean corpuscular volume, Red blood cell count	diseases (ankylosing spondylitis, Crohn's disease, psoriasis, primary sclerosing cholangitis, ulcerative colitis) (pleiotropy), Chronic lymphocytic leukemia; ENSG00000026103: Immunoglobulin A, Chronic lymphocytic leukemia, Mosquito bite size, Blood protein levels, Blood protein levels in cardiovascular risk, Juvenile idiopathic arthritis (oligoarticular or rheumatoid factor-negative polyarticular), Ankylosing spondylitis, Mean corpuscular volume, Red blood cell count; ENSG00000138135; ; ENSG00000107798: Coronary heart disease, Fibrinogen levels, Coronary artery disease (myocardial infarction, percutaneous transluminal coronary angioplasty, coronary artery bypass grafting, angina or chronic ischemic heart disease), Coronary artery disease, Neutrophil count, Blood protein levels, C-reactive protein levels, Sum neutrophil eosinophil counts, Myocardial infarction, Itch intensity from mosquito bite adjusted by bite size, White blood cell count, Red cell distribution width
2	241709643	242343441	242087712	8.54E-11	ENSG00000130294, ENSG00000172482, ENSG00000172478, ENSG00000162804, ENSG00000122085, ENSG00000115687, ENSG00000115685,	KIF1A, AGXT, C2orf54, SNED1, MTERFD2, PASK, PPP1R7 , ANO7, HDLB P,	ENSG00000115685:	ENSG00000130294: Nicotine withdrawal symptom count, Response to placebo treatment in childhood asthma (FVC change), Waist circumference adjusted for body mass index, Self-reported math ability, Height; ENSG00000172482: Blood metabolite levels, Height; ENSG00000172478: Major depressive disorder; ENSG00000162804: Sex hormone-binding globulin levels, Blood protein levels, Growth-regulated protein alpha levels, Pharmacokinetics of antipsychotic drugs in severe mental disorder (concentration drug ratio),

					ENSG00000146205, ENSG00000115677, ENSG00000168385, ENSG00000006607	SEPT2, FARP2		Toxicity response to radiotherapy in prostate cancer (hematuria) (time to event), Height; ENSG00000122085: Sex hormone-binding globulin levels, Blood protein levels, Pharmacokinetics of antipsychotic drugs in severe mental disorder (concentration drug ratio); ENSG00000115687: Height; ENSG00000115685: ; ENSG00000146205: Reaction time; ENSG00000115677: Chronic lymphocytic leukemia, Fibrinogen levels, HDL cholesterol levels, Apolipoprotein B levels, Apolipoprotein A1 levels, Male-pattern baldness, Waist circumference adjusted for BMI (adjusted for smoking behaviour), Waist circumference adjusted for BMI (joint analysis main effects and smoking interaction), Waist circumference adjusted for BMI in non-smokers, Waist circumference adjusted for body mass index, Intraocular pressure, Height, Balding type 1, Eosinophil counts; ENSG00000168385: Height, Cerebrospinal fluid immune biomarker levels, Male-pattern baldness, Balding type 1, Lung function (FEV1/FVC); ENSG00000006607: Chronic lymphocytic leukemia, Prostate cancer, Vitiligo, Triglyceride levels, Cerebrospinal fluid immune biomarker levels, Disability (impaired activities of daily living), C-reactive protein levels, Fibrinogen levels, Fibrinogen, Red blood cell count, Low density lipoprotein cholesterol levels, Systolic blood pressure, Educational attainment (MTAG), Educational attainment (years of education), White blood cell count
1 1	4463476 4	45073 989	44838 496	1.42E-10	ENSG00000085117, ENSG000001	CD82, TSPA N18,	ENSG00000157570: Schizophrenia, Obstructive sleep apnea trait (average	ENSG00000085117: Hemostatic factors and hematological phenotypes, Red cell distribution width, White matter

					57570, ENSG000001 75274	TP53I 11	respiratory event duration), Intraocular pressure, Cortical brain region measurements (area, volume and thickness)	microstructure (axial diuivities), White matter microstructure (mean diuivities), White matter microstructure (radial diuivities), Lymphocyte counts, DNA methylation variation (age effect); ENSG00000157570: Schizophrenia, Obstructive sleep apnea trait (average respiratory event duration), Intraocular pressure, Cortical brain region measurements (area, volume and thickness); ENSG00000175274:
7	9594795 9	96347 959	96147 959	2.23E- 10	ENSG000000 04864, ENSG000001 97851, ENSG000001 27922	SLC25 A13, C7orf 76, SHFM 1	ENSG00000197851:	ENSG00000004864: Height, Pork consumption, Oily fish consumption; ENSG00000197851: ; ENSG00000127922: Bone mineral density (hip), Bone mineral density (spine), Femoral neck bone mineral density, Total body bone mineral density, Heel bone mineral density, Total body bone mineral density (age 45-60), Total body bone mineral density (age over 60), Serum platinum levels after completion of cisplatin chemotherapy, Bone mineral density, Nicotine glucouronidation, Chin dimples, Lumbar spine bone mineral density, Bone ultrasound measurement (velocity of sound), Fractures, Cortical surface area, Facial morphology traits (63 three-dimensional facial segments)
1 5	9338509 6	93786 062	93586 062	3.83E- 10	ENSG000001 73575, ENSG000001 82175	CHD2, RGMA	ENSG00000182175: HDL cholesterol levels x short total sleep time interaction (2df test), Blood protein levels, Heel bone mineral density	ENSG00000173575: IgG glycosylation, Schizophrenia, General cognitive ability, Pulse pressure, Mean corpuscular volume, Cognitive performance (MTAG), Self- reported math ability, Self-reported math ability (MTAG), Educational attainment (years of education), Educational attainment (MTAG), Highest math class taken, Highest math class taken (MTAG), Mean corpuscular hemoglobin, Menarche (age at onset); ENSG00000182175: HDL cholesterol levels x short total sleep time

								interaction (2df test), Blood protein levels, Heel bone mineral density
1 3	8817053 2	88668 426	88411 593	4.40E- 10	ENSG000001 65300	SLITRK 5	ENSG00000165300:	ENSG00000165300:
5	3172542 8	32125 428	31925 428	4.77E- 10	ENSG000001 33401, ENSG000001 13384	PDZD 2, GOLP H3	ENSG00000133401: Obesity-related traits, Myocardial infarction, Height, Chronic obstructive pulmonary disease, Vertical cup-disc ratio (adjusted for vertical disc diameter), Vertical cup-disc ratio (multi-trait analysis), Eotaxin levels, Renal cell carcinoma, Bipolar disorder (body mass index interaction), Response to serotonin reuptake inhibitors in major depressive disorder (plasma drug and metabolite levels), Optic disc size, Vertical cup-disc ratio, Adolescent idiopathic scoliosis, Interleukin-7 levels, Colorectal cancer, Metabolite levels, Corpus callosum central volume, Working memory, Breast cancer specific mortality in estrogen receptor positive breast cancer, Self-reported childhood asthma in adult smokers, Heart rate in heart failure with reduced ejection fraction	ENSG00000133401: Obesity-related traits, Myocardial infarction, Height, Chronic obstructive pulmonary disease, Vertical cup-disc ratio (adjusted for vertical disc diameter), Vertical cup-disc ratio (multi-trait analysis), Eotaxin levels, Renal cell carcinoma, Bipolar disorder (body mass index interaction), Response to serotonin reuptake inhibitors in major depressive disorder (plasma drug and metabolite levels), Optic disc size, Vertical cup-disc ratio, Adolescent idiopathic scoliosis, Interleukin-7 levels, Colorectal cancer, Metabolite levels, Corpus callosum central volume, Working memory, Breast cancer specific mortality in estrogen receptor positive breast cancer, Self-reported childhood asthma in adult smokers, Heart rate in heart failure with reduced ejection fraction; ENSG00000113384: Height
2	2114982 06	21206 1871	21186 1871	5.45E- 10	ENSG000000 21826	CPS1	ENSG00000021826: Chronic kidney disease, Fibrinogen, Body mass index in asthmatics, Homocysteine levels, Metabolite levels, Glomerular filtration rate	ENSG00000021826: Chronic kidney disease, Fibrinogen, Body mass index in asthmatics, Homocysteine levels, Metabolite levels, Glomerular filtration rate (creatinine), Betaine levels in individuals undergoing cardiac evaluation,

						<p>(creatinine), Betaine levels in individuals undergoing cardiac evaluation, Glomerular filtration rate in non diabetics (creatinine), Body mass index, Eosinophil percentage of white cells, Mean corpuscular volume, Platelet count, Macular telangiectasia type 2, Metabolite levels (small molecules and protein measures), Plateletcrit, Amino acid levels, Fibrinogen levels, Mean corpuscular hemoglobin, Creatinine levels, Alanine transaminase levels, Serum metabolite concentrations in chronic kidney disease, Urinary metabolite levels in chronic kidney disease, Urinary metabolite modules (eigenmetabolites) in chronic kidney disease, Fat-free mass, HDL cholesterol levels, Appendicular lean mass, Apolipoprotein A1 levels, HDL cholesterol, Plasma homocysteine levels (post-methionine load test), Blood metabolite levels, Serum metabolite levels, Plasma free amino acid levels (adjusted for twenty other PFAAs), Serum 25-Hydroxyvitamin D levels, Eosinophil counts, HDL cholesterol levels x alcohol consumption (regular vs non-regular drinkers) interaction</p>	<p>Glomerular filtration rate in non diabetics (creatinine), Body mass index, Eosinophil percentage of white cells, Mean corpuscular volume, Platelet count, Macular telangiectasia type 2, Metabolite levels (small molecules and protein measures), Plateletcrit, Amino acid levels, Fibrinogen levels, Mean corpuscular hemoglobin, Creatinine levels, Alanine transaminase levels, Serum metabolite concentrations in chronic kidney disease, Urinary metabolite levels in chronic kidney disease, Urinary metabolite modules (eigenmetabolites) in chronic kidney disease, Fat-free mass, HDL cholesterol levels, Appendicular lean mass, Apolipoprotein A1 levels, HDL cholesterol, Plasma homocysteine levels (post-methionine load test), Blood metabolite levels, Serum metabolite levels, Plasma free amino acid levels (adjusted for twenty other PFAAs), Serum 25-Hydroxyvitamin D levels, Eosinophil counts, HDL cholesterol levels x alcohol consumption (regular vs non-regular drinkers) interaction (2df), Eosinophil percentage of granulocytes, Mean platelet volume, Urinary metabolites, Glomerular filtration rate, Estimated glomerular filtration rate, Blood urea nitrogen levels, Estimated glomerular filtration rate in diabetes, Estimated glomerular filtration rate in non-diabetics, HDL cholesterol levels in current drinkers, HDL cholesterol levels x alcohol consumption (drinkers vs non-drinkers) interaction (2df), Blood protein levels, Urinary albumin-to-creatinine ratio, Red blood cell count, Red cell distribution width, White blood cell count, Neutrophil count, Urinary albumin excretion (no hypertensive medication), Urinary albumin</p>
--	--	--	--	--	--	---	---

							(2df), Eosinophil percentage of granulocytes, Mean platelet volume, Urinary metabolites, Glomerular filtration rate, Estimated glomerular filtration rate, Blood urea nitrogen levels, Estimated glomerular filtration rate in diabetes, Estimated glomerular filtration rate in non-diabetics, HDL cholesterol levels in current drinkers, HDL cholesterol levels x alcohol consumption (drinkers vs non-drinkers) interaction (2df), Blood protein levels, Urinary albumin-to-creatinine ratio, Red blood cell count, Red cell distribution width, White blood cell count, Neutrophil count, Urinary albumin excretion (no hypertensive medication), Urinary albumin excretion, Glycine levels, Lymphocyte counts, Systolic blood pressure, Urinary potassium to creatinine ratio, Urinary sodium to creatinine ratio, Serum uric acid levels, Urate levels, Height	excretion, Glycine levels, Lymphocyte counts, Systolic blood pressure, Urinary potassium to creatinine ratio, Urinary sodium to creatinine ratio, Serum uric acid levels, Urate levels, Height
1 2	1335668 01	13396 6801	13376 6801	6.71E- 10	ENSG000001 98393, ENSG000001 98040, ENSG000001 96387, ENSG000002 14029,	ZNF26 , ZNF84 , ZNF14 0, ZNF89 1,	ENSG00000090612: Waist-hip ratio, Waist-to-hip ratio adjusted for BMI, Type 2 diabetes, Cardiovascular disease	ENSG00000198393: ; ENSG00000198040: Heparin-induced thrombocytopenia; ENSG00000196387: ; ENSG00000214029: Refractive error, Waist-to-hip ratio adjusted for BMI, Intraocular pressure; ENSG00000256223: Type 2 diabetes; ENSG00000090612: Waist-hip ratio, Waist-to-hip ratio adjusted for BMI, Type 2

					ENSG00000256223, ENSG00000090612, ENSG00000227059	ZNF10 , ZNF268 , ANHX		diabetes, Cardiovascular disease; ENSG00000227059:
9	113454364	113913765	113713765	7.08E-10	ENSG00000030304, ENSG00000198121	MUSK , LPAR1	ENSG00000198121: Corneal structure, Post bronchodilator FEV1 in COPD, Eosinophil percentage of white cells, Eosinophil counts, Eosinophil percentage of granulocytes, Metabolite levels, Central corneal thickness, Neutrophil percentage of granulocytes, Sum eosinophil basophil counts, Pursuit maintenance gain, Height, Colorectal cancer or advanced adenoma	ENSG0000030304: Heel bone mineral density, Body mass index, Plasma factor V levels in venous thrombosis (conditioned on rs6027), Maximum stenosis, Mean degree of stenosis, Height; ENSG00000198121: Corneal structure, Post bronchodilator FEV1 in COPD, Eosinophil percentage of white cells, Eosinophil counts, Eosinophil percentage of granulocytes, Metabolite levels, Central corneal thickness, Neutrophil percentage of granulocytes, Sum eosinophil basophil counts, Pursuit maintenance gain, Height, Colorectal cancer or advanced adenoma
4	38047254	38736155	38260451	1.39E-09	ENSG00000065882, ENSG00000109787	TBC1D1, KLF3	ENSG00000065882: Amyotrophic lateral sclerosis in C9orf72 mutation negative individuals, Periodontal microbiota, Facial morphology (factor 21, depth of nasal alae), Neutrophil percentage of white cells, Lymphocyte percentage of white cells, Metabolite levels, Verbal declarative memory, Reaction time, Lymphocyte counts, Weight, Response to SSRI (symptom remission), Response to antidepressants (symptom remission), White blood cell count, Eosinophil counts	ENSG00000065882: Amyotrophic lateral sclerosis in C9orf72 mutation negative individuals, Periodontal microbiota, Facial morphology (factor 21, depth of nasal alae), Neutrophil percentage of white cells, Lymphocyte percentage of white cells, Metabolite levels, Verbal declarative memory, Reaction time, Lymphocyte counts, Weight, Response to SSRI (symptom remission), Response to antidepressants (symptom remission), White blood cell count, Eosinophil counts; ENSG00000109787: Eosinophil percentage of white cells, Eosinophil counts, White blood cell count, Sum eosinophil basophil counts, Eosinophil percentage of granulocytes, Neutrophil percentage of granulocytes, Lymphocyte counts, Body mass index, Hand grip strength, Mean

								platelet volume, Mean corpuscular hemoglobin, Red cell distribution width
6	170396266	170796266	170596266	1.54E-09	ENSG00000198719, ENSG00000112584	DLL1, FAM120B	ENSG00000198719: General risk tolerance (MTAG)	ENSG00000198719: General risk tolerance (MTAG); ENSG00000112584: General risk tolerance (MTAG), Idiopathic dilated cardiomyopathy, Paracentral lobule volume, Menarche (age at onset)
10	94978044	95389525	95189525	1.57E-09	ENSG00000138119, ENSG00000138180, ENSG00000186188, ENSG00000138207, ENSG00000095464	MYOF, CEP55, FFAR4, RBP4, PDE6C	ENSG00000138119: Gut microbiome composition (summer), Facial morphology (factor 17, height of vermilion upper lip), Cerebrospinal fluid immune biomarker levels; ENSG00000138180: Lobe attachment (rater-scored or self-reported), Height; ENSG00000186188: Retinol levels, Optic disc area, Waist-to-hip ratio adjusted for BMI, Blood protein levels, Waist-hip ratio, White blood cell count; ENSG00000138207: Optic disc area, Blood protein levels; ENSG00000095464: Preschool internalizing problems, Urinary tract infection frequency	
1	168973787	169701700	169450264	1.71E-09	ENSG00000143153, ENSG00000143156, ENSG00000117475, ENSG00000117477, ENSG00000117479, ENSG00000198734, ENSG00000174175, ENSG00000000460, ENSG00000188404,	ATP1B1, NME7, BLZF1, CCDC181, SLC19A2, F5, SELP, C1orf112, SELL, SELE	ENSG00000117479: QT interval	ENSG00000143153: QT interval, Coronary artery disease, Venous thromboembolism, Pulse pressure, Electrocardiographic traits, Systolic blood pressure; ENSG00000143156: Venous thromboembolism, D-dimer levels, Coronary artery disease, QT interval, Pulse pressure, Systolic blood pressure, Mumps, QT dynamics during exercise; ENSG00000117475: Blood protein levels; ENSG00000117477: ; ENSG00000117479: QT interval; ENSG00000198734: Hippocampal atrophy, Activated partial thromboplastin time, Venous thromboembolism, Hemostatic factors and hematological phenotypes, Uric acid levels, Inflammatory bowel disease, Thrombosis, Ischemic stroke, Optic disc area, Blood

					ENSG0000007908			protein levels, Prothrombin time, Vertical cup-disc ratio (multi-trait analysis), Cytokine network levels (multivariate analysis), CTACK levels, Optic disc size, Vertical cup-disc ratio, Peripheral artery disease, Stem cell factor levels, Medication use (antithrombotic agents); ENSG00000174175: Soluble levels of adhesion molecules, Activated partial thromboplastin time, Blood protein levels, Optic disc size, Late-onset Alzheimer's disease; ENSG00000000460: Venous thromboembolism, Acne (severe), Blood protein levels, Intrinsic epigenetic age acceleration, Amyotrophic lateral sclerosis, Tonsillectomy, White blood cell count, Monocyte count, Cardiac Troponin-T levels, Age at menopause, Eosinophil counts; ENSG00000188404: Blood protein levels, Amyotrophic lateral sclerosis, Monocyte count, Eosinophil counts; ENSG00000007908: White blood cell count, Blood protein levels, Age at menopause
--	--	--	--	--	----------------	--	--	--

4599
4600
4601
4602
4603

Table S4d.1 Eurasia k3

chr	start	end	bestpos	Min(P)	ensembl	hgnc	disease_trait_best	disease_trait
6	134192815	134628278	134392815	8.49E-16	ENSG00000118526, ENSG00000028839, ENSG00000146411, ENSG00000118515	TCF21, TBPL1, SLC2A12 , SGK1	ENSG00000146411: Coronary artery disease, High chromosomal aberration frequency (chromosome type), FEV1, Lung function (FVC). Mean corpuscular haemoglobin.	ENSG00000118526: Coronary heart disease, Coronary artery disease or ischemic stroke, Coronary artery disease, Coronary artery disease or large artery stroke, PR interval, Medication use (diuretics), Lung function (FEV1/FVC); ENSG00000028839: Mean corpuscular hemoglobin; ENSG00000146411: Coronary artery disease, High chromosomal aberration frequency (chromosome type), FEV1, Lung function (FVC); ENSG00000118515: Immune response to smallpox (secreted IFN-alpha), Alzheimer disease and age of onset, Pelvic organ

								prolapse, Pelvic organ prolapse (moderate/severe), Blond vs. brown/black hair color, Schizophrenia (inflammation and infection response interaction), Body mass index (smoking years interaction), Metabolite levels, Adolescent idiopathic scoliosis, Hair color
15	48226484	48633494	48433494	2.09E-15	ENSG00000188467, ENSG00000104177, ENSG00000233932, ENSG00000074803, ENSG00000128951	SLC24A5, MYEF2, CTXN2, SLC12A1, DUT	ENSG00000188467: Hair color, Skin pigmentation, Body mass index, Skin, hair and eye pigmentation (multivariate analysis), Eye color, Eye color (brightness), Eye color (saturation), Iris color (a* coordinate), Skin reflectance (Melanin index), Iris color (b* coordinate); ENSG00000104177: Skin pigmentation, Hair color; ENSG00000233932: ; ENSG00000074803: Systemic lupus erythematosus, Longevity, Skin reflectance (Melanin index), Lymphocyte counts; ENSG00000128951: Protein quantitative trait loci	ENSG00000188467: Hair color, Skin pigmentation, Body mass index, Skin, hair and eye pigmentation (multivariate analysis), Eye color, Eye color (brightness), Eye color (saturation), Iris color (a* coordinate), Skin reflectance (Melanin index), Iris color (b* coordinate); ENSG00000104177: Skin pigmentation, Hair color; ENSG00000233932: ; ENSG00000074803: Systemic lupus erythematosus, Longevity, Skin reflectance (Melanin index), Lymphocyte counts; ENSG00000128951: Protein quantitative trait loci
2	135084038	137229668	135479980	5.49E-11	ENSG00000152127, ENSG00000152128, ENSG00000153086, ENSG00000082258, ENSG00000176601, ENSG00000115839, ENSG00000121988, ENSG00000048991, ENSG00000144224, ENSG00000115850, ENSG00000076003, ENSG00000115866, ENSG00000121966	MGAT5, TMEM163, ACMSD, CCNT2, MAP3K19, RAB3GAP1, ZRANB3, R3HDM1, UBXN4, LCT, MCM6, DARS, CXCR4	ENSG00000152128: Large artery stroke, Neuroticism, Parkinson's disease or first degree relation to individual with Parkinson's disease, HDL cholesterol levels, Cutaneous melanoma or hair colour, Asthma x air pollution interaction (2df), Hematocrit, Blond vs. brown/black hair color, Spatial memory, Hemoglobin concentration, Low density lipoprotein cholesterol levels, Self-reported math ability, Self-reported math ability (MTAG), Red blood cell count, Hair color	ENSG00000152127: Multiple sclerosis (severity), Subcutaneous adipose tissue, Post bronchodilator FEV1/FVC ratio, Serum alkaline phosphatase levels, N-glycan levels, Blood protein levels, Chronic lymphocytic leukemia or systemic lupus erythematosus, Marginal zone lymphoma or systemic lupus erythematosus, Systemic lupus erythematosus, Eosinophil counts; ENSG00000152128: Large artery stroke, Neuroticism, Parkinson's disease or first degree relation to individual with Parkinson's disease, HDL cholesterol levels, Cutaneous melanoma or hair colour, Asthma x air pollution interaction (2df), Hematocrit, Blond vs. brown/black hair color, Spatial memory, Hemoglobin concentration, Low density lipoprotein cholesterol levels, Self-reported math ability, Self-reported math ability (MTAG), Red blood cell count, Hair color; ENSG00000153086: Obesity-related traits, LDL cholesterol levels, Apolipoprotein B levels, Diisocyanate-induced asthma, Blood metabolite levels, Hematocrit, Free thyroxine concentration, Red blood cell count, Hand grip strength, Diastolic blood pressure; ENSG00000082258: Age at menopause; ENSG00000176601: Colonoscopy-negative controls vs population controls, Corneal structure, Hematocrit, Hemoglobin concentration, Mean corpuscular hemoglobin; ENSG00000115839: Cholesterol, total, Body mass index, Blood metabolite levels, HDL cholesterol levels x alcohol consumption (regular vs non-regular drinkers) interaction

								(2df), LDL cholesterol levels in current drinkers, LDL cholesterol levels x alcohol consumption (drinkers vs non-drinkers) interaction (2df), LDL cholesterol levels x alcohol consumption (regular vs non-regular drinkers) interaction (2df), LDL cholesterol levels, Gut microbiota (bacterial taxa, rank normal transformation method), Sudden cardiac arrest in coronary artery disease; ENSG00000121988: Mosquito bite size, Low density lipoprotein cholesterol levels, Hip circumference, Waist circumference adjusted for body mass index, Sudden cardiac arrest in coronary artery disease, Type 2 diabetes, Height; ENSG00000048991: Mosquito bite size, Blood protein levels, Urinary metabolite levels in chronic kidney disease, Height, LDL cholesterol levels x short total sleep time interaction (2df test), HDL cholesterol levels, Apolipoprotein A1 levels, Red cell distribution width, Hand grip strength; ENSG00000144224: Cholesterol, total, Corneal structure; ENSG00000115850: White blood cell count, Blood protein levels, Osteoarthritis (self-reported); ENSG00000076003: Body mass index, Total cholesterol change in response to fenofibrate in statin-treated type 2 diabetes, Blood protein levels, Hip circumference, 1,5-anhydroglucitol levels, Gut microbiota (bacterial taxa, rank normal transformation method); ENSG00000115866: Mosquito bite size, White blood cell count, Neutrophil count, Monocyte count, Total cholesterol change in response to fenofibrate in statin-treated type 2 diabetes, Systemic lupus erythematosus; ENSG00000121966: Tonsillectomy
5	33699690	34164938	33958959	2.30E-10	ENSG00000151388, ENSG00000182631, ENSG00000164175, ENSG00000242110, ENSG00000082196	ADAMTS12, RXFP3, SLC45A2 , AMACR, C1QTNF3	ENSG00000164175: Tanning, Skin pigmentation, Black vs. blond hair color, Black vs. red hair color, Melanoma, Hair color, Eye color, Squamous cell carcinoma, Skin colour saturation, Perceived skin darkness, Skin sensitivity to sun, Black vs. non-black hair color, Skin aging (microtopography measurement), Basal cell carcinoma, Cutaneous melanoma or hair colour, Cutaneous malignant melanoma, Nevus	ENSG00000151388: Stroke (pediatric), Mortality in heart failure, Aspartate aminotransferase levels in excessive alcohol consumption, Neurofibrillary tangles, Adolescent idiopathic scoliosis, Height, Hair color; ENSG00000182631: ; ENSG00000164175: Tanning, Skin pigmentation, Black vs. blond hair color, Black vs. red hair color, Melanoma, Hair color, Eye color, Squamous cell carcinoma, Skin colour saturation, Perceived skin darkness, Skin sensitivity to sun, Black vs. non-black hair color, Skin aging (microtopography measurement), Basal cell carcinoma, Cutaneous melanoma or hair colour, Cutaneous malignant melanoma, Nevus count or cutaneous melanoma, Rosacea symptom severity, Low tan response, Skin, hair and eye pigmentation (multivariate

							count or cutaneous melanoma, Rosacea symptom severity, Low tan response, Skin, hair and eye pigmentation (multivariate analysis), Brown vs. black hair color, Blond vs. brown/black hair color, Cutaneous squamous cell carcinoma, Eye color (saturation), Eye color (brightness), Monobrow, Keratinocyte cancer (MTAG), Eye color traits, Skin pigmentation traits, Hair morphology traits, Sunburns	analysis), Brown vs. black hair color, Blond vs. brown/black hair color, Cutaneous squamous cell carcinoma, Eye color (saturation), Eye color (brightness), Monobrow, Keratinocyte cancer (MTAG), Eye color traits, Skin pigmentation traits, Hair morphology traits, Sunburns; ENSG00000242110: Blond vs. brown/black hair color, Longevity; ENSG0000082196: Waist-to-hip ratio adjusted for BMI, Waist-hip ratio
2	85369379	85885211	85592841	2.81E-10	ENSG00000152284, ENSG00000152291, ENSG00000042445, ENSG00000115459, ENSG00000042493, ENSG00000152292, ENSG00000168906, ENSG00000115486, ENSG00000118640, ENSG00000168899, ENSG00000168894, ENSG00000168890, ENSG00000168883, ENSG00000168887, ENSG00000168878	TCF7L1, TGOLN2, RETSAT, ELMOD3 , CAPG, SH2D6, MAT2A, GG CX, VAMP8 , VAMP5 , RNF181, TMEM150A, USP39, C2orf68, SFTP B	ENSG00000115459: Blood protein levels.	ENSG00000152284: Total body bone mineral density, Pulse pressure, Heel bone mineral density, Cervical cancer, Red blood cell count, Lung function (FEV1/FVC), Systolic blood pressure; ENSG00000152291: Ear protrusion, White blood cell count; ENSG00000042445: ; ENSG00000115459: Blood protein levels; ENSG00000042493: ; ENSG00000152292: Mean platelet volume, Platelet count; ENSG00000168906: Basophil count, Coronary artery disease; ENSG00000115486: Coronary artery disease (myocardial infarction, percutaneous transluminal coronary angioplasty, coronary artery bypass grafting, angina or chronic ischemic heart disease), Triglyceride levels, Coronary artery disease, Eosinophil counts; ENSG00000118640: Eosinophil counts, Prostate cancer, Fat-free mass, Sum eosinophil basophil counts, Coronary artery disease; ENSG00000168899: Parental longevity (father's age at death), Adolescent idiopathic scoliosis, Height; ENSG00000168894: ; ENSG00000168890: White blood cell count; ENSG00000168883: Neurofibrillary tangles; ENSG00000168887: ; ENSG00000168878: Blood protein levels
3	129487889	129894124	129694124	4.23E-10	ENSG00000172765, ENSG00000170893	TMCC1, TRH	ENSG00000170893: waist-hip circumference, Sitting height, Chronic lower respiratory diseases, asthma. rs10934899	ENSG00000172765: Height, Type 2 diabetes, Waist-to-hip ratio adjusted for BMI; ENSG00000170893:
4	38553198	38965720	38765720	2.95E-09	ENSG00000109787, ENSG00000174123, ENSG00000174125,	KLF3, TLR10 , TLR1, TLR6, FAM114A1	ENSG00000174123: Peripheral arterial disease (traffic-related air pollution interaction), Asthma or	ENSG00000109787: Eosinophil percentage of white cells, Eosinophil counts, White blood cell count, Sum eosinophil basophil counts, Eosinophil percentage of granulocytes, Neutrophil percentage of granulocytes, Lymphocyte counts,

					ENSG00000174130, ENSG00000197712		allergic disease (pleiotropy), Adolescent idiopathic scoliosis	Body mass index, Hand grip strength, Mean platelet volume, Mean corpuscular hemoglobin, Red cell distribution width; ENSG00000174123: Peripheral arterial disease (traffic-related air pollution interaction), Asthma or allergic disease (pleiotropy), Adolescent idiopathic scoliosis; ENSG00000174125: Coronary artery calcified atherosclerotic plaque score in type 2 diabetes, Limited cutaneous systemic scleroderma, Diabetes in response to antihypertensive drug treatment (treatment strategy interaction), Asthma, Allergic sensitization, Self-reported allergy, Asthma and hay fever, Alcohol consumption, Asthma (childhood onset), Allergic rhinitis, Hay fever and/or eczema, Allergy, Allergic disease (asthma, hay fever or eczema), Breast cancer, Asthma onset (childhood vs adult), Composite immunoglobulin trait (IgG/IgM), DNA methylation variation (age effect), Asthma (age of onset), Eczema, Respiratory diseases; ENSG00000174130: Allergic disease (asthma, hay fever or eczema); ENSG00000197712: Alcohol dependence, Moderate or severe diarrhoea in darapladib-treated cardiovascular disease (time to event), Red blood cell count
1	227020437	227877723	227407855	3.77E-09	ENSG00000143801, ENSG00000163050, ENSG00000143776, ENSG00000181450	PSEN2, ADCK3, CDC42BPA, ZNF678	ENSG00000143776: Optic disc area, Optic cup area, Blond vs. brown/black hair color, Metabolite levels, Diastolic blood pressure, Myeloid white cell count, Blood urea nitrogen levels, White blood cell count, Neutrophil count, Cardiovascular death or myocardial infarction in response to clopidogrel treatment, Smoking initiation (ever regular vs never regular) (MTAG), Hair color	ENSG00000143801: Worry, Heel bone mineral density; ENSG00000163050: Granulocyte percentage of myeloid white cells, Cataracts (operation), Neutrophil count, Lymphocyte percentage of white cells, Sum basophil neutrophil counts, Waist-to-hip ratio adjusted for BMI, Worry, Sum neutrophil eosinophil counts, Granulocyte count, Gut microbiota (bacterial taxa, hurdle binary method), PR interval, Waist circumference adjusted for body mass index, Estimated glomerular filtration rate, Highest math class taken; ENSG00000143776: Optic disc area, Optic cup area, Blond vs. brown/black hair color, Metabolite levels, Diastolic blood pressure, Myeloid white cell count, Blood urea nitrogen levels, White blood cell count, Neutrophil count, Cardiovascular death or myocardial infarction in response to clopidogrel treatment, Smoking initiation (ever regular vs never regular) (MTAG), Hair color; ENSG00000181450: Height, Hip circumference adjusted for BMI, Body fat distribution (arm fat ratio), Body fat distribution (leg fat ratio), Body fat distribution (trunk fat ratio), Insular cortex volume, Waist circumference adjusted for body mass index

21	26691192	27196935	26935513	1.27E-08	ENSG00000154719, ENSG00000154721, ENSG00000154723, ENSG00000154727	MRPL39, JAM2, ATP5J, GABPA	ENSG00000154719: pork intake, arthrosis, Comparative height size at 10, Abdominal hernia	ENSG00000154719 ; ENSG00000154721: Longitudinal change in brain amyloid TP burden, Age at loss of ambulation in Duchenne muscular dystrophy; ENSG00000154723: Alzheimer's disease in hypertension-negative individuals; ENSG00000154727: Gout (normal type), Nicotine dependence symptom count
9	27009422	27434948	27209422	1.56E-08	ENSG00000096872, ENSG00000120156, ENSG00000120160, ENSG00000120162	IFT74, TEK, EQTN, MOB3B	ENSG00000120156: Comparative height at 10. Sitting height. IgG glycosylation, Coronary artery disease, Blood protein levels, Blood protein levels in cardiovascular risk, G, Trans fatty acid levels, Endothelial growth factor levels, Cognitive decline (age-related), Clostridium difficile infection in multiple myeloma.	ENSG00000096872: Emphysema annual change measurement in smokers (adjusted lung density); ENSG00000120156: IgG glycosylation, Coronary artery disease, Blood protein levels, Blood protein levels in cardiovascular risk, Schizophrenia, Trans fatty acid levels, Endothelial growth factor levels, Cognitive decline (age-related), Clostridium difficile infection in multiple myeloma; ENSG00000120160: ; ENSG00000120162: Response to TNF antagonist treatment, Urinary symptoms in response to radiotherapy in prostate cancer, Amyotrophic lateral sclerosis, Plasma trimethylamine N-oxide levels, Breast cancer specific mortality in estrogen receptor negative breast cancer, Height, Hair color
15	38464638	38992430	38792430	1.65E-08	ENSG00000166068, ENSG00000171262, ENSG00000172575, ENSG00000175779	SPRED1, FAM98B, RASGRP1, C15orf53	ENSG00000172575: Type 1 diabetes, Crohn's disease, Rheumatoid arthritis (ACPA-positive), Rheumatoid arthritis, Multiple sclerosis, Type 2 diabetes, Carboplatin disposition in epithelial ovarian cancer, Autoimmune thyroid disease, Medication use (thyroid preparations), Autoimmune traits (pleiotropy), Autoimmune traits, Hypothyroidism	ENSG00000166068: Birth weight, HDL cholesterol levels x long total sleep time interaction (2df test), Adolescent idiopathic scoliosis, Cognitive performance (processing speed); ENSG00000171262: Systemic lupus erythematosus; ENSG00000172575: Type 1 diabetes, Crohn's disease, Rheumatoid arthritis (ACPA-positive), Rheumatoid arthritis, Multiple sclerosis, Type 2 diabetes, Carboplatin disposition in epithelial ovarian cancer, Autoimmune thyroid disease, Medication use (thyroid preparations), Autoimmune traits (pleiotropy), Autoimmune traits, Hypothyroidism; ENSG00000175779: Cardiac hypertrophy, Metabolic traits, Bipolar disorder or major depressive disorder, Bipolar disorder, Bipolar disorder and schizophrenia, Metabolite levels, Developmental language disorder, Platelet count, Plateletcrit, Parental extreme longevity (95 years and older), Triglyceride change in response to fenofibrate in statin-treated type 2 diabetes, Bipolar I disorder, Alcohol dependence (age at onset), Response to abacavir-containing treatment in HIV-1 infection (virologic failure), Brain region volumes, Bone mineral density, Age at loss of ambulation in Duchenne muscular dystrophy, Follicle stimulating hormone levels in polycystic ovary syndrome, Chronic kidney disease, Estimated

							glomerular filtration rate, Estimated glomerular filtration rate in diabetes, Triglycerides, Circulating fibroblast growth factor 23 levels, Cortical surface area, Leukocyte telomere length, Lymphocyte counts, Non-alcoholic fatty liver disease activity score, Monocyte count, Systemic lupus erythematosus, anorexia nervosa, attention-deficit/hyperactivity disorder, autism spectrum disorder, bipolar disorder, major depression, obsessive-compulsive disorder, schizophrenia, or Tourette syndrome (pleiotropy), Lung function (FVC), Height
11	61264124	62097073	61604967	2.34E-08	ENSG00000204950, LRRC10B, SYT7, ENSG0000011347, DAGLA, ENSG00000134780, MYRF, ENSG00000124920, TMEM258, ENSG00000134825, FEN1, ENSG00000168496, FADS2, ENSG00000134824, FADS1, ENSG00000149485, FADS3, ENSG00000221968, RAB31L1, ENSG00000167994, BEST1, ENSG00000167995, FTH1, ENSG00000167996, INCENP, ENSG00000149503, SCGB1D1, ENSG00000168515, SCGB2A1, ENSG00000124939, SCGB1D2, ENSG00000124935, SCGB2A2, ENSG00000110484, SCGB1D4, ENSG00000197745	ENSG00000134824: HDL cholesterol, Triglycerides, Fasting blood glucose, Homeostasis model assessment of beta-cell function, Cholesterol, total, LDL cholesterol, Heart rate, Metabolic syndrome, Resting heart rate, Lipid metabolism phenotypes, Hematology traits, Liver enzyme levels (alkaline phosphatase), Comprehensive strength and appendicular lean mass, Metabolite levels, Response to statin therapy, Metabolic traits, Plasma omega-3 polyunsaturated fatty acid levels (docosapentaenoic acid), Plasma omega-3 polyunsaturated fatty acid levels (alpha-linolenic acid), Plasma omega-3 polyunsaturated fatty acid level (eicosapentaenoic acid), Platelet count, Fasting blood glucose (BMI interaction), Inflammatory bowel disease, Plasma omega-6 polyunsaturated fatty acid levels (gamma-linolenic acid), Plasma omega-6 polyunsaturated fatty acid levels (dihomo-gamma-linolenic acid), Plasma omega-6 polyunsaturated fatty acid levels (arachidonic	ENSG00000204950: Diastolic blood pressure, Systolic blood pressure, Systolic blood pressure (cigarette smoking interaction), Diastolic blood pressure (cigarette smoking interaction), Mean arterial pressure, Mean arterial pressure x alcohol consumption interaction (2df test), Diastolic blood pressure x alcohol consumption interaction (2df test), Hypertension, Medication use (agents acting on the renin-angiotensin system), Heel bone mineral density; ENSG0000011347: Intelligence (MTAG), Phosphatidylcholine levels, Cholesteryl ester levels, Educational attainment (years of education), Cognitive performance (MTAG), Cognitive performance; ENSG00000134780: Immune response to smallpox (secreted IL-2), Plasma omega-3 polyunsaturated fatty acid levels (docosapentaenoic acid), Plasma omega-6 polyunsaturated fatty acid levels (linoleic acid), 3-hydroxypropylmercapturic acid levels in smokers, Refractive error, Cerebrospinal fluid sTREM-2 levels, Spherical equivalent, C-reactive protein levels, Neuroticism, Positive affect, Depressive symptoms, Well-being spectrum (multivariate analysis), Life satisfaction, Depression, Lung function (FVC); ENSG00000124920: Plasma omega-3 polyunsaturated fatty acid levels (docosapentaenoic acid), Plasma omega-3 polyunsaturated fatty acid levels (alpha-linolenic acid), Plasma omega-3 polyunsaturated fatty acid level (eicosapentaenoic acid), Metabolite levels, Crohn's disease, Plasma omega-6 polyunsaturated fatty acid levels (gamma-linolenic acid), Plasma omega-6 polyunsaturated fatty acid levels (dihomo-gamma-linolenic acid), Plasma omega-6 polyunsaturated fatty acid levels (arachidonic acid), Hematocrit, Triglycerides, HDL cholesterol, Glycerophospholipid levels, Total cholesterol levels, Moyamoya disease, Resting heart rate, Red blood cell

						<p>acid), Delta-6 desaturase activity, Plasma omega-6 polyunsaturated fatty acid levels (linoleic acid), Eosinophil counts, C-reactive protein levels or LDL-cholesterol levels (pleiotropy), C-reactive protein levels or HDL-cholesterol levels (pleiotropy), C-reactive protein levels or triglyceride levels (pleiotropy), Metabolite levels (lipid measures), Gestational age at birth (child effect), Plateletcrit, Granulocyte percentage of myeloid white cells, Monocyte percentage of white cells, Glycerophospholipid levels, Sphingolipid levels, Vitiligo, Glycated hemoglobin levels, Total cholesterol levels, Heel bone mineral density, Non-albumin protein levels, Albumin-globulin ratio, Hemoglobin A1c levels, Alanine transaminase levels, Triglyceride levels, Breast milk fatty acid composition (maternal genotype effect), Breast milk fatty acid composition (infant genotype effect), Serum metabolite concentrations in chronic kidney disease, Serum metabolite ratios in chronic kidney disease, Low density lipoprotein cholesterol levels, Height, Apolipoprotein B levels, Pulse pressure, LDL cholesterol levels, LDL cholesterol levels x long total sleep time interaction (2df test), HDL cholesterol levels, Triglyceride levels x short total sleep time interaction (2df test), Serum</p>	<p>count, Serum total protein level, Serum metabolite concentrations in chronic kidney disease, Serum metabolite ratios in chronic kidney disease, Triglyceride levels x long total sleep time interaction (2df test), Serum metabolite levels (CMS), Heel bone mineral density, Hemoglobin levels, Spherical equivalent, Blood metabolite levels, Colorectal cancer, Serum omega-3 polyunsaturated fatty acid concentration in metabolic syndrome, Low density lipoprotein cholesterol levels, Serum metabolite levels, Chronic inflammatory diseases (ankylosing spondylitis, Crohn's disease, psoriasis, primary sclerosing cholangitis, ulcerative colitis) (pleiotropy), Hemoglobin concentration, Vaccenic acid (18:1n-7) levels, Gondoic acid (20:1n-9) levels, LDL cholesterol levels, Iron status biomarkers (total iron binding capacity), Trans fatty acid levels, High density lipoprotein cholesterol levels, Red blood cell fatty acid levels, Colorectal cancer or advanced adenoma, Phosphatidylcholine-ether levels, Phosphatidylethanolamine-ether levels, Asthma (adult onset), Asthma, Nasal polyps, Systolic blood pressure, Triglyceride levels, Glycemic traits (pleiotropy), Educational attainment (years of education), Respiratory diseases; ENSG00000134825: Crohn's disease, Metabolic syndrome, Palmitoleic acid (16:1n-7) levels, Stearic acid (18:0) levels, Oleic acid (18:1n-9) levels, Phospholipid levels (plasma), Plasma omega-3 polyunsaturated fatty acid levels (docosapentaenoic acid), Plasma omega-3 polyunsaturated fatty acid level (eicosapentaenoic acid), Plasma omega-3 polyunsaturated fatty acid levels (alpha-linolenic acid), Metabolite levels, Plasma omega-6 polyunsaturated fatty acid levels (gamma-linolenic acid), Plasma omega-6 polyunsaturated fatty acid levels (dihomo-gamma-linolenic acid), Plasma omega-6 polyunsaturated fatty acid levels (arachidonic acid), Hematocrit, Triglycerides, HDL cholesterol, Glycerophospholipid levels, Total cholesterol levels, Resting heart rate, Red blood cell count, Serum total protein level, Irritable mood, Serum metabolite concentrations in chronic kidney disease, Serum metabolite ratios in chronic kidney disease, LDL cholesterol levels x short total sleep time interaction (2df test), Triglyceride levels x long total sleep time interaction (2df test), Serum metabolite levels (CMS), Heel bone mineral density, Hemoglobin levels, Spherical equivalent,</p>
--	--	--	--	--	--	--	--

						<p>metabolite levels (CMS), HDL cholesterol levels x long total sleep time interaction (2df test), HDL cholesterol levels x short total sleep time interaction (2df test), HDL cholesterol levels x alcohol consumption (regular vs non-regular drinkers) interaction (2df), Triglyceride levels x alcohol consumption (regular vs non-regular drinkers) interaction (2df), Bipolar I disorder, Apolipoprotein A1 levels, Neutrophil count, Colorectal cancer, QRS duration, Crohn's disease, Age-related disease endophenotypes, Sum basophil neutrophil counts, Red cell distribution width, P wave duration, Blood metabolite levels, Blood metabolite ratios, Rheumatoid arthritis, QT interval, Iron status biomarkers (transferrin levels), Serum metabolite levels, Change in serum metabolite levels (CMS), Serum omega-3 polyunsaturated fatty acid concentration in metabolic syndrome, Serum omega-6 to omega-3 polyunsaturated fatty acid ratio in metabolic syndrome, Serum docosahexaenoic fatty acid concentration in metabolic syndrome, Delta-5 desaturase activity response to n3-polyunsaturated fat supplement, Change in serum metabolite levels, LDL cholesterol x physical activity interaction (2df test), High density lipoprotein</p>	<p>Blood metabolite levels, Colorectal cancer, Serum omega-3 polyunsaturated fatty acid concentration in metabolic syndrome, Serum metabolite levels, Blond vs. brown/black hair color, Low density lipoprotein cholesterol levels, Bipolar disorder or major depressive disorder, Chronic inflammatory diseases (ankylosing spondylitis, Crohn's disease, psoriasis, primary sclerosing cholangitis, ulcerative colitis) (pleiotropy), Hemoglobin concentration, Vaccenic acid (18:1n-7) levels, Gondoic acid (20:1n-9) levels, LDL cholesterol levels, Iron status biomarkers (total iron binding capacity), Trans fatty acid levels, High density lipoprotein cholesterol levels, Carboplatin disposition in epithelial ovarian cancer, Red blood cell fatty acid levels, Colorectal cancer or advanced adenoma, Phosphatidylcholine-ether levels, Phosphatidylcholine levels, Phosphatidylethanolamine-ether levels, Cholesteryl ester levels, Asthma (adult onset), Asthma, Nasal polyps, Triglyceride levels, Glycemic traits (pleiotropy), Educational attainment (years of education), Respiratory diseases, Hair color; ENSG00000168496: Plasma omega-3 polyunsaturated fatty acid levels (docosapentaenoic acid), Plasma omega-3 polyunsaturated fatty acid levels (alpha-linolenic acid), Plasma omega-3 polyunsaturated fatty acid level (eicosapentaenoic acid), Platelet count, Inflammatory bowel disease, Crohn's disease, Trans fatty acid levels, Metabolite levels, Red blood cell fatty acid levels, Colorectal cancer; ENSG00000134824: HDL cholesterol, Triglycerides, Fasting blood glucose, Homeostasis model assessment of beta-cell function, Cholesterol, total, LDL cholesterol, Heart rate, Metabolic syndrome, Resting heart rate, Lipid metabolism phenotypes, Hematology traits, Liver enzyme levels (alkaline phosphatase), Comprehensive strength and appendicular lean mass, Metabolite levels, Response to statin therapy, Metabolic traits, Plasma omega-3 polyunsaturated fatty acid levels (docosapentaenoic acid), Plasma omega-3 polyunsaturated fatty acid levels (alpha-linolenic acid), Plasma omega-3 polyunsaturated fatty acid level (eicosapentaenoic acid), Platelet count, Fasting blood glucose (BMI interaction), Inflammatory bowel disease, Plasma omega-6 polyunsaturated fatty acid levels (gamma-linolenic acid), Plasma omega-6 polyunsaturated fatty acid levels (dihomo-gamma-linolenic acid), Plasma omega-6 polyunsaturated fatty</p>
--	--	--	--	--	--	---	---

						<p>cholesterol levels, Asthma, Bipolar disorder, Male-pattern baldness, Nonatopic asthma, Osteoporosis-related phenotypes (MTAG), Granulocyte count, Age-related diseases, mortality and associated endophenotypes, Sum neutrophil eosinophil counts, Gondoic acid (20:1n-9) levels, Sum eosinophil basophil counts, Myeloid white cell count, White blood cell count, Mean platelet volume, Trans fatty acid levels, Red blood cell fatty acid levels, Plasma omega-6 polyunsaturated fatty acid levels (adrenic acid), Laryngeal squamous cell carcinoma, Lung cancer in ever smokers, Triacylglycerol 56:6 levels, Metabolite risk score for predicting weight gain, Fatty acid desaturase activity (serum), Fatty acid desaturase activity (adipose tissue), Phosphatidylcholine levels, Phosphatidylcholine-ether levels, Phosphatidylethanolamine levels, Lysophosphatidylethanolamine levels, Phosphatidylethanolamine-ether levels, Phosphatidylinositol levels, Cholesteryl ester levels, Sphingomyelin levels, Lysophosphatidylcholine levels, Triacylglyceride levels, HDL cholesterol levels in current drinkers, LDL cholesterol levels in current drinkers, Triglyceride levels in current drinkers, HDL cholesterol levels x alcohol consumption (drinkers vs non-</p>	<p>acid levels (arachidonic acid), Delta-6 desaturase activity, Plasma omega-6 polyunsaturated fatty acid levels (linoleic acid), Eosinophil counts, C-reactive protein levels or LDL-cholesterol levels (pleiotropy), C-reactive protein levels or HDL-cholesterol levels (pleiotropy), C-reactive protein levels or triglyceride levels (pleiotropy), Metabolite levels (lipid measures), Gestational age at birth (child effect), Plateletcrit, Granulocyte percentage of myeloid white cells, Monocyte percentage of white cells, Glycerophospholipid levels, Sphingolipid levels, Vitiligo, Glycated hemoglobin levels, Total cholesterol levels, Heel bone mineral density, Non-albumin protein levels, Albumin-globulin ratio, Hemoglobin A1c levels, Alanine transaminase levels, Triglyceride levels, Breast milk fatty acid composition (maternal genotype effect), Breast milk fatty acid composition (infant genotype effect), Serum metabolite concentrations in chronic kidney disease, Serum metabolite ratios in chronic kidney disease, Low density lipoprotein cholesterol levels, Height, Apolipoprotein B levels, Pulse pressure, LDL cholesterol levels, LDL cholesterol levels x long total sleep time interaction (2df test), HDL cholesterol levels, Triglyceride levels x short total sleep time interaction (2df test), Serum metabolite levels (CMS), HDL cholesterol levels x long total sleep time interaction (2df test), HDL cholesterol levels x short total sleep time interaction (2df test), HDL cholesterol levels x alcohol consumption (regular vs non-regular drinkers) interaction (2df), Triglyceride levels x alcohol consumption (regular vs non-regular drinkers) interaction (2df), Bipolar I disorder, Apolipoprotein A1 levels, Neutrophil count, Colorectal cancer, QRS duration, Crohn's disease, Age-related disease endophenotypes, Sum basophil neutrophil counts, Red cell distribution width, P wave duration, Blood metabolite levels, Blood metabolite ratios, Rheumatoid arthritis, QT interval, Iron status biomarkers (transferrin levels), Serum metabolite levels, Change in serum metabolite levels (CMS), Serum omega-3 polyunsaturated fatty acid concentration in metabolic syndrome, Serum omega-6 to omega-3 polyunsaturated fatty acid ratio in metabolic syndrome, Serum docosahexaenoic fatty acid concentration in metabolic syndrome, Delta-5 desaturase activity response to n3-polyunsaturated fat supplement, Change in serum metabolite levels, LDL cholesterol x physical activity</p>
--	--	--	--	--	--	---	---

						<p>drinkers) interaction (2df), Triglyceride levels x alcohol consumption (drinkers vs non-drinkers) interaction (2df), LDL cholesterol levels x alcohol consumption (drinkers vs non-drinkers) interaction (2df), LDL cholesterol levels x alcohol consumption (regular vs non-regular drinkers) interaction (2df), Red blood cell count, Aortic valve stenosis, Sleep duration, Hematocrit, Hemoglobin concentration, Total triglycerides levels, Lymphocyte counts, Mean corpuscular hemoglobin concentration, PR interval, Type 2 diabetes, Mean corpuscular volume, IgA levels, Adult onset asthma or type 2 diabetes, Nonatopic asthma or type 2 diabetes, Adult onset asthma or fasting glucose levels, Nonatopic asthma or fasting glucose levels, Asthma (adult onset), Medication use (thyroid preparations), Medication use (adrenergics, inhalants), Urate levels, Gallstone disease, anorexia nervosa, attention-deficit/hyperactivity disorder, autism spectrum disorder, bipolar disorder, major depression, obsessive-compulsive disorder, schizophrenia, or Tourette syndrome (pleiotropy), QT dynamics during exercise, Balding type 1, Hypothyroidism</p>	<p>interaction (2df test), High density lipoprotein cholesterol levels, Asthma, Bipolar disorder, Male-pattern baldness, Nonatopic asthma, Osteoporosis-related phenotypes (MTAG), Granulocyte count, Age-related diseases, mortality and associated endophenotypes, Sum neutrophil eosinophil counts, Gondoic acid (20:1n-9) levels, Sum eosinophil basophil counts, Myeloid white cell count, White blood cell count, Mean platelet volume, Trans fatty acid levels, Red blood cell fatty acid levels, Plasma omega-6 polyunsaturated fatty acid levels (adrenic acid), Laryngeal squamous cell carcinoma, Lung cancer in ever smokers, Triacylglycerol 56:6 levels, Metabolite risk score for predicting weight gain, Fatty acid desaturase activity (serum), Fatty acid desaturase activity (adipose tissue), Phosphatidylcholine levels, Phosphatidylcholine-ether levels, Phosphatidylethanolamine levels, Lysophosphatidylethanolamine levels, Phosphatidylethanolamine-ether levels, Phosphatidylinositol levels, Cholesteryl ester levels, Sphingomyelin levels, Lysophosphatidylcholine levels, Triacylglyceride levels, HDL cholesterol levels in current drinkers, LDL cholesterol levels in current drinkers, Triglyceride levels in current drinkers, HDL cholesterol levels x alcohol consumption (drinkers vs non-drinkers) interaction (2df), Triglyceride levels x alcohol consumption (drinkers vs non-drinkers) interaction (2df), LDL cholesterol levels x alcohol consumption (drinkers vs non-drinkers) interaction (2df), LDL cholesterol levels x alcohol consumption (regular vs non-regular drinkers) interaction (2df), Red blood cell count, Aortic valve stenosis, Sleep duration, Hematocrit, Hemoglobin concentration, Total triglycerides levels, Lymphocyte counts, Mean corpuscular hemoglobin concentration, PR interval, Type 2 diabetes, Mean corpuscular volume, IgA levels, Adult onset asthma or type 2 diabetes, Nonatopic asthma or type 2 diabetes, Adult onset asthma or fasting glucose levels, Nonatopic asthma or fasting glucose levels, Asthma (adult onset), Medication use (thyroid preparations), Medication use (adrenergics, inhalants), Urate levels, Gallstone disease, anorexia nervosa, attention-deficit/hyperactivity disorder, autism spectrum disorder, bipolar disorder, major depression, obsessive-compulsive disorder, schizophrenia, or Tourette syndrome (pleiotropy), QT dynamics during exercise, Balding type 1, Hypothyroidism;</p>
--	--	--	--	--	--	--	--

								<p>ENSG00000149485: HDL cholesterol, Triglycerides, Fasting blood glucose, Homeostasis model assessment of beta-cell function, LDL cholesterol, Heart rate, Metabolic syndrome, Resting heart rate, Lipid metabolism phenotypes, Hematology traits, Comprehensive strength and appendicular lean mass, Metabolite levels, Cholesterol, total, Metabolic traits, Plasma omega-3 polyunsaturated fatty acid levels (alpha-linolenic acid), Plasma omega-3 polyunsaturated fatty acid level (eicosapentaenoic acid), Plasma omega-3 polyunsaturated fatty acid levels (docosapentaenoic acid), Fasting blood glucose (BMI interaction), Plasma omega-6 polyunsaturated fatty acid levels (gamma-linolenic acid), Plasma omega-6 polyunsaturated fatty acid levels (dihomo-gamma-linolenic acid), Delta-6 desaturase activity, Plasma omega-6 polyunsaturated fatty acid levels (linoleic acid), Eosinophil counts, Platelet count, C-reactive protein levels or HDL-cholesterol levels (pleiotropy), C-reactive protein levels or triglyceride levels (pleiotropy), Metabolite levels (lipid measures), Granulocyte percentage of myeloid white cells, Monocyte percentage of white cells, Glycerophospholipid levels, Sphingolipid levels, Vitiligo, Total cholesterol levels, Alanine transaminase levels, Triglyceride levels, Breast milk fatty acid composition (maternal genotype effect), Breast milk fatty acid composition (infant genotype effect), Serum metabolite concentrations in chronic kidney disease, Serum metabolite ratios in chronic kidney disease, Low density lipoprotein cholesterol levels, Apolipoprotein B levels, Pulse pressure, LDL cholesterol levels, LDL cholesterol levels x long total sleep time interaction (2df test), HDL cholesterol levels, Triglyceride levels x short total sleep time interaction (2df test), Serum metabolite levels (CMS), HDL cholesterol levels x long total sleep time interaction (2df test), HDL cholesterol levels x short total sleep time interaction (2df test), HDL cholesterol levels x alcohol consumption (regular vs non-regular drinkers) interaction (2df), Triglyceride levels x alcohol consumption (regular vs non-regular drinkers) interaction (2df), Bipolar I disorder, Apolipoprotein A1 levels, Age-related disease endophenotypes, Red cell distribution width, Height, Blood metabolite levels, Blood metabolite ratios, Rheumatoid arthritis, Serum metabolite levels, Change in serum metabolite levels (CMS), Serum omega-3 polyunsaturated fatty acid</p>
--	--	--	--	--	--	--	--	--

								concentration in metabolic syndrome, Serum omega-6 to omega-3 polyunsaturated fatty acid ratio in metabolic syndrome, Serum docosahexaenoic fatty acid concentration in metabolic syndrome, Delta-5 desaturase activity response to n3-polyunsaturated fat supplement, Change in serum metabolite levels, LDL cholesterol x physical activity interaction (2df test), High density lipoprotein cholesterol levels, Asthma, Bipolar disorder, Osteoporosis-related phenotypes (MTAG), Granulocyte count, Age-related diseases, mortality and associated endophenotypes, Sum neutrophil eosinophil counts, Sum eosinophil basophil counts, Myeloid white cell count, White blood cell count, Mean platelet volume, Trans fatty acid levels, Red blood cell fatty acid levels, Plasma omega-6 polyunsaturated fatty acid levels (arachidonic acid), Plasma omega-6 polyunsaturated fatty acid levels (adrenic acid), Laryngeal squamous cell carcinoma, Triacylglycerol 56:6 levels, Metabolite risk score for predicting weight gain, Fatty acid desaturase activity (serum), Fatty acid desaturase activity (adipose tissue), Phosphatidylcholine-ether levels, Lysophosphatidylethanolamine levels, Phosphatidylcholine levels, Phosphatidylethanolamine-ether levels, Phosphatidylinositol levels, Cholesteryl ester levels, Sphingomyelin levels, Lysophosphatidylcholine levels, Triacylglyceride levels, HDL cholesterol levels in current drinkers, LDL cholesterol levels in current drinkers, Triglyceride levels in current drinkers, HDL cholesterol levels x alcohol consumption (drinkers vs non-drinkers) interaction (2df), Triglyceride levels x alcohol consumption (drinkers vs non-drinkers) interaction (2df), LDL cholesterol levels x alcohol consumption (drinkers vs non-drinkers) interaction (2df), LDL cholesterol levels x alcohol consumption (regular vs non-regular drinkers) interaction (2df), Red blood cell count, QT interval, Aortic valve stenosis, Sleep duration, Neutrophil count, Total triglycerides levels, Mean corpuscular hemoglobin concentration, Mean corpuscular volume, IgA levels, Gallstone disease; ENSG00000221968: Sphingolipid levels, Metabolite levels, Phosphatidylcholine levels, Cholesteryl ester levels, Sphingolipid d18:1/d18:2 ratio, Urate levels; ENSG00000167994: Sphingolipid levels, Metabolite levels, Schizophrenia, Mean corpuscular hemoglobin; ENSG00000167995: Plasma omega-3 polyunsaturated fatty
--	--	--	--	--	--	--	--	--

							acid level (eicosapentaenoic acid), Metabolite levels; ENSG00000167996; ; ENSG00000149503: Prostate cancer, Breast cancer; ENSG00000168515; ; ENSG00000124939; ; ENSG00000124935; ; ENSG00000110484; ; ENSG00000197745:	
16	66852047	67871804	67473389	5.25E-08	ENSG00000159593, ENSG00000168748, ENSG00000172840, ENSG00000166589, ENSG00000166592, ENSG00000166595, ENSG00000172831, ENSG00000172828, ENSG00000172824, ENSG00000067955, ENSG00000125149, ENSG00000237172, ENSG00000102871, ENSG00000135722, ENSG00000102878, ENSG00000140939, ENSG00000196123, ENSG00000179044, ENSG00000205250, ENSG00000102890, ENSG00000125122, ENSG00000168701, ENSG00000135723, ENSG00000135740, ENSG00000196155, ENSG00000168676, ENSG00000159708, ENSG00000159713, ENSG00000159714, ENSG00000176387, ENSG00000159720 , ENSG00000159723, ENSG00000039523, ENSG00000102974, ENSG00000159753,	NAE1, CA7, PDP2, CDH16, RRAD, FAM96B, CES2, CES3, CES4A, CBF, C16orf70, B3GNT9, TRADD, FBXL8, HSF4, NOL3, KIAA0895L, EXOC3L1, E2F4, ELMO3, LRRC29, TMEM208, FHOD1, SLC9A5, PLEKHG4, KCTD19, LRRC36, TPPP3, ZDHHC1, HSD11B2, ATP6VOD1 , AGRP, FAM65A, CTCF, RLTPR, ACD, PARD6A, ENKD1, C16orf86,	ENSG00000159720: Disease progression in age-related macular degeneration, Medication use (thyroid preparations)	ENSG00000159593: HDL cholesterol levels, Apolipoprotein A1 levels; ENSG00000168748; ; ENSG00000172840: Diastolic blood pressure; ENSG00000166589: Blood protein levels in cardiovascular risk, HDL cholesterol levels x alcohol consumption (regular vs non-regular drinkers) interaction (2df), HDL cholesterol levels in current drinkers, HDL cholesterol levels x alcohol consumption (drinkers vs non-drinkers) interaction (2df); ENSG00000166592; ; ENSG00000166595; ; ENSG00000172831; ; ENSG00000172828: Heel bone mineral density, Blood protein levels; ENSG00000172824: Hemoglobin levels; ENSG00000067955: Mean corpuscular hemoglobin, Breast cancer specific mortality in estrogen receptor positive breast cancer, Blood protein levels, Red cell distribution width; ENSG00000125149: Adolescent idiopathic scoliosis, Height; ENSG00000237172; ; ENSG00000102871; ; ENSG00000135722; ; ENSG00000102878; ; ENSG00000140939; ; ENSG00000196123: Mean corpuscular volume, Heel bone mineral density; ENSG00000179044: Hemoglobin levels, Monocyte count; ENSG00000205250: Mean corpuscular volume, Mean corpuscular hemoglobin; ENSG00000102890: Intraocular pressure; ENSG00000125122: HDL cholesterol, Male-pattern baldness; ENSG00000168701; ; ENSG00000135723: Waist circumference adjusted for body mass index, Waist-to-hip ratio adjusted for BMI; ENSG00000135740; ; ENSG00000196155: Heel bone mineral density, Blood protein levels; ENSG00000168676: Waist circumference adjusted for body mass index, Waist circumference adjusted for BMI (adjusted for smoking behaviour), Waist circumference adjusted for BMI (joint analysis main effects and smoking interaction), Waist circumference adjusted for BMI in non-smokers, Body mass index, Height, Hypothyroidism; ENSG00000159708: Waist circumference adjusted for body mass index, Adolescent idiopathic scoliosis, Blood protein levels, Waist circumference

					ENSG00000102977, ENSG00000102981, ENSG00000124074, ENSG00000159761, ENSG00000141098, ENSG00000141084, ENSG00000102904, ENSG00000102901	GFOD2, RANBP10, TSNAXIP1, CENPT		adjusted for BMI (joint analysis main effects and physical activity interaction), Waist circumference adjusted for BMI in active individuals, Waist-to-hip ratio adjusted for BMI, Waist-hip ratio, Waist-to-hip ratio adjusted for BMI (additive genetic model), Mean corpuscular hemoglobin, Eosinophil counts; ENSG00000159713; ; ENSG00000159714: Bone mineral density; ENSG00000176387; ; ENSG00000159720: Disease progression in age-related macular degeneration, Medication use (thyroid preparations); ENSG00000159723; ; ENSG00000039523: Heel bone mineral density, HDL cholesterol, Hemoglobin levels, Mean platelet volume, Lung function (FVC), Lung function (FEV1/FVC); ENSG00000102974: Emotional recognition, Mean corpuscular hemoglobin concentration; ENSG00000159753: HDL cholesterol levels x long total sleep time interaction (2df test), HDL cholesterol, Myopia (age of diagnosis), Lymphocyte counts, Eosinophil counts; ENSG00000102977: Obesity-related traits, Hemoglobin levels; ENSG00000102981; ; ENSG00000124074: Blood protein levels; ENSG00000159761; ; ENSG00000141098: HDL cholesterol, Male-pattern baldness, High density lipoprotein cholesterol levels; ENSG00000141084: Hematocrit, HDL cholesterol, Empathy quotient, Hemoglobin concentration, High density lipoprotein cholesterol levels; ENSG00000102904: Balding type 1; ENSG00000102901: Red blood cell count, Mean corpuscular hemoglobin
2	154307589	154707589	154507589	6.80E-08	ENSG00000177519	RPRM	ENSG00000177519:	ENSG00000177519:
5	94974878	95385933	95174878	8.07E-08	ENSG00000175449, ENSG00000145757, ENSG00000164292, ENSG00000173221, ENSG00000236882, ENSG00000118985	RFESD, SPATA9, RHOBTB3, GLRX, C5orf27, ELL2	ENSG00000236882:	ENSG00000175449: Blood protein levels, Pharmacokinetics of antiepileptic drugs in severe mental disorder (concentration drug ratio); ENSG00000145757: Blood protein levels, Lung function (FEV1/FVC); ENSG00000164292: Metabolite levels, Adolescent idiopathic scoliosis; ENSG00000173221: Metabolite levels, Adolescent idiopathic scoliosis; ENSG00000236882; ; ENSG00000118985: IgG glycosylation, Multiple myeloma, Multiple myeloma and monoclonal gammopathy, Coronary artery calcified atherosclerotic plaque score in type 2 diabetes, Non-albumin protein levels, Albumin-globulin ratio, B-cell malignancies (chronic lymphocytic leukemia, Hodgkin lymphoma or multiple myeloma) (pleiotropy), Serum total protein level, IgG digalactosylation phenotypes (multivariate analysis), IgG sialylation phenotypes

								(multivariate analysis), Mean corpuscular hemoglobin, IgA levels, Mean corpuscular volume
11	60728079	61159993	60928079	9.36E-08	ENSG00000013725, ENSG00000110448 , ENSG00000167987, ENSG00000229859, ENSG00000229183, ENSG00000256713, ENSG00000167992, ENSG00000167986, ENSG00000149476, ENSG00000162144, ENSG00000149483, ENSG00000187049	CD6, CD5, VPS37C , PGA3, PGA4, PGA5, VWCE, DDB1, DAK, CYB561A3, TMEM138, TMEM216	ENSG00000167987: Rheumatoid arthritis (ACPA-positive), Rheumatoid arthritis, Adolescent idiopathic scoliosis, Periventricular white matter hyperintensities	ENSG00000013725: Multiple sclerosis, Inflammatory bowel disease, Crohn's disease, Ulcerative colitis, Chronic inflammatory diseases (ankylosing spondylitis, Crohn's disease, psoriasis, primary sclerosing cholangitis, ulcerative colitis) (pleiotropy), CD6 levels; ENSG00000110448: ; ENSG00000167987: Rheumatoid arthritis (ACPA-positive), Rheumatoid arthritis, Adolescent idiopathic scoliosis, Periventricular white matter hyperintensities; ENSG00000229859: ; ENSG00000229183: ; ENSG00000256713: Height; ENSG00000167992: ; ENSG00000167986: ; ENSG00000149476: ; ENSG00000162144: ; ENSG00000149483: ; ENSG00000187049:
9	98012608	98509513	98296403	1.14E-07	ENSG00000158169, ENSG00000185920	FANCC, PTCH1	ENSG00000185920: Pulmonary function, Height, Pulmonary function (smoking interaction), Bone mineral density (spine), Bone mineral density (hip), Waist circumference adjusted for body mass index, Hip circumference adjusted for BMI, Nonsyndromic cleft lip with cleft palate, Birth weight, Heel bone mineral density, Neuroticism, Neuroticism, Feeling hurt, Feeling worry, Alanine aminotransferase levels in low alcohol consumption, Appendicular lean mass, Cortical surface area (visual PC2), Birth weight (MTAG), Birth length (MTAG), Monobrow, Intelligence, Worry, Lung function (FVC), Waist circumference adjusted for BMI (adjusted for smoking behaviour), Waist circumference adjusted for BMI (joint analysis main effects and smoking interaction), Waist circumference adjusted for BMI in non-smokers, 3-month functional outcome in ischaemic stroke (modified Rankin score), Cortical surface area, Depressive symptoms, Positive affect, Offspring birth weight, Lung function (FEV1/FVC), Well-being spectrum (multivariate analysis), White matter microstructure (fractional anisotropy), Type 2 diabetes, Life satisfaction, Macular thickness, Sensitivity to environmental stress and adversity, Cognitive performance, Cognitive performance (MTAG), Highest math class taken (MTAG), White blood cell count, Sunburns	ENSG00000158169: Cortical surface area (visual PC2), Hematocrit, Waist circumference adjusted for body mass index, Heel bone mineral density, Height; ENSG00000185920: Pulmonary function, Height, Pulmonary function (smoking interaction), Bone mineral density (spine), Bone mineral density (hip), Waist circumference adjusted for body mass index, Hip circumference adjusted for BMI, Nonsyndromic cleft lip with cleft palate, Birth weight, Heel bone mineral density, Neuroticism, Neuroticism, Feeling hurt, Feeling worry, Alanine aminotransferase levels in low alcohol consumption, Appendicular lean mass, Cortical surface area (visual PC2), Birth weight (MTAG), Birth length (MTAG), Monobrow, Intelligence, Worry, Lung function (FVC), Waist circumference adjusted for BMI (adjusted for smoking behaviour), Waist circumference adjusted for BMI (joint analysis main effects and smoking interaction), Waist circumference adjusted for BMI in non-smokers, 3-month functional outcome in ischaemic stroke (modified Rankin score), Cortical surface area, Depressive symptoms, Positive affect, Offspring birth weight, Lung function (FEV1/FVC), Well-being spectrum (multivariate analysis), White matter microstructure (fractional anisotropy), Type 2 diabetes, Life satisfaction, Macular thickness, Sensitivity to environmental stress and adversity, Cognitive performance, Cognitive performance (MTAG), Highest math class taken (MTAG), White blood cell count, Sunburns

							circumference adjusted for BMI in non-smokers, 3-month functional outcome in ischaemic stroke (modified Rankin score), Cortical surface area, Depressive symptoms, Positive affect, Offspring birth weight, Lung function (FEV1/FVC), Well-being spectrum (multivariate analysis), White matter microstructure (fractional anisotropy), Type 2 diabetes, Life satisfaction, Macular thickness, Sensitivity to environmental stress and adversity, Cognitive performance, Cognitive performance (MTAG), Highest math class taken (MTAG), White blood cell count, Sunburns	
2	101127642	101528728	101328728	1.32E-07	ENSG00000115539, ENSG00000170485	PDCL3, NPAS2	ENSG00000115539:	ENSG00000115539: ; ENSG00000170485: 3-hydroxy-1-methylpropylmercapturic acid levels in smokers, Facial morphology (factor 15, philtrum width), Intraocular pressure, Serum 25-Hydroxyvitamin D levels, Intraocular pressure and central corneal thickness (multi-trait analysis), Pulse pressure, Chronotype, Educational attainment (years of education), Educational attainment (MTAG), Waist-hip ratio, White blood cell count
13	97953799	98367921	98153799	1.33E-07	ENSG00000139793, ENSG00000125249	MBNL2, RAP2A	ENSG00000125249:	ENSG00000139793: Alcoholism (alcohol use disorder factor score), Alcoholism (alcohol dependence factor score), Energy expenditure (24h), Platelet reactivity in response to clopidogrel treatment, Diastolic blood pressure, Pre-treatment viral load in HIV-1 infection, Self-reported math ability, Self-reported math ability (MTAG), Highest math class taken (MTAG), Menarche (age at onset); ENSG00000125249:
19	54178627	54585820	54385437	1.43E-07	ENSG00000142405, ENSG00000179820, ENSG00000126583, ENSG00000105605, ENSG00000142408, ENSG00000130433,	NLRP12, MYADM, PRKCG, CACNG7, CACNG8, CACNG6,	ENSG00000179820:	ENSG00000142405: Granulocyte percentage of myeloid white cells, Blood protein levels, Macrophage Migration Inhibitory Factor levels, Monocyte percentage of white cells, Monocyte count; ENSG00000179820: ; ENSG00000126583: ; ENSG00000105605: General cognitive ability, Educational attainment (MTAG), Cognitive performance (MTAG), Self-reported math ability, Self-reported math ability (MTAG),

					ENSG00000189068, ENSG00000248385	VSTM1, TARM1		Highest math class taken (MTAG); ENSG00000142408; ; ENSG00000130433; ; ENSG00000189068: Blood protein levels; ENSG00000248385:
4	7809737	8218920	8016164	2.53E-07	ENSG00000196526, ENSG00000163995, ENSG00000125089	AFAP1, ABLIM2, SH3TC1	ENSG00000163995: Post bronchodilator FEV1 in COPD, Cognitive function, Triacylglyceride levels, Early onset periodontitis x smoking status interaction, Irritable bowel syndrome, Diffusing capacity of carbon monoxide, Alcohol consumption (drinks per week) (MTAG), Response to cognitive- behavioural therapy in major depressive disorder	ENSG00000196526: Post bronchodilator FEV1/FVC ratio, Intraocular pressure, Pulse pressure, Glaucoma (primary open- angle), Glaucoma, Glaucoma (multi-trait analysis), Alanine aminotransferase levels in excessive alcohol consumption, HDL cholesterol levels x short total sleep time interaction (2df test), Granulocyte-colony stimulating factor levels, Lung function (FEV1/FVC), Diverticular disease, Mean platelet volume, Platelet count, Medication use (antiglaucoma preparations and miotics), Height, Systolic blood pressure; ENSG00000163995: Post bronchodilator FEV1 in COPD, Cognitive function, Triacylglyceride levels, Early onset periodontitis x smoking status interaction, Irritable bowel syndrome, Diffusing capacity of carbon monoxide, Alcohol consumption (drinks per week) (MTAG), Response to cognitive-behavioural therapy in major depressive disorder; ENSG00000125089: Response to metformin (IC50), Idiopathic downbeat nystagmus, Low density lipoprotein cholesterol levels
4	76831352	77231352	77031352	3.04E-07	ENSG00000138744, ENSG00000198301, ENSG00000138755, ENSG00000156219, ENSG00000169245, ENSG00000169248, ENSG00000138750, ENSG00000138760, ENSG00000189157, ENSG00000118804	NAAA, SDAD1, CXCL9, ART3, CXCL10, CXCL11, NUP54, SCARB2 , FAM47E, FAM47E- STBD1	ENSG00000138750: Deliberate self-harm	ENSG00000138744: Coronary artery calcified atherosclerotic plaque (130 HU threshold) in type 2 diabetes, Blood protein levels, Neurological blood protein biomarker levels; ENSG00000198301: Longevity, Monokine induced by gamma interferon levels, Hippocampal volume, Interferon gamma- induced protein 10 levels; ENSG00000138755: Blood protein levels; ENSG00000156219: Blood protein levels, C-X-C motif chemokine 10 levels, Monokine induced by gamma interferon levels, Hippocampal volume, Neonatal cytokine/chemokine levels (fetal genetic effect), Type 2 diabetes, Body mass index, Mean corpuscular hemoglobin; ENSG00000169245: Blood protein levels, C-X-C motif chemokine 10 levels, Neonatal cytokine/chemokine levels (fetal genetic effect); ENSG00000169248: Blood protein levels; ENSG00000138750: Deliberate self-harm; ENSG00000138760: Body mass index, Parkinson's disease or first degree relation to individual with Parkinson's disease, Body mass index (joint analysis main effects and smoking interaction), BMI (adjusted for smoking behaviour); ENSG00000189157: Parkinson's disease,

								Parkinson's disease or first degree relation to individual with Parkinson's disease, HDL cholesterol levels x long total sleep time interaction (2df test), Platelet count, Mean corpuscular hemoglobin, Mean corpuscular volume, Blood protein levels; ENSG00000118804:
3	167638903	168093774	167852482	4.62E-07	ENSG00000173905	GOLIM4	ENSG00000173905: Response to metformin (IC50), Metabolite levels, Cerebrospinal fluid t-tau levels in Alzheimer's disease dementia, Self-reported math ability, Height, Self-reported math ability (MTAG)	ENSG00000173905: Response to metformin (IC50), Metabolite levels, Cerebrospinal fluid t-tau levels in Alzheimer's disease dementia, Self-reported math ability, Height, Self-reported math ability (MTAG)
5	85631480	86067680	85867680	4.80E-07	ENSG00000127184	COX7C	ENSG00000127184:	ENSG00000127184:
15	28144238	28735266	28344238	5.46E-07	ENSG00000104044, ENSG00000128731, ENSG00000153684	OCA2, HERC2, GOLGA8F	ENSG00000128731: Black vs. blond hair color, Black vs. red hair color, Eye color, Vitiligo, Hair color, Tanning, Eye color traits, Blond vs. brown hair color, Blue vs. green eyes, Blue vs. brown eyes, Iris color, Multiple myeloma (IgH translocation), Alzheimer disease and age of onset, Squamous cell carcinoma, Skin colour saturation, Perceived skin darkness, Skin sensitivity to sun, Blond vs non-blond hair color, Brown vs. non-brown hair color, Light vs. dark hair color, Basal cell carcinoma, Osteoarthritis of the hip (with total joint replacement), Intraocular pressure, Skin pigmentation (conditioned on rs1426654 and rs35397), Skin pigmentation, Cutaneous melanoma or hair colour, Cutaneous malignant melanoma, Nevus count or cutaneous melanoma, Diisocyanate-induced asthma, Refractive astigmatism, Corneal	ENSG00000104044: Black vs. blond hair color, Black vs. red hair color, Eye color, Lung function (forced expiratory flow during mid-portion (25% and 75%) of forced vital capacity), Squamous cell carcinoma, Post bronchodilator FEV1/FVC ratio, Blond vs non-blond hair color, Brown vs. non-brown hair color, Light vs. dark hair color, Red vs non-red hair color, Facial morphology (factor 19), Uveal melanoma, Iris color (b* coordinate), Iris color (L* coordinate), Iris color (a* coordinate), Skin pigmentation (conditioned on rs1426654 and rs35397), Cataracts (operation), Cutaneous malignant melanoma, Central retinal vein equivalent, Central retinal arteriolar equivalent, Low tan response, Eye color (brightness), Eye color (hue), Skin, hair and eye pigmentation (multivariate analysis), Blond vs. brown/black hair color, Brown vs. black hair color, Cutaneous squamous cell carcinoma, Skin pigmentation, Eye color (saturation), Melanoma, Shingles, Eye color traits, Skin pigmentation traits, Macular thickness, Hair color, Sunburns; ENSG00000128731: Black vs. blond hair color, Black vs. red hair color, Eye color, Vitiligo, Hair color, Tanning, Eye color traits, Blond vs. brown hair color, Blue vs. green eyes, Blue vs. brown eyes, Iris color, Multiple myeloma (IgH translocation), Alzheimer disease and age of onset, Squamous cell carcinoma, Skin colour saturation, Perceived skin darkness, Skin sensitivity to sun, Blond vs non-blond hair color, Brown vs. non-brown hair color, Light vs. dark hair color, Basal cell carcinoma, Osteoarthritis of the hip (with total joint replacement), Intraocular pressure, Skin pigmentation

							astigmatism, Bone mineral content, Low tan response, Rosacea symptom severity, Eye color (hue), Skin, hair and eye pigmentation (multivariate analysis), Brown vs. black hair color, Red vs. brown/black hair color, Blond vs. brown/black hair color, Eye color (saturation), Eye color (brightness), Monobrow, Colorectal cancer, Colonoscopy-negative controls vs population controls, Keratinocyte cancer (MTAG), Glaucoma, Iris color (b* coordinate), Iris heterochromicity, Iris color (L* coordinate), Type 2 diabetes, Skin pigmentation traits, Hair morphology traits, Sunburns	(conditioned on rs1426654 and rs35397), Skin pigmentation, Cutaneous melanoma or hair colour, Cutaneous malignant melanoma, Nevus count or cutaneous melanoma, Diisocyanate-induced asthma, Refractive astigmatism, Corneal astigmatism, Bone mineral content, Low tan response, Rosacea symptom severity, Eye color (hue), Skin, hair and eye pigmentation (multivariate analysis), Brown vs. black hair color, Red vs. brown/black hair color, Blond vs. brown/black hair color, Eye color (saturation), Eye color (brightness), Monobrow, Colorectal cancer, Colonoscopy-negative controls vs population controls, Keratinocyte cancer (MTAG), Glaucoma, Iris color (b* coordinate), Iris heterochromicity, Iris color (L* coordinate), Type 2 diabetes, Skin pigmentation traits, Hair morphology traits, Sunburns; ENSG00000153684:
--	--	--	--	--	--	--	--	--

4604

4605 **Table S4d.2.** Top 25 regions from west-eurasia scan.

4606

4607

chr	start	end	bestpos	minpvalue	ensembl	hgnc	disease_trait_best	disease_trait
15	48233494	48633494	48433494	2.02E-31	ENSG00000188467, ENSG00000104177, ENSG00000233932, ENSG00000074803, ENSG00000128951	SLC24A5, MYEF2, CTXN2, SLC12A1, DUT	ENSG00000188467: Hair color, Skin pigmentation, Body mass index, Skin, hair and eye pigmentation (multivariate analysis), Eye color, Eye color (brightness), Eye color (saturation), Iris color (a* coordinate), Skin reflectance (Melanin index), Iris color (b* coordinate); ENSG00000104177: Skin pigmentation, Hair color; ENSG00000233932: ; ENSG00000074803: Systemic lupus erythematosus, Longevity, Skin reflectance (Melanin index), Lymphocyte counts; ENSG00000128951: Protein quantitative trait loci	ENSG00000188467: Hair color, Skin pigmentation, Body mass index, Skin, hair and eye pigmentation (multivariate analysis), Eye color, Eye color (brightness), Eye color (saturation), Iris color (a* coordinate), Skin reflectance (Melanin index), Iris color (b* coordinate); ENSG00000104177: Skin pigmentation, Hair color; ENSG00000233932: ; ENSG00000074803: Systemic lupus erythematosus, Longevity, Skin reflectance (Melanin index), Lymphocyte counts; ENSG00000128951: Protein quantitative trait loci
6	134192815	134829070	134392815	4.09E-26	ENSG00000118526, ENSG00000028839, ENSG00000146411, ENSG00000118515	TCF21, TBPL1, SLC2A12, SGK1	ENSG00000146411: Coronary artery disease, High chromosomal aberration frequency (chromosome type), FEV1, Lung function (FVC)	ENSG00000118526: Coronary heart disease, Coronary artery disease or ischemic stroke, Coronary artery disease, Coronary artery disease or large artery stroke, PR interval, Medication use (diuretics), Lung function (FEV1/FVC); ENSG00000028839: Mean corpuscular hemoglobin; ENSG00000146411: Coronary artery disease, High chromosomal aberration frequency (chromosome type), FEV1, Lung function (FVC);

								ENSG00000118515: Immune reponse to smallpox (secreted IFN-alpha), Alzheimer disease and age of onset, Pelvic organ prolapse, Pelvic organ prolapse (moderate/severe), Blond vs. brown/black hair color, Schizophrenia (inflammation and infection response interaction), Body mass index (smoking years interaction), Metabolite levels, Adolescent idiopathic scoliosis, Hair color
11	61264124	62263869	61770303	8.70E-22	ENSG00000204950, ENSG00000011347, ENSG00000134780, ENSG00000124920, ENSG00000134825, ENSG00000168496, ENSG00000134824, ENSG00000149485, ENSG00000221968, ENSG00000167994, ENSG00000167995, ENSG00000167996, ENSG00000149503, ENSG00000168515, ENSG00000124939, ENSG00000124935, ENSG00000110484, ENSG00000197745, ENSG00000162174, ENSG00000149021, ENSG00000124942	LRR10B, SYT7, DAGLA, MYRF, TMEM258, FEN1, FADS2, FADS1, FADS3, RAB31L1, BEST1, FTH1, INCENP, SCGB1D1, SCGB2A1, SCGB1D2, SCGB2A2, SCGB1D4, ASRGL1, SCGB1A1, AHNAK	ENSG00000167996:	ENSG00000204950: Diastolic blood pressure, Systolic blood pressure, Systolic blood pressure (cigarette smoking interaction), Diastolic blood pressure (cigarette smoking interaction), Mean arterial pressure, Mean arterial pressure x alcohol consumption interaction (2df test), Diastolic blood pressure x alcohol consumption interaction (2df test), Hypertension, Medication use (agents acting on the renin-angiotensin system), Heel bone mineral density; ENSG0000011347: Intelligence (MTAG), Phosphatidylcholine levels, Cholesteryl ester levels, Educational attainment (years of education), Cognitive performance (MTAG), Cognitive performance; ENSG00000134780: Immune reponse to smallpox (secreted IL-2), Plasma omega-3 polyunsaturated fatty acid levels (docosapentaenoic acid), Plasma omega-6 polyunsaturated fatty acid levels (linoleic acid), 3-hydroxypropylmercapturic acid levels in smokers, Refractive error, Cerebrospinal fluid sTREM-2 levels, Spherical equivalent, C-reactive protein levels, Neuroticism, Positive affect, Depressive symptoms, Well-being spectrum (multivariate analysis), Life satisfaction, Depression, Lung function (FVC); ENSG00000124920: Plasma omega-3 polyunsaturated fatty acid levels (docosapentaenoic acid), Plasma omega-3 polyunsaturated fatty acid levels (alpha-linolenic acid), Plasma omega-3 polyunsaturated fatty acid level (eicosapentaenoic acid), Metabolite levels, Crohn's disease, Plasma omega-6 polyunsaturated fatty acid levels (gamma-linolenic acid), Plasma omega-6 polyunsaturated fatty acid levels (dihomo-gamma-linolenic acid), Plasma omega-6 polyunsaturated fatty acid levels (arachidonic acid), Hematocrit, Triglycerides, HDL cholesterol, Glycerophospholipid levels, Total cholesterol levels, Moyamoya disease, Resting heart rate, Red blood cell count, Serum total protein level, Serum metabolite concentrations in chronic kidney disease, Serum metabolite ratios in chronic kidney disease, Triglyceride levels x long total sleep time interaction (2df test), Serum metabolite levels (CMS), Heel bone mineral density, Hemoglobin levels, Spherical equivalent, Blood metabolite levels, Colorectal cancer, Serum omega-3 polyunsaturated fatty acid concentration in metabolic syndrome, Low density lipoprotein cholesterol levels, Serum metabolite levels, Chronic inflammatory diseases (ankylosing spondylitis, Crohn's disease, psoriasis, primary sclerosing

								<p> cholangitis, ulcerative colitis) (pleiotropy), Hemoglobin concentration, Vaccenic acid (18:1n-7) levels, Gondoic acid (20:1n-9) levels, LDL cholesterol levels, Iron status biomarkers (total iron binding capacity), Trans fatty acid levels, High density lipoprotein cholesterol levels, Red blood cell fatty acid levels, Colorectal cancer or advanced adenoma, Phosphatidylcholine-ether levels, Phosphatidylethanolamine-ether levels, Asthma (adult onset), Asthma, Nasal polyps, Systolic blood pressure, Triglyceride levels, Glycemic traits (pleiotropy), Educational attainment (years of education), Respiratory diseases; ENSG00000134825: Crohn's disease, Metabolic syndrome, Palmitoleic acid (16:1n-7) levels, Stearic acid (18:0) levels, Oleic acid (18:1n-9) levels, Phospholipid levels (plasma), Plasma omega-3 polyunsaturated fatty acid levels (docosapentaenoic acid), Plasma omega-3 polyunsaturated fatty acid level (eicosapentaenoic acid), Plasma omega-3 polyunsaturated fatty acid levels (alpha-linolenic acid), Metabolite levels, Plasma omega-6 polyunsaturated fatty acid levels (gamma-linolenic acid), Plasma omega-6 polyunsaturated fatty acid levels (dihomo-gamma-linolenic acid), Plasma omega-6 polyunsaturated fatty acid levels (arachidonic acid), Hematocrit, Triglycerides, HDL cholesterol, Glycerophospholipid levels, Total cholesterol levels, Resting heart rate, Red blood cell count, Serum total protein level, Irritable mood, Serum metabolite concentrations in chronic kidney disease, Serum metabolite ratios in chronic kidney disease, LDL cholesterol levels x short total sleep time interaction (2df test), Triglyceride levels x long total sleep time interaction (2df test), Serum metabolite levels (CMS), Heel bone mineral density, Hemoglobin levels, Spherical equivalent, Blood metabolite levels, Colorectal cancer, Serum omega-3 polyunsaturated fatty acid concentration in metabolic syndrome, Serum metabolite levels, Blond vs. brown/black hair color, Low density lipoprotein cholesterol levels, Bipolar disorder or major depressive disorder, Chronic inflammatory diseases (ankylosing spondylitis, Crohn's disease, psoriasis, primary sclerosing cholangitis, ulcerative colitis) (pleiotropy), Hemoglobin concentration, Vaccenic acid (18:1n-7) levels, Gondoic acid (20:1n-9) levels, LDL cholesterol levels, Iron status biomarkers (total iron binding capacity), Trans fatty acid levels, High density lipoprotein cholesterol levels, Carboplatin disposition in epithelial ovarian cancer, Red blood cell fatty acid levels, Colorectal cancer or advanced adenoma, Phosphatidylcholine-ether levels, Phosphatidylcholine levels, Phosphatidylethanolamine-ether levels, Cholesteryl ester levels, Asthma (adult onset), Asthma, Nasal polyps, Triglyceride levels, Glycemic traits (pleiotropy), Educational attainment (years of education), Respiratory diseases, Hair color; ENSG00000168496: Plasma omega-3 polyunsaturated fatty acid levels </p>
--	--	--	--	--	--	--	--	--

								<p>Gallstone disease, anorexia nervosa, attention-deficit/hyperactivity disorder, autism spectrum disorder, bipolar disorder, major depression, obsessive-compulsive disorder, schizophrenia, or Tourette syndrome (pleiotropy), QT dynamics during exercise, Balding type 1, Hypothyroidism; ENSG00000149485: HDL cholesterol, Triglycerides, Fasting blood glucose, Homeostasis model assessment of beta-cell function, LDL cholesterol, Heart rate, Metabolic syndrome, Resting heart rate, Lipid metabolism phenotypes, Hematology traits, Comprehensive strength and appendicular lean mass, Metabolite levels, Cholesterol, total, Metabolic traits, Plasma omega-3 polyunsaturated fatty acid levels (alpha-linolenic acid), Plasma omega-3 polyunsaturated fatty acid level (eicosapentaenoic acid), Plasma omega-3 polyunsaturated fatty acid levels (docosapentaenoic acid), Fasting blood glucose (BMI interaction), Plasma omega-6 polyunsaturated fatty acid levels (gamma-linolenic acid), Plasma omega-6 polyunsaturated fatty acid levels (dihomo-gamma-linolenic acid), Delta-6 desaturase activity, Plasma omega-6 polyunsaturated fatty acid levels (linoleic acid), Eosinophil counts, Platelet count, C-reactive protein levels or HDL-cholesterol levels (pleiotropy), C-reactive protein levels or triglyceride levels (pleiotropy), Metabolite levels (lipid measures), Granulocyte percentage of myeloid white cells, Monocyte percentage of white cells, Glycerophospholipid levels, Sphingolipid levels, Vitiligo, Total cholesterol levels, Alanine transaminase levels, Triglyceride levels, Breast milk fatty acid composition (maternal genotype effect), Breast milk fatty acid composition (infant genotype effect), Serum metabolite concentrations in chronic kidney disease, Serum metabolite ratios in chronic kidney disease, Low density lipoprotein cholesterol levels, Apolipoprotein B levels, Pulse pressure, LDL cholesterol levels, LDL cholesterol levels x long total sleep time interaction (2df test), HDL cholesterol levels, Triglyceride levels x short total sleep time interaction (2df test), Serum metabolite levels (CMS), HDL cholesterol levels x long total sleep time interaction (2df test), HDL cholesterol levels x short total sleep time interaction (2df test), HDL cholesterol levels x alcohol consumption (regular vs non-regular drinkers) interaction (2df), Triglyceride levels x alcohol consumption (regular vs non-regular drinkers) interaction (2df), Bipolar I disorder, Apolipoprotein A1 levels, Age-related disease endophenotypes, Red cell distribution width, Height, Blood metabolite levels, Blood metabolite ratios, Rheumatoid arthritis, Serum metabolite levels, Change in serum metabolite levels (CMS), Serum omega-3 polyunsaturated fatty acid concentration in metabolic syndrome, Serum omega-6 to omega-3 polyunsaturated fatty acid ratio in metabolic syndrome, Serum docosahexaenoic fatty acid concentration in metabolic</p>
--	--	--	--	--	--	--	--	---

								<p>syndrome, Delta-5 desaturase activity response to n3-polyunsaturated fat supplement, Change in serum metabolite levels, LDL cholesterol x physical activity interaction (2df test), High density lipoprotein cholesterol levels, Asthma, Bipolar disorder, Osteoporosis-related phenotypes (MTAG), Granulocyte count, Age-related diseases, mortality and associated endophenotypes, Sum neutrophil eosinophil counts, Sum eosinophil basophil counts, Myeloid white cell count, White blood cell count, Mean platelet volume, Trans fatty acid levels, Red blood cell fatty acid levels, Plasma omega-6 polyunsaturated fatty acid levels (arachidonic acid), Plasma omega-6 polyunsaturated fatty acid levels (adrenic acid), Laryngeal squamous cell carcinoma, Triacylglycerol 56:6 levels, Metabolite risk score for predicting weight gain, Fatty acid desaturase activity (serum), Fatty acid desaturase activity (adipose tissue), Phosphatidylcholine-ether levels, Lysophosphatidylethanolamine levels, Phosphatidylcholine levels, Phosphatidylethanolamine-ether levels, Phosphatidylinositol levels, Cholesteryl ester levels, Sphingomyelin levels, Lysophosphatidylcholine levels, Triacylglyceride levels, HDL cholesterol levels in current drinkers, LDL cholesterol levels in current drinkers, Triglyceride levels in current drinkers, HDL cholesterol levels x alcohol consumption (drinkers vs non-drinkers) interaction (2df), Triglyceride levels x alcohol consumption (drinkers vs non-drinkers) interaction (2df), LDL cholesterol levels x alcohol consumption (drinkers vs non-drinkers) interaction (2df), LDL cholesterol levels x alcohol consumption (regular vs non-regular drinkers) interaction (2df), Red blood cell count, QT interval, Aortic valve stenosis, Sleep duration, Neutrophil count, Total triglycerides levels, Mean corpuscular hemoglobin concentration, Mean corpuscular volume, IgA levels, Gallstone disease; ENSG00000221968: Sphingolipid levels, Metabolite levels, Phosphatidylcholine levels, Cholesteryl ester levels, Sphingolipid d18:1/d18:2 ratio, Urate levels; ENSG00000167994: Sphingolipid levels, Metabolite levels, Schizophrenia, Mean corpuscular hemoglobin; ENSG00000167995: Plasma omega-3 polyunsaturated fatty acid level (eicosapentaenoic acid), Metabolite levels; ENSG00000167996: ; ENSG00000149503: Prostate cancer, Breast cancer; ENSG00000168515: ; ENSG00000124939: ; ENSG00000124935: ; ENSG00000110484: ; ENSG00000197745: ; ENSG00000162174: 3-hydroxypropylmercapturic acid levels in smokers, Trunk fat mass, Body fat mass; ENSG00000149021: Waist-to-hip ratio adjusted for BMI, Waist-hip ratio; ENSG00000124942: Alzheimer disease and age of onset, Heel bone mineral density, Alzheimer's disease (late onset), HDL cholesterol levels, Cutaneous melanoma or hair colour, Waist-to-hip ratio adjusted for BMI, Blond vs. brown/black hair color, C-reactive protein levels, Lung function (FEV1),</p>
--	--	--	--	--	--	--	--	--

								Waist-hip ratio, Waist circumference adjusted for body mass index, Hair color
17	79055998	79469799	79255998	1.91E-20	ENSG00000175866, ENSG00000181409, ENSG00000141577, ENSG00000167302, ENSG00000224877, ENSG00000157637, ENSG00000185332	BAIAP2, AATK, AZI1, ENTHD2, C17orf89, SLC38A10, TMEM105	ENSG00000157637: Longevity, IgG glycosylation, Serum 25-Hydroxyvitamin D levels, Menarche (age at onset)	ENSG00000175866: Serum uric acid levels in response to allopurinol in gout, Neuroticism, Neuroticism, Feeling tense, Fat-free mass, Memory performance, Depression, Worry, Depressed affect, Body mass index, General factor of neuroticism, Positive affect, Well-being spectrum (multivariate analysis), Breast cancer, Depressive symptoms (MTAG), Depressive symptoms, General cognitive ability, Sensitivity to environmental stress and adversity, Cognitive performance (MTAG), Self-reported math ability, Self-reported math ability (MTAG); ENSG00000181409: Obesity-related traits, Neuroticism, Feeling tense, Beef consumption, Fat-free mass, Worry, Neuroticism, Body mass index, Depressive symptoms, Neutrophil count, Life satisfaction, Sensitivity to environmental stress and adversity, Height, Systolic blood pressure; ENSG00000141577: IgG glycosylation, Frontotemporal dementia, IgG galactosylation phenotypes (multivariate analysis), IgG N-glycosylation phenotypes (multivariate analysis), Red blood cell count; ENSG00000167302: Blood protein levels; ENSG00000224877: IgG fucosylation phenotypes (multivariate analysis); ENSG00000157637: Longevity, IgG glycosylation, Serum 25-Hydroxyvitamin D levels, Menarche (age at onset); ENSG00000185332: Intake of total sugars
11	60562850	61159993	60762850	4.62E-19	ENSG00000172689, ENSG00000110104, ENSG00000183134, ENSG00000149506, ENSG00000110107, ENSG00000110108, ENSG00000006118, ENSG00000110446, ENSG0000013725, ENSG00000110448, ENSG00000167987, ENSG00000229859, ENSG00000229183, ENSG00000256713, ENSG00000167992, ENSG00000167986, ENSG00000149476, ENSG00000162144,	MS4A10, CCDC86, PTGDR2, ZP1, PRPF19, TMEM109, TMEM132A, SLC15A3, CD6, CD5, VPS37C, PGA3, PGA4, PGA5, VWCE, DDB1, DAK, CYB561A3, TMEM138, TMEM216	ENSG0000013725: Multiple sclerosis, Inflammatory bowel disease, Crohn's disease, Ulcerative colitis, Chronic inflammatory diseases (ankylosing spondylitis, Crohn's disease, psoriasis, primary sclerosing cholangitis, ulcerative colitis) (pleiotropy), CD6 levels	ENSG00000172689: ; ENSG00000110104: ; ENSG00000183134: ; ENSG00000149506: Response to ziprazidone in schizophrenia; ENSG00000110107: ; ENSG00000110108: Blood protein levels; ENSG00000006118: Blood protein levels; ENSG00000110446: ; ENSG0000013725: Multiple sclerosis, Inflammatory bowel disease, Crohn's disease, Ulcerative colitis, Chronic inflammatory diseases (ankylosing spondylitis, Crohn's disease, psoriasis, primary sclerosing cholangitis, ulcerative colitis) (pleiotropy), CD6 levels; ENSG00000110448: ; ENSG00000167987: Rheumatoid arthritis (ACPA-positive), Rheumatoid arthritis, Adolescent idiopathic scoliosis, Periventricular white matter hyperintensities; ENSG00000229859: ; ENSG00000229183: ; ENSG00000256713: Height; ENSG00000167992: ; ENSG00000167986: ; ENSG00000149476: ; ENSG00000162144: ; ENSG00000149483: ; ENSG00000187049:

					ENSG00000149483, ENSG00000187049			
1	147576333	148018267	147805232	6.54E-19	ENSG00000203836, ENSG00000255963, ENSG00000122497	NBPF24, PPIAL4A, NBPF14	ENSG00000255963:	ENSG00000203836: ; ENSG00000255963: ; ENSG00000122497:
4	100054570	100521443	100254570	3.90E-18	ENSG00000198099, ENSG00000172955, ENSG00000187758, ENSG00000196616, ENSG00000196344, ENSG00000138813, ENSG00000145331, ENSG00000138823	ADH4, ADH6, ADH1A, ADH1B, ADH7, C4orf17, TRMT10A, MTTP	ENSG00000196616: Conduct disorder (maternal expressed emotions interaction), Esophageal cancer (alcohol interaction), Oral cavity and pharyngeal cancer, Alcohol dependence, Oropharynx cancer, Oral cavity cancer, Cerebrospinal fluid clusterin levels, Body mass index, Low density lipoprotein cholesterol levels, Risk-taking tendency (4-domain principal component model), LDL cholesterol levels, Relative fat intake, Relative protein intake, Apolipoprotein B levels, HDL cholesterol levels, Esophageal cancer, Alcohol consumption, Serum 25-Hydroxyvitamin D levels, C-reactive protein levels, Heel bone mineral density, Alcohol consumption (drinks per week), Alcohol consumption (drinkers vs non-drinkers), Alcohol consumption (heavy vs. light/non-drinkers), Pulse pressure, Alcohol consumption over the past year, Major depression and alcohol dependence, Lung cancer, Alcohol consumption in current drinkers, Blood urea nitrogen levels, Waist-hip ratio, Maximum habitual alcohol consumption, Regular attendance at a pub or social club, Urinary sodium excretion, Predicted visceral adipose tissue, Total cholesterol levels, Alcohol use disorder, Alcohol use	ENSG00000198099: Esophageal cancer (alcohol interaction), Alcohol dependence, Triglyceride levels, Serum metabolite levels, Blood protein levels, Alcohol use disorder, Platelet count, Eosinophil counts; ENSG00000172955: QRS duration; ENSG00000187758: Urinary metabolite levels in chronic kidney disease, Serum 25-Hydroxyvitamin D levels; ENSG00000196616: Conduct disorder (maternal expressed emotions interaction), Esophageal cancer (alcohol interaction), Oral cavity and pharyngeal cancer, Alcohol dependence, Oropharynx cancer, Oral cavity cancer, Cerebrospinal fluid clusterin levels, Body mass index, Low density lipoprotein cholesterol levels, Risk-taking tendency (4-domain principal component model), LDL cholesterol levels, Relative fat intake, Relative protein intake, Apolipoprotein B levels, HDL cholesterol levels, Esophageal cancer, Alcohol consumption, Serum 25-Hydroxyvitamin D levels, C-reactive protein levels, Heel bone mineral density, Alcohol consumption (drinks per week), Alcohol consumption (drinkers vs non-drinkers), Alcohol consumption (heavy vs. light/non-drinkers), Pulse pressure, Alcohol consumption over the past year, Major depression and alcohol dependence, Lung cancer, Alcohol consumption in current drinkers, Blood urea nitrogen levels, Waist-hip ratio, Maximum habitual alcohol consumption, Regular attendance at a pub or social club, Urinary sodium excretion, Predicted visceral adipose tissue, Total cholesterol levels, Alcohol use disorder, Alcohol use disorder (consumption score), Hemoglobin concentration, Problematic alcohol use, Alcohol consumption (drinks per week) (MTAG), Problematic alcohol use (MTAG), Alcohol use disorder (total score), Alcohol use disorder (dependence and problematic use scores), Bitter alcoholic beverage consumption, Alcohol dependence symptom count, Alcohol dependence (tolerance), Alcohol dependence (desire to cut drinking), Cardiovascular disease, Red blood cell count, Mean corpuscular hemoglobin, Systolic blood pressure; ENSG00000196344: Oral cavity and pharyngeal cancer, Blood protein levels, Maximum habitual alcohol consumption; ENSG00000138813: Metabolite levels (MHPG), Alcohol consumption (drinks per week), Pre-treatment viral load in HIV-1 infection, Height; ENSG00000145331: Gut microbiome composition (summer); ENSG00000138823: Celiac disease,

							disorder (consumption score), Hemoglobin concentration, Problematic alcohol use, Alcohol consumption (drinks per week) (MTAG), Problematic alcohol use (MTAG), Alcohol use disorder (total score), Alcohol use disorder (dependence and problematic use scores), Bitter alcoholic beverage consumption, Alcohol dependence symptom count, Alcohol dependence (tolerance), Alcohol dependence (desire to cut drinking), Cardiovascular disease, Red blood cell count, Mean corpuscular hemoglobin, Systolic blood pressure	Triglyceride levels, Maximum habitual alcohol consumption, HDL cholesterol, Alcohol use disorder, Lung function (FVC)
2	226448377	226848377	226648377	1.06E-17	ENSG00000144460	NYAP2	ENSG00000144460: Neurociticism, Docetaxel-induced peripheral neuropathy in metastatic castrate-resistant prostate cancer, HDL cholesterol levels, Neuroticism, Smoking initiation (ever regular vs never regular), Age of smoking initiation (MTAG), Positive affect, Well-being spectrum (multivariate analysis), Depressive symptoms, Gut microbiota (bacterial taxa, hurdle binary method), Smoking cessation (MTAG), Life satisfaction, Smoking initiation (ever regular vs never regular) (MTAG), Educational attainment (MTAG), Cognitive performance (MTAG), Highest math class taken (MTAG), Smoking status	ENSG00000144460: Neurociticism, Docetaxel-induced peripheral neuropathy in metastatic castrate-resistant prostate cancer, HDL cholesterol levels, Neuroticism, Smoking initiation (ever regular vs never regular), Age of smoking initiation (MTAG), Positive affect, Well-being spectrum (multivariate analysis), Depressive symptoms, Gut microbiota (bacterial taxa, hurdle binary method), Smoking cessation (MTAG), Life satisfaction, Smoking initiation (ever regular vs never regular) (MTAG), Educational attainment (MTAG), Cognitive performance (MTAG), Educational attainment (years of education), Highest math class taken (MTAG), Smoking status
14	73102086	73504565	73302086	2.25E-17	ENSG00000205683, ENSG00000119599, ENSG00000165861	DPF3, DCAF4, ZFYVE1	ENSG00000205683: 3-hydroxy-1-methylpropylmercapturic acid levels in smokers, Disease progression in age-related macular degeneration,	ENSG00000205683: 3-hydroxy-1-methylpropylmercapturic acid levels in smokers, Disease progression in age-related macular degeneration, Pulse pressure, Metabolite levels, Rosacea symptom severity, Adolescent idiopathic scoliosis, C-reactive protein levels, Stem cell factor levels, Renal

							<p>Pulse pressure, Metabolite levels, Rosacea symptom severity, Adolescent idiopathic scoliosis, C-reactive protein levels, Stem cell factor levels, Renal cell carcinoma, Intake of sweets, Body mass index, Atrial fibrillation, Hematocrit, Central corneal thickness, Hemoglobin concentration, Tuberculosis, Systemic lupus erythematosus, Intraocular pressure, Red cell distribution width</p>	<p>cell carcinoma, Intake of sweets, Body mass index, Atrial fibrillation, Hematocrit, Central corneal thickness, Hemoglobin concentration, Tuberculosis, Systemic lupus erythematosus, Intraocular pressure, Red cell distribution width; ENSG00000119599: Leukocyte telomere length, Eosinophil counts, Body mass index, Systolic blood pressure; ENSG00000165861: LDL cholesterol levels, Resistance to Mycobacterium tuberculosis in HIV-positive individuals measured by a negative tuberculin skin test (continuous), Red blood cell count, Body mass index, Intelligence, Mean corpuscular hemoglobin, Red cell distribution width</p>
11	27418490	27933143	27632440	5.74E-17	ENSG00000205213, ENSG00000148943, ENSG00000176697	LGR4, LIN7C, BDNF	<p>ENSG00000176697: Obesity, Body mass index, Smoking behavior, Weight, Childhood body mass index, Menarche (age at onset), Menopause (age at onset), Coronary artery disease, Feeling nervous, Smoking initiation, Snoring, Risk-taking tendency (4-domain principal component model), General risk tolerance (MTAG), Triglyceride levels, Smoking status (ever vs never smokers), Body fat percentage, Fat-free mass, Diastolic blood pressure, Body mass index (joint analysis main effects and physical activity interaction), Coffee consumption, Worry, C-reactive protein levels, Systolic blood pressure, Hip circumference, Body mass index in physically active individuals, Body mass index in physically inactive individuals, Body mass index (SNP x SNP interaction), BMI (adjusted for smoking behaviour), BMI in non-smokers, BMI in smokers, Body mass index (joint analysis main effects and smoking interaction), Waist-hip ratio, Smoking initiation (ever regular vs</p>	<p>ENSG00000205213: Male-pattern baldness, Heel bone mineral density, Spontaneous preterm birth with premature rupture of membranes, Urate levels in obese individuals, Blond vs. brown/black hair color, Waist-hip ratio, Body mass index, Balding type 1, Lung function (FVC), Hair color; ENSG00000148943: Feeling miserable, Hip circumference, Survival in pancreatic cancer, Heel bone mineral density, Body mass index; ENSG00000176697: Obesity, Body mass index, Smoking behavior, Weight, Childhood body mass index, Menarche (age at onset), Menopause (age at onset), Coronary artery disease, Feeling nervous, Smoking initiation, Snoring, Risk-taking tendency (4-domain principal component model), General risk tolerance (MTAG), Triglyceride levels, Smoking status (ever vs never smokers), Body fat percentage, Fat-free mass, Diastolic blood pressure, Body mass index (joint analysis main effects and physical activity interaction), Coffee consumption, Worry, C-reactive protein levels, Systolic blood pressure, Hip circumference, Body mass index in physically active individuals, Body mass index in physically inactive individuals, Body mass index (SNP x SNP interaction), BMI (adjusted for smoking behaviour), BMI in non-smokers, BMI in smokers, Body mass index (joint analysis main effects and smoking interaction), Waist-hip ratio, Smoking initiation (ever regular vs never regular), Metabolic syndrome, Predicted visceral adipose tissue, Basal metabolic rate variance, Body mass index variance, Basal metabolic rate, Hand grip strength, Type 2 diabetes, Problematic alcohol use (MTAG), Diverticular disease, Alcohol consumption, Smoking cessation (MTAG), Age of smoking initiation (MTAG), Cigarettes smoked per day (MTAG), Smoking initiation (ever regular vs never regular) (MTAG), Educational attainment (years of education), Heel bone mineral density, Educational attainment (MTAG), Highest math class taken, Highest math class taken (MTAG), Smoking status</p>

							never regular), Metabolic syndrome, Predicted visceral adipose tissue, Basal metabolic rate variance, Body mass index variance, Basal metabolic rate, Hand grip strength, Type 2 diabetes, Problematic alcohol use (MTAG), Diverticular disease, Alcohol consumption, Smoking cessation (MTAG), Age of smoking initiation (MTAG), Cigarettes smoked per day (MTAG), Smoking initiation (ever regular vs never regular) (MTAG), Educational attainment (years of education), Heel bone mineral density, Educational attainment (MTAG), Highest math class taken, Highest math class taken (MTAG), Smoking status	
1	110586660	111157709	110793599	8.40E-17	ENSG00000143093, ENSG00000156150, ENSG00000186150, ENSG00000197106, ENSG00000116396, ENSG00000162775, ENSG00000168679, ENSG00000134248, ENSG00000143125, ENSG00000143105, ENSG00000177301	STRIP1, ALX3, UBL4B, SLC6A17, KCNC4, RBM15, SLC16A4, LAMTOR5, PROK1, KCNA10, KCNA2	ENSG00000116396: Self-rated health, Oropharynx cancer, Gut microbiota (functional units), Red blood cell count, Educational attainment (MTAG), Cognitive performance (MTAG), Self-reported math ability, Self-reported math ability (MTAG), Educational attainment (years of education), Highest math class taken (MTAG)	ENSG00000143093: Optic disc size, Cognitive performance (MTAG), Educational attainment (years of education); ENSG00000156150: Optic disc size; ENSG00000186150: ; ENSG00000197106: Keratinocyte cancer (MTAG), Basal cell carcinoma, Educational attainment (years of education), Sunburns; ENSG00000116396: Self-rated health, Oropharynx cancer, Gut microbiota (functional units), Red blood cell count, Educational attainment (MTAG), Cognitive performance (MTAG), Self-reported math ability, Self-reported math ability (MTAG), Educational attainment (years of education), Highest math class taken (MTAG); ENSG00000162775: ; ENSG00000168679: Adolescent idiopathic scoliosis; ENSG00000134248: ; ENSG00000143125: ; ENSG00000143105: Bitter taste perception (phenylthiocarbamide) in obesity with metabolic syndrome; ENSG00000177301: Nonsyndromic cleft lip with cleft palate
3	98038565	98438565	98238565	8.59E-17	ENSG00000196098, ENSG00000206536, ENSG00000232382, ENSG00000231861, ENSG00000080822, ENSG00000080819, ENSG00000154165	OR5K4, OR5K3, OR5K1, OR5K2, CLDND1, CPOX, GPR15	ENSG00000080822: Post bronchodilator FEV1/FVC ratio	ENSG00000196098: ; ENSG00000206536: ; ENSG00000232382: ; ENSG00000231861: ; ENSG00000080822: Post bronchodilator FEV1/FVC ratio; ENSG00000080819: ; ENSG00000154165:
17	8655348	9223981	8875399	9.12E-17	ENSG00000183318, ENSG00000185156,	SPDYE4, MFS6L,	ENSG00000141506: 3-hydroxypropylmercapturic acid	ENSG00000183318: ; ENSG00000185156: ; ENSG00000174083: ; ENSG00000141506: 3-hydroxypropylmercapturic acid levels in smokers,

					ENSG00000174083, ENSG00000141506, ENSG00000065320, ENSG00000170310	PIK3R6, PIK3R5, NTN1, STX8	levels in smokers, Metabolite levels, Eosinophil counts, Autoimmune thyroid disease, Atypical femoral fracture in phosphonate treatment, Hypothyroidism	Metabolite levels, Eosinophil counts, Autoimmune thyroid disease, Atypical femoral fracture in phosphonate treatment, Hypothyroidism; ENSG00000065320: Orofacial clefts, Breast cancer (prognosis), Neuroticism, Percentage gas trapping, 3-hydroxypropylmercapturic acid levels in smokers, Lobe attachment (rater-scored or self-reported), Heel bone mineral density, Feeling hurt, Blood protein levels, Nonsyndromic cleft lip with or without cleft palate, Cleft lip with or without cleft palate, Colonoscopy-negative controls vs population controls, Rostral middle frontal gyrus volume, Total grey matter volume, Cortex volume, White matter microstructure (fractional anisotropy), Bisphosphonate-associated atypical femoral fracture, Sensitivity to environmental stress and adversity; ENSG00000170310: antipsychotic drug dosage in schizophrenia or schizoaffective disorder, HDL cholesterol change in response to fenofibrate in statin-treated type 2 diabetes, Diabetic kidney disease, Obstructive sleep apnea trait (apnea hypopnea index), Borderline personality disorder, White matter microstructure (mode of anisotropy), Gut microbiota (bacterial taxa, hurdle binary method), Blood protein levels, Height
8	133319098	133719098	133519098	3.22E-16	ENSG00000184156, ENSG00000129295, ENSG00000165071	KCNQ3, LRR6, TMEM71	ENSG00000129295: Adolescent idiopathic scoliosis, Cognitive performance (MTAG)	ENSG00000184156: Coronary artery calcification, LDL peak particle diameter (total fat intake interaction), QT interval, Interleukin-8 levels, Spontaneous adipocyte lipolysis, Height, Stimulated adipocyte lipolysis, Major depressive disorder, Rate of cognitive decline in Alzheimer's disease, Type 2 diabetes; ENSG00000129295: Adolescent idiopathic scoliosis, Cognitive performance (MTAG); ENSG00000165071: Thyroid stimulating hormone levels, Hyperthyroidism, Smoking initiation (ever regular vs never regular) (MTAG), Age of smoking initiation (MTAG), Educational attainment (MTAG), Educational attainment (years of education), Height
3	100737908	101365506	101082807	6.97E-16	ENSG00000081148, ENSG00000138468, ENSG00000174173, ENSG00000081154	IMPG2, SENP7, TRMT10C, PCNP	ENSG00000138468: Risk-taking tendency (4-domain principal component model), Mosaic loss of chromosome Y (Y chromosome dosage), Appendicular lean mass, Mean corpuscular hemoglobin, Alcohol consumption (drinks per week), Diastolic blood pressure, Red blood cell count, Platelet count, Alzheimer's disease (late onset), Bitter alcoholic beverage	ENSG00000081148: Inflammatory bowel disease, Crohn's disease, Mean corpuscular hemoglobin; ENSG00000138468: Risk-taking tendency (4-domain principal component model), Mosaic loss of chromosome Y (Y chromosome dosage), Appendicular lean mass, Mean corpuscular hemoglobin, Alcohol consumption (drinks per week), Diastolic blood pressure, Red blood cell count, Platelet count, Alzheimer's disease (late onset), Bitter alcoholic beverage consumption, Mean corpuscular volume; ENSG00000174173: Eosinophil counts; ENSG00000081154:

							consumption, Mean corpuscular volume	
18	46132358	46550865	46350865	7.05E-16	ENSG00000134030, ENSG00000101665	CTIF, SMAD7	ENSG00000101665: Colorectal cancer, Hematocrit, Creatinine levels, Glomerular filtration rate, Normal facial asymmetry (angle of surface orientation score), Heel bone mineral density, Red blood cell count, Hemoglobin concentration, Hemoglobin levels, Parental longevity (at least one long-lived parent), Colorectal cancer or advanced adenoma, Estimated glomerular filtration rate, Atrial fibrillation, Eosinophil counts, Asthma, Mean corpuscular hemoglobin, Mean corpuscular volume, Red cell distribution width, Eczema	ENSG00000134030: IgG glycosylation, Lung function (FVC), Response to carboplatin in ovarian cancer (MTT IC50), Male-pattern baldness, Red blood cell count, Adolescent idiopathic scoliosis, Automobile speeding propensity, Circulating odd-numbered chain saturated fatty acid levels (C15:0), Migraine with aura, White blood cell count, Iris heterochromicity, Lumbar spine bone mineral density (integral), Rate of cognitive decline in Alzheimer's disease, Height, Smoking status; ENSG00000101665: Colorectal cancer, Hematocrit, Creatinine levels, Glomerular filtration rate, Normal facial asymmetry (angle of surface orientation score), Heel bone mineral density, Red blood cell count, Hemoglobin concentration, Hemoglobin levels, Parental longevity (at least one long-lived parent), Colorectal cancer or advanced adenoma, Estimated glomerular filtration rate, Atrial fibrillation, Eosinophil counts, Asthma, Mean corpuscular hemoglobin, Mean corpuscular volume, Red cell distribution width, Eczema
22	31356103	31756103	31556103	1.03E-15	ENSG00000133422, ENSG00000183963, ENSG00000185133, ENSG00000100078, ENSG00000138942, ENSG00000182541, ENSG00000100100, ENSG00000100105	MORC2, SMTN, INPP5J, PLA2G3, RNF185, LIMK2, PIK3IP1, PATZ1	ENSG00000100078: Serum 25-Hydroxyvitamin D levels	ENSG00000133422: ; ENSG00000183963: Lung function (FEV1/FVC); ENSG00000185133: Gut microbiota (beta diversity); ENSG00000100078: Serum 25-Hydroxyvitamin D levels; ENSG00000138942: Apolipoprotein A1 levels, Waist circumference, Waist-to-hip ratio adjusted for BMI, Neutrophil count, Waist circumference adjusted for body mass index, Waist-hip ratio; ENSG00000182541: Type 2 diabetes nephropathy, Eosinophil percentage of white cells, Eosinophil counts, Paclitaxel-induced neuropathy, Multiple sclerosis, Sum eosinophil basophil counts, Eosinophil percentage of granulocytes, Neutrophil percentage of granulocytes, Schizophrenia, Metabolite levels, Reaction time; ENSG00000100100: ; ENSG00000100105: Height, Reticulocyte fraction of red cells, Reticulocyte count, Red cell distribution width, Eosinophil counts
3	101491814	101891814	101691814	1.40E-15	ENSG00000144815, ENSG00000144802, ENSG00000170044	NXPE3, NFKBIZ, ZPLD1	ENSG00000144802: Ulcerative colitis, Colorectal cancer, Spatial memory	ENSG00000144815: ; ENSG00000144802: Ulcerative colitis, Colorectal cancer, Spatial memory; ENSG00000170044: Post bronchodilator FEV1/FVC ratio, Adolescent idiopathic scoliosis, C-reactive protein levels, Waist-to-hip circumference ratio (alcohol intake interaction), Response to ranibizumab in age-related macular degeneration (exudative)
12	27574693	27974693	27774693	1.53E-15	ENSG00000029153, ENSG00000165935, ENSG00000110841, ENSG00000174236,	ARNTL2, SMCO2, PPFIBP1, REP15,	ENSG00000110841: Cerebral cortical growth, Tinnitus in cisplatin-treated testicular cancer	ENSG00000029153: Body fat percentage, Waist-to-hip ratio adjusted for BMI, Waist-hip ratio; ENSG00000165935: Obstructive sleep apnea trait (apnea hypopnea index), Lung cancer in ever smokers; ENSG00000110841: Cerebral cortical growth, Tinnitus in cisplatin-treated

					ENSG0000061794, ENSG00000205693, ENSG0000087448	MRPS35, MANSC4, KLHL42		testicular cancer; ENSG0000174236: ; ENSG0000061794: Type 2 diabetes, Type 2 diabetes (adjusted for BMI), Lung function (FEV1/FVC); ENSG00000205693: Blood protein levels; ENSG0000087448: Pulse pressure, Resting heart rate
--	--	--	--	--	---	------------------------------	--	---

4608

4609 **Table S4d.3 West Eurasia.cw_hg**

4610

4611
4612

4613 References

- 4614 1. Luu, K., Bazin, E. & Blum, M. G. B. pcadapt: an R package to perform genome scans
4615 for selection based on principal component analysis. *Mol. Ecol. Resour.* **17**, 67–77
4616 (2017).
- 4617 2. Durinck, S., Spellman, P. T., Birney, E. & Huber, W. Mapping identifiers for the
4618 integration of genomic datasets with the R/Bioconductor package biomaRt. *Nat. Protoc.*
4619 **4**, 1184–1191 (2009).
- 4620 3. Wilde, S. *et al.* Direct evidence for positive selection of skin, hair, and eye
4621 pigmentation in Europeans during the last 5,000 y. *Proc. Natl. Acad. Sci. U. S. A.* **111**,
4622 4832–4837 (2014).
- 4623 4. Pickrell, J. K. *et al.* Signals of recent positive selection in a worldwide sample of
4624 human populations. *Genome Research* vol. 19 826–837 (2009).
- 4625 5. Lao, O., de Gruijter, J. M., van Duijn, K., Navarro, A. & Kayser, M. Signatures of
4626 positive selection in genes associated with human skin pigmentation as revealed from
4627 analyses of single nucleotide polymorphisms. *Ann. Hum. Genet.* **71**, 354–369 (2007).
- 4628 6. Herchuelz, A. *Sodium-calcium exchange and the plasma membrane Ca²⁺-ATPase*
4629 *in cell function: fifth international conference.* (Wiley-Blackwell, 2007).
- 4630 7. Bryk, J. *et al.* Positive Selection in East Asians for an EDAR Allele that Enhances
4631 NF-κB Activation. *PLoS ONE* vol. 3 e2209 (2008).
- 4632 8. Hider, J. L. *et al.* Exploring signatures of positive selection in pigmentation candidate
4633 genes in populations of East Asian ancestry. *BMC Evol. Biol.* **13**, 150 (2013).
- 4634 9. Duforet-Frebourg, N., Luu, K., Laval, G., Bazin, E. & Blum, M. G. B. Detecting
4635 Genomic Signatures of Natural Selection with Principal Component Analysis:
4636 Application to the 1000 Genomes Data. *Mol. Biol. Evol.* **33**, 1082–1093 (2016).
- 4637 10. Prashanth, S. & Deshmukh, S. Ectodermal Dysplasia: A Genetic Review.
4638 *International Journal of Clinical Pediatric Dentistry* vol. 5 197–202 (2012).

- 4639 11. Kimura, R. *et al.* A common variation in EDAR is a genetic determinant of
4640 shovel-shaped incisors. *Am. J. Hum. Genet.* **85**, 528–535 (2009).
- 4641 12. Izawa, T. *et al.* ASXL2 Regulates Glucose, Lipid, and Skeletal Homeostasis.
4642 *Cell Rep.* **11**, 1625–1637 (2015).
- 4643 13. Zou, W. *et al.* Myeloid-specific Asxl2 deletion limits diet-induced obesity by
4644 regulating energy expenditure. *J. Clin. Invest.* **130**, 2644–2656 (2020).
- 4645 14. Park, U.-H., Yoon, S. K., Park, T., Kim, E.-J. & Um, S.-J. Additional Sex
4646 Comb-like (ASXL) Proteins 1 and 2 Play Opposite Roles in Adipogenesis via Reciprocal
4647 Regulation of Peroxisome Proliferator-activated Receptor γ . *Journal of Biological*
4648 *Chemistry* vol. 286 1354–1363 (2011).
- 4649 15. Ponsuksili, S. *et al.* Epigenome-wide skeletal muscle DNA methylation
4650 profiles at the background of distinct metabolic types and ryanodine receptor variation
4651 in pigs. *BMC Genomics* **20**, 492 (2019).
- 4652 16. Samad, M. B. *et al.* [6]-Gingerol, from *Zingiber officinale*, potentiates GLP-1
4653 mediated glucose-stimulated insulin secretion pathway in pancreatic β -cells and
4654 increases RAB8/RAB10-regulated membrane presentation of GLUT4 transporters in
4655 skeletal muscle to improve hyperglycemia in *Lepr^{db/db}* type 2 diabetic mice. *BMC*
4656 *Complementary and Alternative Medicine* vol. 17 (2017).
- 4657 17. Vazirani, R. P. *et al.* Disruption of Adipose Rab10-Dependent Insulin
4658 Signaling Causes Hepatic Insulin Resistance. *Diabetes* **65**, 1577–1589 (2016).
- 4659 18. Hsieh, P. *et al.* Exome Sequencing Provides Evidence of Polygenic
4660 Adaptation to a Fat-Rich Animal Diet in Indigenous Siberian Populations. *Mol. Biol.*
4661 *Evol.* **34**, 2913–2926 (2017).
- 4662 19. Thapa, D. *et al.* The protein acetylase GCN5L1 modulates hepatic fatty acid
4663 oxidation activity via acetylation of the mitochondrial β -oxidation enzyme HADHA. *J.*
4664 *Biol. Chem.* **293**, 17676–17684 (2018).
- 4665 20. Baloni, P. *et al.* Genome-scale metabolic model of the rat liver predicts effects
4666 of diet restriction. *Sci. Rep.* **9**, 9807 (2019).

- 4667 21. Ong, H. S. & Yim, H. C. H. Microbial Factors in Inflammatory Diseases and
4668 Cancers. *Regulation of Inflammatory Signaling in Health and Disease* 153–174 (2017)
4669 doi:10.1007/978-981-10-5987-2_7.
- 4670 22. Logsdon, B. A., Hoffman, G. E. & Mezey, J. G. Mouse obesity network
4671 reconstruction with a variational Bayes algorithm to employ aggressive false positive
4672 control. *BMC Bioinformatics* **13**, 53 (2012).
- 4673 23. Pei, Y.-F. *et al.* Genomic variants at 20p11 associated with body fat mass in
4674 the European population. *Obesity* **25**, 757–764 (2017).
- 4675 24. Sabir, J. S. M. *et al.* Unraveling the role of salt-sensitivity genes in obesity
4676 with integrated network biology and co-expression analysis. *PLoS One* **15**, e0228400
4677 (2020).
- 4678 25. Wu, H. *et al.* Transcriptome Sequencing to Detect the Potential Role of Long
4679 Noncoding RNAs in Salt-Sensitive Hypertensive Rats. *Biomed Res. Int.* **2019**, 2816959
4680 (2019).
- 4681 26. Wang, L. *et al.* Peakwide Mapping on Chromosome 3q13 Identifies the Kalirin
4682 Gene as a Novel Candidate Gene for Coronary Artery Disease. *The American Journal*
4683 *of Human Genetics* vol. 80 650–663 (2007).
- 4684 27. Ikram, M. A., Seshadri, S. & Bis, J. C. Genomewide Association Studies of
4685 Stroke. *Journal of Vascular Surgery* vol. 50 467 (2009).
- 4686 28. Krug, T. *et al.* Kalirin: a novel genetic risk factor for ischemic stroke. *Hum.*
4687 *Genet.* **127**, 513–523 (2010).
- 4688 29. Zang, X.-L. *et al.* Association of a SNP in SLC35F3 Gene with the Risk of
4689 Hypertension in a Chinese Han Population. *Frontiers in Genetics* vol. 7 (2016).
- 4690 30. Zhang, K. *et al.* Genetic implication of a novel thiamine transporter in human
4691 hypertension. *J. Am. Coll. Cardiol.* **63**, 1542–1555 (2014).
- 4692 31. Pacheu-Grau, D. *et al.* COA6 Facilitates Cytochrome c Oxidase Biogenesis
4693 as Thiol-reductase for Copper Metallochaperones in Mitochondria. *J. Mol. Biol.* **432**,
4694 2067–2079 (2020).

- 4695 32. Russo, L. *et al.* Cholesterol 25-hydroxylase (CH25H) as a promoter of
4696 adipose tissue inflammation in obesity and diabetes. *Mol Metab* **39**, 100983 (2020).
- 4697 33. Zhao, J., Chen, J., Li, M., Chen, M. & Sun, C. Multifaceted Functions of
4698 CH25H and 25HC to Modulate the Lipid Metabolism, Immune Responses, and Broadly
4699 Antiviral Activities. *Viruses* **12**, (2020).
- 4700 34. Demir, A., Kahraman, R., Candan, G. & Ergen, A. The role of FAS gene
4701 variants in inflammatory bowel disease. *Turk. J. Gastroenterol.* **31**, 356–361 (2020).
- 4702 35. Rieux-Laucat, F., Magérus-Chatinet, A. & Neven, B. The Autoimmune
4703 Lymphoproliferative Syndrome with Defective FAS or FAS-Ligand Functions. *Journal of*
4704 *Clinical Immunology* vol. 38 558–568 (2018).
- 4705 36. Karis, K. *et al.* Altered Expression Profile of IgLON Family of Neural Cell
4706 Adhesion Molecules in the Dorsolateral Prefrontal Cortex of Schizophrenic Patients.
4707 *Front. Mol. Neurosci.* **11**, 8 (2018).
- 4708 37. Liu, J., Li, M. & Su, B. GWAS-identified schizophrenia risk SNPs
4709 atTSPAN18 are highly diverged between Europeans and East Asians. *American Journal*
4710 *of Medical Genetics Part B: Neuropsychiatric Genetics* vol. 171 1032–1040 (2016).
- 4711 38. Fu, H.-Y. *et al.* The mutation spectrum of the SLC25A13 gene in Chinese
4712 infants with intrahepatic cholestasis and aminoacidemia. *J. Gastroenterol.* **46**, 510–518
4713 (2011).
- 4714 39. Chen, J.-L., Zhang, Z.-H., Li, B.-X., Cai, Z. & Zhou, Q.-H. Bioinformatic and
4715 functional analysis of promoter region of human SLC25A13 gene. *Gene* **693**, 69–75
4716 (2019).
- 4717 40. Vitoria, I. *et al.* Citrin deficiency in a Romanian child living in Spain highlights
4718 the worldwide distribution of this defect and illustrates the value of nutritional therapy.
4719 *Molecular Genetics and Metabolism* vol. 110 181–183 (2013).
- 4720 41. Fiermonte, G. *et al.* An adult with type 2 citrullinemia presenting in Europe. *N.*
4721 *Engl. J. Med.* **358**, 1408–1409 (2008).
- 4722 42. Sarkar, A. & Nandineni, M. R. Association of common genetic variants with

- 4723 human skin color variation in Indian populations. *Am. J. Hum. Biol.* **30**, (2018).
- 4724 43. Edwards, M. *et al.* Association of the OCA2 polymorphism His615Arg with
4725 melanin content in east Asian populations: further evidence of convergent evolution of
4726 skin pigmentation. *PLoS Genet.* **6**, e1000867 (2010).
- 4727 44. Eiberg, H. *et al.* Blue eye color in humans may be caused by a perfectly
4728 associated founder mutation in a regulatory element located within the HERC2 gene
4729 inhibiting OCA2 expression. *Human Genetics* vol. 123 177–187 (2008).
- 4730 45. Sturm, R. A. *et al.* A Single SNP in an Evolutionary Conserved Region within
4731 Intron 86 of the HERC2 Gene Determines Human Blue-Brown Eye Color. *The*
4732 *American Journal of Human Genetics* vol. 82 424–431 (2008).
- 4733 46. Allentoft, M. E. *et al.* Population genomics of Bronze Age Eurasia. *Nature*
4734 **522**, 167–172 (2015).
- 4735 47. Mathieson, I. *et al.* Genome-wide patterns of selection in 230 ancient
4736 Eurasians. *Nature* **528**, 499–503 (2015).
- 4737 48. Barreiro, L. B. *et al.* Evolutionary dynamics of human Toll-like receptors and
4738 their different contributions to host defense. *PLoS Genet.* **5**, e1000562 (2009).
- 4739 49. Astiz, M. & Oster, H. GLUT12-A promising new target for the treatment of
4740 insulin resistance in obesity and type 2 diabetes. *Acta physiologica* vol. 226 e13329
4741 (2019).
- 4742 50. Waller, A. P. *et al.* GLUT12 functions as a basal and insulin-independent
4743 glucose transporter in the heart. *Biochim. Biophys. Acta* **1832**, 121–127 (2013).
- 4744 51. Joehanes, R. *et al.* Integrated genome-wide analysis of expression
4745 quantitative trait loci aids interpretation of genomic association studies. *Genome Biol.*
4746 **18**, 16 (2017).
- 4747 52. CARDIoGRAMplusC4D Consortium *et al.* Large-scale association analysis
4748 identifies new risk loci for coronary artery disease. *Nat. Genet.* **45**, 25–33 (2013).
- 4749 53. Limaye, N. *et al.* Somatic mutations in angiotensin receptor gene TEK cause
4750 solitary and multiple sporadic venous malformations. *Nat. Genet.* **41**, 118–124 (2009).

- 4751 54. Jones, N., Iljin, K., Dumont, D. J. & Alitalo, K. Tie receptors: new modulators
4752 of angiogenic and lymphangiogenic responses. *Nature Reviews Molecular Cell Biology*
4753 vol. 2 257–267 (2001).
- 4754 55. Dumont, D. J. *et al.* Dominant-negative and targeted null mutations in the
4755 endothelial receptor tyrosine kinase, tek, reveal a critical role in vasculogenesis of the
4756 embryo. *Genes Dev.* **8**, 1897–1909 (1994).
- 4757 56. Gál, Z. *et al.* Investigation of the Possible Role of Tie2 Pathway and TEK
4758 Gene in Asthma and Allergic Conjunctivitis. *Frontiers in Genetics* vol. 11 (2020).
- 4759 57. Cmejla, R. *et al.* Human MRCKalpha is regulated by cellular iron levels and
4760 interferes with transferrin iron uptake. *Biochem. Biophys. Res. Commun.* **395**, 163–167
4761 (2010).
- 4762 58. Richard, C. *et al.* Myotonic dystrophy kinase-related CDC42-binding kinase α ,
4763 a new transferrin receptor type 2-binding partner, is a regulator of erythropoiesis. *Am. J.*
4764 *Hematol.* (2021) doi:10.1002/ajh.26104.
- 4765 59. Molineros, J. E. *et al.* Mechanistic Characterization of RASGRP1 Variants
4766 Identifies an hnRNP-K-Regulated Transcriptional Enhancer Contributing to SLE
4767 Susceptibility. *Frontiers in Immunology* vol. 10 (2019).
- 4768 60. Potier, M. L. *et al.* RasGRP1 and RasGRP3 expression in lymphocytes of
4769 rheumatoid arthritis patients. *Annals of the Rheumatic Diseases* vol. 71 A54.2–A54
4770 (2012).
- 4771 61. Somekh, I. *et al.* Correction to: Novel Mutations in RASGRP1 Are Associated
4772 with Immunodeficiency, Immune Dysregulation, and EBV-Induced Lymphoma. *J. Clin.*
4773 *Immunol.* **38**, 711 (2018).
- 4774 62. Karlas, A. *et al.* Genome-wide RNAi screen identifies human host factors
4775 crucial for influenza virus replication. *Nature* **463**, 818–822 (2010).
- 4776 63. Hao, L. *et al.* Drosophila RNAi screen identifies host genes important for
4777 influenza virus replication. *Nature* **454**, 890–893 (2008).
- 4778 64. Mariniello, B. *et al.* Analysis of the 11 β -Hydroxysteroid Dehydrogenase Type

- 4779 2 Gene (HSD11B2) in Human Essential Hypertension*. *Am. J. Hypertens.* **18**, 1091–
4780 1098 (2005).
- 4781 65. Zeller, T. *et al.* Transcriptome-Wide Analysis Identifies Novel Associations
4782 With Blood Pressure. *Hypertension* **70**, 743–750 (2017).
- 4783 66. Gunaratnam, K., Vidal, C., Gimble, J. M. & Duque, G. Mechanisms of
4784 palmitate-induced lipotoxicity in human osteoblasts. *Endocrinology* **155**, 108–116
4785 (2014).
- 4786 67. Lumish, H. S., O'Reilly, M. & Reilly, M. P. Sex Differences in Genomic Drivers
4787 of Adipose Distribution and Related Cardiometabolic Disorders: Opportunities for
4788 Precision Medicine. *Arterioscler. Thromb. Vasc. Biol.* **40**, 45–60 (2020).
- 4789 68. Ng, M. C. Y. *et al.* Discovery and fine-mapping of adiposity loci using high
4790 density imputation of genome-wide association studies in individuals of African
4791 ancestry: African Ancestry Anthropometry Genetics Consortium. *PLoS Genet.* **13**,
4792 e1006719 (2017).
- 4793 69. Ueda, K. *et al.* Renal Dysfunction Induced by Kidney-Specific Gene Deletion
4794 of *as* as a Primary Cause of Salt-Dependent Hypertension. *Hypertension* **70**, 111–118
4795 (2017).
- 4796 70. Agarwal, A. K. *et al.* CA-Repeat polymorphism in intron 1 of HSD11B2 :
4797 effects on gene expression and salt sensitivity. *Hypertension* **36**, 187–194 (2000).
- 4798 71. Buckley, M. T. *et al.* Selection in Europeans on Fatty Acid Desaturases
4799 Associated with Dietary Changes. *Mol. Biol. Evol.* **34**, 1307–1318 (2017).
- 4800 72. Ye, K., Gao, F., Wang, D., Bar-Yosef, O. & Keinan, A. Dietary adaptation of
4801 FADS genes in Europe varied across time and geography. *Nat Ecol Evol* **1**, 167 (2017).
- 4802 73. Mathieson, S. & Mathieson, I. FADS1 and the Timing of Human Adaptation to
4803 Agriculture. *Mol. Biol. Evol.* **35**, 2957–2970 (2018).
- 4804 74. Pan, G. *et al.* PATZ1 down-regulates FADS1 by binding to rs174557 and is
4805 opposed by SP1/SREBP1c. *Nucleic Acids Res.* **45**, 2408–2422 (2017).
- 4806 75. Tabur, S. *et al.* Evidence for elevated (LIMK2 and CFL1) and suppressed

- 4807 (ICAM1, EZR, MAP2K2, and NOS3) gene expressions in metabolic syndrome.
4808 *Endocrine* **53**, 465–470 (2016).
- 4809 76. Fairn, G. D. & McMaster, C. R. Emerging roles of the oxysterol-binding
4810 protein family in metabolism, transport, and signaling. *Cell. Mol. Life Sci.* **65**, 228–236
4811 (2008).
- 4812 77. Lehto, M. & Olkkonen, V. M. The OSBP-related proteins: a novel protein
4813 family involved in vesicle transport, cellular lipid metabolism, and cell signalling.
4814 *Biochim. Biophys. Acta* **1631**, 1–11 (2003).
- 4815 78. Sánchez-Solana, B., Li, D.-Q. & Kumar, R. Cytosolic functions of MORC2 in
4816 lipogenesis and adipogenesis. *Biochim. Biophys. Acta* **1843**, 316–326 (2014).
- 4817 79. Sartipy, P., Camejo, G., Svensson, L. & Hurt-Camejo, E. Phospholipase A2
4818 modification of lipoproteins: potential effects on atherogenesis. *Adv. Exp. Med. Biol.*
4819 **507**, 3–7 (2002).
- 4820 80. Dennis, E. A., Cao, J., Hsu, Y.-H., Magrioti, V. & Kokotos, G. Phospholipase
4821 A2 Enzymes: Physical Structure, Biological Function, Disease Implication, Chemical
4822 Inhibition, and Therapeutic Intervention. *Chem. Rev.* **111**, 6130 (2011).
- 4823 81. Wyles, J. P., Perry, R. J. & Ridgway, N. D. Characterization of the sterol-
4824 binding domain of oxysterol-binding protein (OSBP)-related protein 4 reveals a novel
4825 role in vimentin organization. *Experimental Cell Research* vol. 313 1426–1437 (2007).
- 4826 82. Tintle, N. L. *et al.* A genome-wide association study of saturated, mono- and
4827 polyunsaturated red blood cell fatty acids in the Framingham Heart Offspring Study.
4828 *Prostaglandins Leukot. Essent. Fatty Acids* **94**, 65–72 (2015).
- 4829 83. Thomas, C. *et al.* LPCAT3 deficiency in hematopoietic cells alters cholesterol
4830 and phospholipid homeostasis and promotes atherosclerosis. *Atherosclerosis* **275**,
4831 409–418 (2018).
- 4832 84. Liu, N. *et al.* Hyperuricemia induces lipid disturbances mediated by LPCAT3
4833 upregulation in the liver. *FASEB J.* (2020) doi:10.1096/fj.202000950R.
- 4834 85. Godahewa, G. I., Bathige, S. D. N. K., Herath, H. M. L. P. B., Noh, J. K. &

4835 Lee, J. Characterization of rock bream (*Oplegnathus fasciatus*) complement
4836 components C1r and C1s in terms of molecular aspects, genomic modulation, and
4837 immune responsive transcriptional profiles following bacterial and viral pathogen
4838 exposure. *Fish Shellfish Immunol.* **46**, 656–668 (2015).

4839 86. Dai, D.-F. *et al.* Plasma concentration of SCUBE1, a novel platelet protein, is
4840 elevated in patients with acute coronary syndrome and ischemic stroke. *J. Am. Coll.*
4841 *Cardiol.* **51**, 2173–2180 (2008).

4842 87. Huang, X., Liu, G., Guo, J. & Su, Z. The PI3K/AKT pathway in obesity and
4843 type 2 diabetes. *Int. J. Biol. Sci.* **14**, 1483–1496 (2018).

4844 88. Zhong, X. *et al.* LNK deficiency decreases obesity-induced insulin resistance
4845 by regulating GLUT4 through the PI3K-Akt-AS160 pathway in adipose tissue. *Aging*
4846 **12**, 17150–17166 (2020).

4847 89. López-Gómez, C. *et al.* Oleic Acid Protects Against Insulin Resistance by
4848 Regulating the Genes Related to the PI3K Signaling Pathway. *J. Clin. Med. Res.* **9**,
4849 (2020).

4850 90. Sandrini, L. *et al.* Association between Obesity and Circulating Brain-Derived
4851 Neurotrophic Factor (BDNF) Levels: Systematic Review of Literature and Meta-
4852 Analysis. *Int. J. Mol. Sci.* **19**, (2018).

4853 91. Lommatzsch, M. *et al.* The impact of age, weight and gender on BDNF levels
4854 in human platelets and plasma. *Neurobiol. Aging* **26**, 115–123 (2005).

4855 92. Tsuchida, A. *et al.* Brain-derived neurotrophic factor ameliorates lipid
4856 metabolism in diabetic mice. *Diabetes Obes. Metab.* **4**, 262–269 (2002).

4857 93. Wu, A., Molteni, R., Ying, Z. & Gomez-Pinilla, F. A saturated-fat diet
4858 aggravates the outcome of traumatic brain injury on hippocampal plasticity and
4859 cognitive function by reducing brain-derived neurotrophic factor. *Neuroscience* **119**,
4860 365–375 (2003).

4861 94. Molteni, R., Barnard, R. J., Ying, Z., Roberts, C. K. & Gómez-Pinilla, F. A
4862 high-fat, refined sugar diet reduces hippocampal brain-derived neurotrophic factor,

4863 neuronal plasticity, and learning. *Neuroscience* **112**, 803–814 (2002).

4864 95. Graber, T. G., Borack, M. S., Reidy, P. T., Volpi, E. & Rasmussen, B. B.
4865 Essential amino acid ingestion alters expression of genes associated with amino acid
4866 sensing, transport, and mTORC1 regulation in human skeletal muscle. *Nutr. Metab.* **14**,
4867 35 (2017).

4868 96. Hellsten, S. V., Hägglund, M. G., Eriksson, M. M. & Fredriksson, R. The
4869 neuronal and astrocytic protein SLC38A10 transports glutamine, glutamate, and
4870 aspartate, suggesting a role in neurotransmission. *FEBS Open Bio* **7**, 730–746 (2017).

4871 97. Tripathi, R., Hosseini, K., Arapi, V., Fredriksson, R. & Bagchi, S. SLC38A10
4872 (SNAT10) is Located in ER and Golgi Compartments and Has a Role in Regulating
4873 Nascent Protein Synthesis. *International Journal of Molecular Sciences* vol. 20 6265
4874 (2019).

4875 98. Schmidt, S. *et al.* Effect of omega-3 polyunsaturated fatty acids on the
4876 cytoskeleton: an open-label intervention study. *Lipids Health Dis.* **14**, 4 (2015).

4877 99. Zulkafli, I. S., Waddell, B. J. & Mark, P. J. Postnatal Dietary Omega-3 Fatty
4878 Acid Supplementation Rescues Glucocorticoid-Programmed Adiposity, Hypertension,
4879 and Hyperlipidemia in Male Rat Offspring Raised on a High-Fat Diet. *Endocrinology* vol.
4880 154 3110–3117 (2013).

4881 100. Aqil, M., Mallik, S., Bandyopadhyay, S., Maulik, U. & Jameel, S.
4882 Transcriptomic Analysis of mRNAs in Human Monocytic Cells Expressing the HIV-1 Nef
4883 Protein and Their Exosomes. *Biomed Res. Int.* **2015**, (2015).

4884 101. Kirpich, I. A. *et al.* Integrated hepatic transcriptome and proteome analysis of
4885 mice with high-fat diet-induced nonalcoholic fatty liver disease. *J. Nutr. Biochem.* **22**,
4886 38–45 (2011).

4887 102. Li, H. *et al.* Diversification of the ADH1B gene during expansion of modern
4888 humans. *Ann. Hum. Genet.* **75**, 497–507 (2011).

4889 103. Polimanti, R. & Gelernter, J. ADH1B: From alcoholism, natural selection, and
4890 cancer to the human phenome. *Am. J. Med. Genet. B Neuropsychiatr. Genet.* **177**,

- 4891 113–125 (2018).
- 4892 104. Zhong, X. *et al.* The zinc-finger protein ZFYVE1 modulates TLR3-mediated
4893 signaling by facilitating TLR3 ligand binding. *Cell. Mol. Immunol.* **17**, 741–752 (2020).
- 4894 105. Yang, Y. *et al.* The RNA-binding protein Mex3B is a coreceptor of Toll-like
4895 receptor 3 in innate antiviral response. *Cell Res.* **26**, 288–303 (2016).
- 4896 106. Kim, S. V. *et al.* GPR15-mediated homing controls immune homeostasis in
4897 the large intestine mucosa. *Science* **340**, 1456–1459 (2013).
- 4898 107. Nguyen, L. P. *et al.* Role and species-specific expression of colon T cell
4899 homing receptor GPR15 in colitis. *Nature Immunology* vol. 16 207–213 (2015).
- 4900 108. Monteleone, G., Boirivant, M., Pallone, F. & MacDonald, T. T. TGF- β 1 and
4901 Smad7 in the regulation of IBD. *Mucosal Immunology* vol. 1 S50–S53 (2008).
- 4902 109. Kennedy, B. W. C. Mongersen, an Oral SMAD7 Antisense Oligonucleotide,
4903 and Crohn’s Disease. *The New England journal of medicine* vol. 372 2461 (2015).
- 4904 110. Garo, L. P. *et al.* Smad7 Controls Immunoregulatory PDL2/1-PD1 Signaling in
4905 Intestinal Inflammation and Autoimmunity. *Cell Rep.* **28**, 3353–3366.e5 (2019).
- 4906 111. Iqbal, Z. *et al.* Homozygous SLC6A17 mutations cause autosomal-recessive
4907 intellectual disability with progressive tremor, speech impairment, and behavioral
4908 problems. *Am. J. Hum. Genet.* **96**, 386–396 (2015).
- 4909 112. Park, J. *et al.* KCNC1-related disorders: new de novo variants expand the
4910 phenotypic spectrum. *Ann Clin Transl Neurol* **6**, 1319–1326 (2019).

4911 **4e) Correlation between components of variation in population**
4912 **structure and components of variation in SNP-trait association**

4913
4914
4915
4916
4917
4918

Alba Refoyo Martínez¹, Fernando Racimo¹

¹Lundbeck Foundation GeoGenetics Centre, GLOBE Institute, University of Copenhagen,
Copenhagen, Denmark

4919 Introduction

4920 Ancient Western and Central Eurasian populations show strong patterns of genome-wide
4921 differentiation, which are even stronger than those observed among present-day Western
4922 and Central Eurasians (Supplementary Note S4d). We were interested in determining
4923 whether variants that contribute to differentiation between populations may also be
4924 associated with complex traits. This could perhaps serve to evince the genetic component of
4925 trait differences among these ancient populations.

4926 Recently, Tanigawa et al., 2019¹ developed a method (DeGAs) by which SNP-trait
4927 associations from thousands of phenotypes could be represented in a low-dimensional
4928 space, yielding a smaller set of latent components of genetic associations. They performed
4929 GWAS in 337,199 White British individuals in the UK Biobank study for 2,138 phenotypes,
4930 obtained Z-scores from each SNP and each trait, and then performed truncated singular-
4931 value decomposition (TSVD) on the matrix of Z-scores, which yielded 100 components of
4932 variation. The decomposition allows one to see how much a particular variant or trait
4933 contributes to the variation explained by that component, and to rank the components by
4934 their amount of contribution to the total variation in SNP-trait associations.

4935 By looking at these components of genetic association in the UK Biobank and comparing
4936 them with our components of variation in ancient population structure, we can begin to
4937 understand which sets of trait-associated variants may have been important during recent
4938 human evolution.

4939 Methods

4940 We carried out a principal component analysis (PCA) on the 1,165 West-Eurasian imputed
4941 genomes, using *pcadapt*². At the variant level we filtered out, 1) variants with MAF < 5%, 2)
4942 genotype missingness rate > 50% and, 3) variants with a genotype probability lower than 0.8
4943 in more than 10% of the samples (Figure 1).

4944 Our aim is to study the correlation between our population structure components, captured in
4945 the first two components, and the 100 components of SNP-trait associations from the
4946 DeGAs analysis, which were obtained from ¹. Henceforth, we will refer to the components of
4947 population structure from the ancient genome PCA with the label “PS” (i.e. the first principal
4948 component is PS1, the second is PS2, etc.). We will refer to the components of trait-
4949 association variation from Tanigawa et al., 2019¹ as “DG” (i.e. the first component is DG1,
4950 the second is DG2, etc.). We applied the Pearson correlation between the loadings of the
4951 two analyses (Figure 3).

4952 We aimed to test whether the correlations observed were significantly different from those
4953 that would be observed under a model in which there was no association between a
4954 particular DG component and a particular PS component. For this, we obtained p-values
4955 using a randomization scheme. We randomised the sign of the 1000 DGs loadings
4956 accounting for the structure of linkage disequilibrium (LD) along the genome, by dividing our
4957 genome into 1Mb blocks and randomising the sign of all SNPs within each block in unison
4958 recalculating the correlation with our first two components (PS1-PS2).

4959

4960 For each combination of PS_x and DG_y (where x and y are indices over all PS and DG
4961 components, respectively), we obtained a P-value using the following equation:

4962

$$4963 \quad P = \frac{1 + \sum_i^N I(|cor_x^i| > |cor_x^y|)}{1 + N}$$

4964 Here, cor_x^y is the true correlation between PS_x and DG_y, cor_x^y is the correlation between
4965 PS_x and DG_y after randomising signs, $I()$ is an indicator function and N is the number of
4966 randomised samples used, which was set to 1,000.

4967 Results

4968 We first performed a PCA on the genotypes of the 1,165 imputed ancient West-Eurasian
4969 individuals (Supplementary Note S4d). The first component (1.8% of total variance
4970 explained) represents a gradient separating Neolithic farmer populations from hunter-
4971 gatherer genomes, while the second component (1.1%) captures a gradient separating
4972 ancient East Asians, ancient Steppe populations and ancient Western Eurasians (Figure
4973 **S4d.4**).

4974 We computed the correlation between the first 8 components of population structure and the
4975 100 trait-association components. Figure **S4e.1** shows the distribution of all DG correlations
4976 with the first 8 PS components (PS1 to PS8). For this analysis, we focused on PS1 and PS2.
4977 PS1 captures a gradient separating ancient Neolithic farmers from Mesolithic hunter-
4978 gatherer genomes. The trait components most correlated with PS1 are DG89, DG71, DG22,
4979 DG15, DG52, DG5, DG82, DG2, DG51, DG99, DG69, DG12, DG84, DG83 and DG79
4980 (Table XX and top four 4 in Figure **S4e.2**). These components have in common that are
4981 driven by anthropometric and lifestyle and environmental traits, with an important
4982 contribution from diet-related measures. They are mainly correlated with five different
4983 broader categories: 1) lifestyle and environment, mainly representing food intake, cereal and
4984 fruit, and time spent outdoor and sun exposure; 2) verbal interview about mental health

4985 described by worrier/anxious/nervous feelings and sleep duration; 3) impedance, body
4986 measurement and bone mineral density; 4) lung capacity and asthma, mostly represented
4987 by the FEV1 measure; 5) blood pressure measures, hypertension and cholesterol.

4988 PS2 separates ancient East Asian and European samples. The significant top correlations of
4989 DG components with PS2, top four shown in Figure **S4e.3**, are DG38, DG12, DG40, DG84,
4990 DG56, DG82, DG36, DG15, DG28. and DG68. One of them was also significant for PS1.
4991 These components are mainly driven by anthropometric measures and cardiovascular
4992 measurement and disorders. They can be classified into the same categories previously
4993 described. The main difference with PS1 is that more components are mainly explained by
4994 the lifestyle and environmental relation and impedance measure while in PS2, blood
4995 pressure and cardiovascular disorders and skin-sun exposure relation have higher
4996 contributions.

4997 We computed p-values using two randomization schemes to test whether these correlations
4998 were significantly different from zero. The first test relies on the randomization of the sign of
4999 the loadings before calculating the correlations between the loadings representing
5000 components of structure and the loadings from the DeGAs analysis. We first divided the
5001 genome into 5Mb blocks and then randomised the signs of the SNPs in each block in the
5002 same direction, to account for LD. We also use a bootstrapping approach also dividing the
5003 genome into blocks.

5004 We observe there are 22 significant correlations with PS1 while there are 7 significant
5005 correlations with PS2. Significant p-values for the block-based randomization for PS1 and
5006 PS2 are shown in **Table S4e.1 and Table S4e.2** respectively.

5007 Discussion

5008 Among the significantly PS1-associated trait components, there are many components
5009 involving traits related to lifestyle and environment: eg. water intake, cereal intake and time
5010 spent outdoors. This signal may perhaps be related to differences in diet and lifestyle
5011 between ancient hunter-gatherers and farmers^{3,4}. We also find traits related to
5012 haematological measurements, forced expiratory volume (see also Supplementary Note
5013 S4d), body fat, anxiety and hypertension, among others.

5014

5015 However, we should also keep in mind that population stratification in the original UK
5016 Biobank dataset may affect these results. For example, if present-day British individuals with
5017 different amounts of Neolithic farmer ancestry happen to live in different areas of Britain with
5018 markedly different lifestyles, one could presume that a spurious signal of correlation between

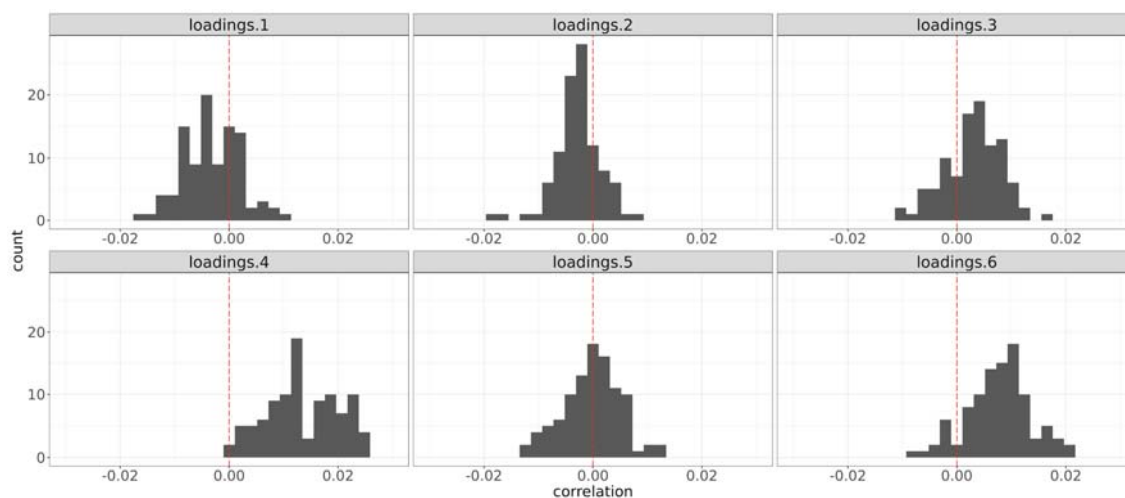
5019 those lifestyle differences and amount of Neolithic ancestry could be generated. Indeed,
5020 while non-significant, we find that some of the trait association components that are most
5021 correlated with PS1 have to do with lifestyle and environmental traits such as tea intake, and
5022 mental health, and these may, in turn, may be correlated with variation in PS1 in present-day
5023 British people (due, perhaps, to more recent immigrants from Asia in urban settings, for
5024 example).

5025

5026 The strongest DG correlations with the first two PCs performed on the European populations
5027 from the 1000 Genomes Project (FIN, CEU, GBR, TSI and IBS) are quite different to the
5028 ones with the ancient samples. Only DG52, which is mainly explained by asthma and bone
5029 mineral density traits, is shared with the West-Eurasian study. The first component, which
5030 separates Finnish from the other European populations, is strongly correlated with those
5031 DGs in which haematological measurements contribute between 40-92% (Table S4e.3). The
5032 second component, that separates the British and Utah residents (CEPH) with Northern and
5033 Western European ancestry from the Southern European and Finnish populations. This
5034 component is correlated with DG components in which spirometry traits that measure lung
5035 capacity, lifestyle and environment, such as water, coffee and tea intake, and body
5036 measurement and impedance traits, are the ones contributing most (Table S4e.4).

5037

5038 Figures

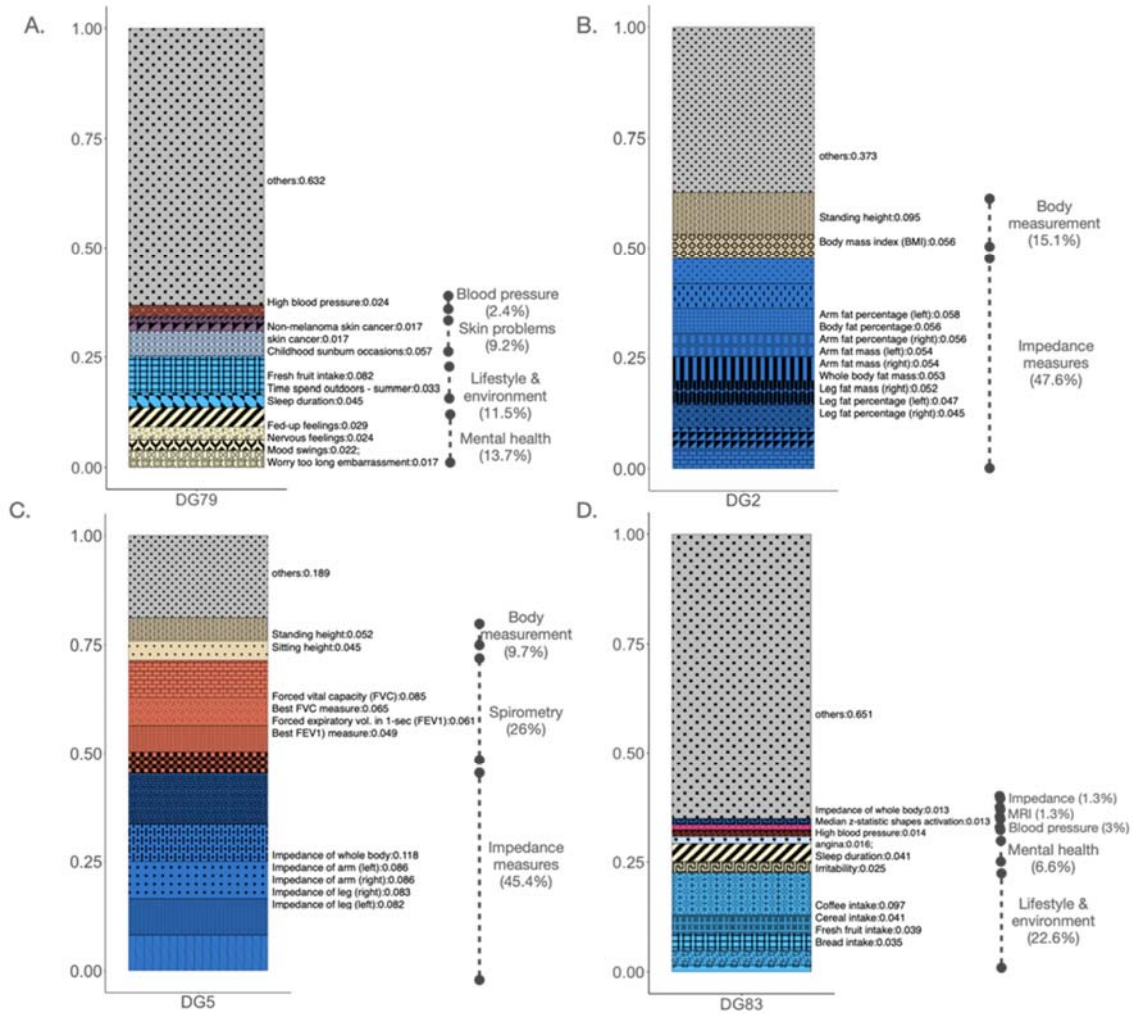


5039

5040

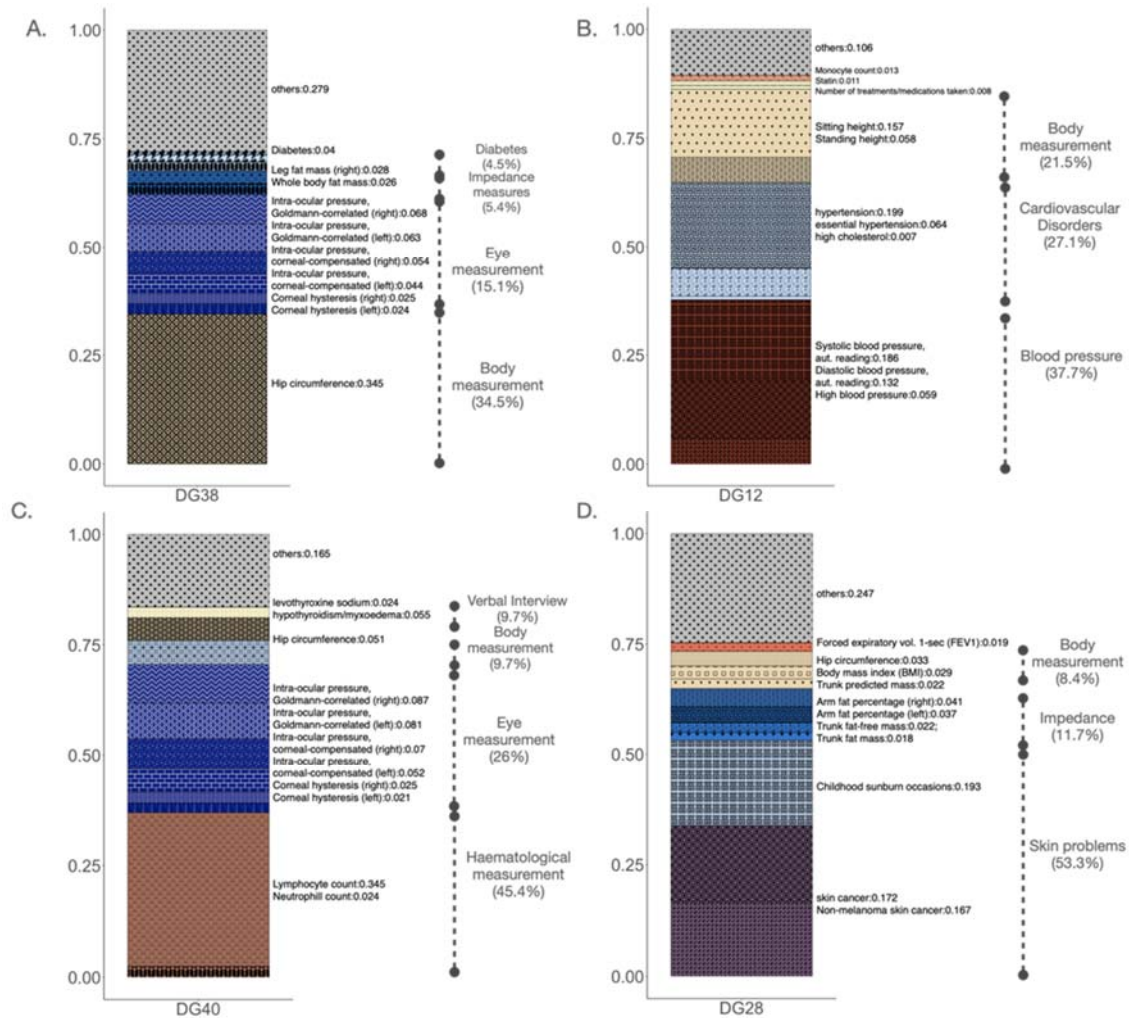
5041 **Figure S4e.1. Pearson correlation between the 6 loadings from the PCA analysis (PS1-**
5042 **PS6) and the 100 loadings from the DeGAs analysis (DG1-DG100).**

5043



5044

5045 **Figure S4e.2. PS1.** Top 4 UK Biobank trait-association components from DeGAs that have
 5046 the highest correlation ($P < 5e-3$) with the principal component separating ancient farmer
 5047 hunter-gatherer populations



5048

5049 **Figure S4e.3. PS2.** Top 4 UK Biobank trait-association components from DeGAs that have
 5050 the highest correlation ($P < 5e-3$) with the principal component separating East Asian and
 5051 European samples.

5052

5053

Neolithic farmer vs. hunter-gatherer genomes gradient				
D G	Correlation	P-value	contribution_category	contribution_phenotypes
79	-0.017	0.001	others:0.632; Blood pressure:0.024; Cancer:0.034; Non-cancer Illness:0.057; Lifestyle and environment:0.115; Verbal Interview:0.137	others:0.632; High blood pressure:0.024; Non-melanoma skin cancer:0.017; skin cancer:0.017; Childhood sunburn occasions:0.057; Fresh fruit intake:0.082; Time spend outdoors in summer:0.033; Sleep duration:0.045; Fed-up feelings:0.029; Nervous feelings:0.024; Mood swings:0.022; Worry too long after embarrassment:0.017
2	-0.014	0.001	others:0.373; Body measurement:0.151; Impedance measures:0.476	others:0.373; Standing height:0.095; Body mass index (BMI):0.056; Arm fat percentage (left):0.058; Body fat percentage:0.056; Arm fat percentage (right):0.056; Arm fat mass (left):0.054; Arm fat mass (right):0.054; Whole body fat mass:0.053; Leg fat mass (right):0.052; Leg fat percentage (left):0.047; Leg fat percentage (right):0.045
5	-0.013	0.003	others:0.189; Body measurement:0.097; Spirometry:0.26; Impedance measures:0.454	others:0.189; Standing height:0.052; Sitting height:0.045; Forced vital capacity (FVC):0.085; Forced vital capacity (FVC), Best measure:0.065; Forced expiratory volume in 1-second (FEV1):0.061; Forced expiratory volume in 1-second (FEV1), Best measure:0.049; Impedance of whole body:0.118; Impedance of arm (left):0.086; Impedance of arm (right):0.086; Impedance of leg (right):0.083; Impedance of leg (left):0.082
83	-0.013	0.003	others:0.651; Impedance measures:0.013; Brain MRI:0.013; Blood pressure:0.014; Non-cancer Illness:0.016; Verbal Interview:0.066; Lifestyle and environment:0.226	others:0.651; Impedance of whole body:0.013; Median z-statistic (in group-defined mask) for shapes activation:0.013; High blood pressure:0.014; angina:0.016; Sleep duration:0.041; Irritability:0.025; Coffee intake:0.097; Cereal intake:0.041; Fresh fruit intake:0.039; Bread intake:0.035; Water intake:0.015
71	-0.012	0.003	others:0.506; Non-cancer Illness:0.041; Hematological measurement:0.073; Verbal Interview:0.108; Lifestyle and environment:0.273	others:0.506; cholelithiasis/gall stones:0.024; cholecystitis:0.017; Monocyte percentage:0.036; Eosinophill percentage:0.02; Eosinophill count:0.016; Number of treatments/medications taken:0.108; Fresh fruit intake:0.104; Time spend outdoors in summer:0.074; Tea intake:0.045; Time spent outdoors in winter:0.033; Coffee intake:0.016
89	-0.012	0.005	others:0.58; Arterial stiffness:0.034; Spirometry:0.061; Lifestyle and environment:0.067; Non-cancer Illness:0.077; Verbal Interview:0.181	others:0.58; Pulse wave peak to peak time:0.017; Pulse wave Arterial Stiffness index:0.017; Forced expiratory volume in 1-second (FEV1), predicted percentage:0.061; Average weekly beer plus cider intake:0.067; hypertension:0.042; essential hypertension:0.035; Worrier/anxious feelings:0.053; Sleep duration:0.053; Fed-up feelings:0.033; Age at menopause (last menstrual period):0.022; Number of treatments/medications taken:0.021

22	0.014	0.00999	others:0.225; Hematological measurement:0.018; Impedance measures:0.029; Blood pressure:0.083; Non-cancer illness:0.317; Verbal Interview:0.329	others:0.225; Neutrophil count:0.018; Impedance of arm (left):0.016; Leg fat percentage (right):0.013; Systolic blood pressure, automated reading:0.03; Diastolic blood pressure, automated reading:0.026; Pulse rate, automated reading:0.026; high cholesterol:0.303 ; Alzheimer's disease/dementia:0.014; Statin:0.237 ; simvastatin:0.062; atorvastatin:0.03
99	-0.01	0.01698	others:0.63; Spirometry:0.032; Lifestyle and environment:0.061; Impedance measures:0.067; Verbal Interview:0.087; Eye measurement:0.122	others:0.63; Forced expiratory volume in 1-second (FEV1), predicted percentage:0.032; Average weekly beer plus cider intake:0.061; Impedance of arm (left):0.041; Impedance of arm (right):0.026; Age started wearing glasses or contact lenses:0.071; Seen doctor (GP) for nerves, anxiety, tension or depression:0.016; logMAR in round (left):0.039; logMAR, initial (left):0.027; logMAR, final (left):0.024; Corneal resistance factor (left):0.017; Corneal hysteresis (right):0.015
80	-0.01	0.01698	others:0.652; Hematological measurement:0.017; Non-cancer illness:0.017; Sex-specific factors:0.021; Lifestyle and environment:0.046; Verbal Interview:0.053; Impedance measures:0.061; Body measurement:0.131	others:0.652; Monocyte percentage:0.017; hypertension:0.017; Birth weight of first child:0.021; Water intake:0.027; Coffee intake:0.019; Worrier/anxious feelings:0.02; gout:0.018; Sleep duration:0.015; Arm predicted mass (left):0.032; Arm fat-free mass (left):0.029; Weight:0.131
69	-0.01	0.02198	others:0.529; Blood pressure:0.038; Impedance measures:0.068; Body measurement:0.089; Hematological measurement:0.132; Lifestyle and environment:0.144	others:0.529; High blood pressure:0.038; Impedance of whole body:0.045; Impedance of arm (right):0.023; Body mass index (BMI):0.062; Weight:0.027; Monocyte percentage:0.11; Eosinophil percentage:0.022; Cereal intake:0.047; Water intake:0.034; Time spend outdoors in summer:0.033; Time spent outdoors in winter:0.029
77	-0.009	0.01798	others:0.748; Urine assays:0.042; Lifestyle and environment:0.043; Verbal Interview:0.049; Brain MRI:0.118	others:0.748; Sodium in urine:0.042; Fresh fruit intake:0.043; Sleep duration:0.028; Number of treatments/medications taken:0.021; Weighted-mean FA in tract superior longitudinal fasciculus (right):0.022; Weighted-mean FA in tract superior longitudinal fasciculus (left):0.019; Mean FA in superior longitudinal fasciculus on FA skeleton (right):0.019; Weighted-mean MD in tract anterior thalamic radiation (right):0.018; Weighted-mean MD in tract anterior thalamic radiation (left):0.015; Weighted-mean L1 in tract anterior thalamic radiation (left):0.013; Mean FA in superior longitudinal fasciculus on FA skeleton (left):0.012
22	0.009	0.01898	others:0.225; Hematological measurement:0.018; Impedance measures:0.029; Blood pressure:0.083; Non-cancer illness:0.317; Verbal Interview:0.329	others:0.225; Neutrophil count:0.018; Impedance of arm (left):0.016; Leg fat percentage (right):0.013; Systolic blood pressure, automated reading:0.03; Diastolic blood pressure, automated reading:0.026; Pulse rate, automated reading:0.026; high cholesterol:0.303; Alzheimer's disease/dementia:0.014; Statin:0.237 ; simvastatin:0.062; atorvastatin:0.03
90	-0.009	0.03197	others:0.631; Impedance measures:0.021; Eye	others:0.631; Impedance of arm (right):0.021; logMAR in round (left):0.027; Creatinine (enzymatic) in urine:0.027; Nucleated red blood cell percentage:0.03;

			measurement:0.027; Urine assays:0.027; Blood pressure:0.03; Spirometry:0.034; Non-cancer illness:0.057; Verbal Interview:0.075; Lifestyle and environment:0.098	Forced expiratory volume in 1-second (FEV1), predicted percentage:0.034; Childhood sunburn occasions:0.034; essential hypertension:0.022; Number of treatments/medications taken:0.05; Irritability:0.025; Fresh fruit intake:0.05; Bread intake:0.048
7 5	-0.009	0.0279 7	others:0.509; Verbal Interview:0.071; Non-cancer illness:0.084; Hematological measurement:0.107; Impedance measures:0.108; Lifestyle and environment:0.12	others:0.509; Number of treatments/medications taken:0.071; Alzheimer's disease/dementia:0.061; essential hypertension:0.023; Monocyte percentage:0.088; Monocyte count:0.019; Whole body water mass:0.041; Whole body fat-free mass:0.04; Basal metabolic rate:0.027; Fresh fruit intake:0.048; Time spend outdoors in summer:0.043; Time spent outdoors in winter:0.029
5 1	0.011	0.0389 6	others:0.425; Verbal Interview:0.024; DXA assessment:0.26; Non-cancer illness:0.291	others:0.425; Wheeze or whistling in the chest in last year:0.024; Femur total BMD (bone mineral density) (left):0.038; Femur total BMD (bone mineral density) T-score (left):0.035; Femur total BMD (bone mineral density) (right):0.033; Femur total BMD (bone mineral density) T-score (right):0.031; Femur troch BMD (bone mineral density) (left):0.028; Femur troch BMD (bone mineral density) (right):0.026; Femur shaft BMD (bone mineral density) (left):0.024; Femur troch BMD (bone mineral density) T-score (left):0.023; Femur troch BMD (bone mineral density) T-score (right):0.021; asthma:0.291
6 1	-0.009	0.0379 6	others:0.612; Urine assays:0.016; Lifestyle and environment:0.049; Impedance measures:0.055; Body measurement:0.127; Spirometry:0.14	others:0.612; Sodium in urine:0.016; Water intake:0.036; Time spend outdoors in summer:0.013; Leg fat mass (right):0.027; Arm fat percentage (right):0.015; Trunk fat percentage:0.013; Body mass index (BMI):0.063; Weight:0.035; Waist circumference:0.028; Peak expiratory flow (PEF):0.11; Forced expiratory volume in 1-second (FEV1):0.03
8 2	-0.012	0.0409 6	others:0.641; Urine assays:0.027; Body measurement:0.073; Lifestyle and environment:0.078; Impedance measures:0.09; Verbal Interview:0.092	others:0.641; Creatinine (enzymatic) in urine:0.027; Hand grip strength (right):0.038; Hand grip strength (left):0.035; Tea intake:0.063; Water intake:0.015; Impedance of leg (left):0.035; Impedance of leg (right):0.022; Trunk fat mass:0.017; Trunk fat mass:0.017; Sleep duration:0.066; Nervous feelings:0.026
5 2	0.013	0.0229 8	others:0.538; Verbal Interview:0.019; Impedance measures:0.027; DXA assessment:0.171; Non-cancer illness:0.245	others:0.538; Wheeze or whistling in the chest in last year:0.019; Impedance of whole body:0.027; Femur total BMD (bone mineral density) (left):0.027; Femur total BMD (bone mineral density) T-score (left):0.025; Femur total BMD (bone mineral density) (right):0.024; Femur total BMD (bone mineral density) T-score (right):0.022; Femur troch BMD (bone mineral density) (left):0.02; Femur troch BMD (bone mineral density) (right):0.019; Femur shaft BMD (bone mineral density) (left):0.017; Femur troch BMD (bone mineral density) T-score (left):0.016; asthma:0.245
6 8	-0.01	0.0479 5	others:0.697; Verbal Interview:0.015; Hematological measurement:0.022; Impedance measures:0.036; Lifestyle and	others:0.697; Long-standing illness, disability or infirmity:0.015; Monocyte percentage:0.022; Leg fat-free mass (left):0.018; Leg predicted mass (left):0.018; Tea intake:0.037; High blood pressure:0.061; Alzheimer's

			environment:0.037; Blood pressure:0.061; Non-cancer illness:0.132	disease/dementia:0.038; cholelithiasis/gall stones:0.03; hypertension:0.03; cholecystitis:0.019; Childhood sunburn occasions:0.015
84	-0.015	0.004	others:0.648; Urine assays:0.051; Blood pressure:0.07; Verbal Interview:0.077; Lifestyle and environment:0.154	others:0.648; Creatinine (enzymatic) in urine:0.035; Potassium in urine:0.016; High blood pressure:0.07; Sleep duration:0.022; Fed-up feelings:0.02; Worrier/anxious feelings:0.019; Nervous feelings:0.017; Tea intake:0.072; Time spend outdoors in summer:0.042; Coffee intake:0.022; Time spent outdoors in winter:0.018
46	-0.011	0.01499	others:0.096; Body measurement:0.045; Non-cancer illness:0.055; Blood pressure:0.181; Spirometry:0.219; Verbal Interview:0.406	others:0.096; Waist circumference:0.045; diabetes:0.02; Diabetes diagnosed by doctor:0.018; Diabetes:0.017; Diastolic blood pressure, automated reading:0.111; Systolic blood pressure, automated reading:0.07; Forced vital capacity (FVC):0.069; Forced vital capacity (FVC), Best measure:0.066; Forced expiratory volume in 1-second (FEV1), Best measure:0.061; Forced expiratory volume in 1-second (FEV1):0.023; Age when periods started (menarche):0.406
70	-0.008	0.03896	others:0.57; Body measurement:0.033; Urine assays:0.047; Non-cancer illness:0.053; Impedance measures:0.135; Lifestyle and environment:0.162	others:0.57; Weight:0.033; Sodium in urine:0.047; Childhood sunburn occasions:0.026; Alzheimer's disease/dementia:0.026; Impedance of leg (right):0.037; Impedance of leg (left):0.028; Arm predicted mass (left):0.026; Arm fat-free mass (left):0.023; Arm predicted mass (right):0.021; Cereal intake:0.111; Water intake:0.051

5054

5055

5056

Table S4e.1. Significant correlations between DG components with PS1 1) randomising the sign of each block of 5Mb in the DG loadings and bootstrapping the same blocks.

Ancient East Asians, ancient Steppe populations and ancient Western Eurasians gradient				
D G	Correlation	P-value	contribution_category	contribution_phenotypes
38	-0.024	0.001	others:0.279; Non-cancer illness:0.045; Impedance measures:0.054; Eye measurement:0.277; Body measurement:0.345	others:0.279; diabetes:0.024; Diabetes:0.02; Leg fat mass (right):0.028; Whole body fat mass:0.026; Intra-ocular pressure, Goldmann-correlated (right):0.068; Intra-ocular pressure, Goldmann-correlated (left):0.063; Intra-ocular pressure, corneal-compensated (right):0.054; Intra-ocular pressure, corneal-compensated (left):0.044; Corneal hysteresis (right):0.025; Corneal hysteresis (left):0.024; Hip circumference:0.345
12	-0.024	0.001	others:0.106; Hematological measurement:0.013; Verbal Interview:0.019; Body measurement:0.215; Non-cancer illness:0.271; Blood pressure:0.377	others:0.106; Monocyte count:0.013; Statin:0.011; Number of treatments/medications taken:0.008; Sitting height :0.157; Standing height:0.058; hypertension :0.199; essential hypertension:0.064; high cholesterol:0.007; Systolic blood pressure , automated reading:0.186; Diastolic blood pressure, automated reading:0.132; High blood pressure:0.059
40	-0.015	0.00799	others:0.165; Verbal Interview:0.024; Body measurement:0.051; Non-cancer illness:0.055; Eye measurement:0.336; Hematological measurement:0.369	others:0.165; levothyroxine sodium:0.024; Hip circumference:0.051; hypothyroidism/myxoedema:0.055; Intra-ocular pressure, Goldmann-correlated (right):0.087; Intra-ocular pressure, Goldmann-correlated (left):0.081; Intra-ocular pressure, corneal-compensated (right):0.07; Intra-ocular pressure, corneal-compensated (left):0.052; Corneal hysteresis (right):0.025; Corneal hysteresis (left):0.021; Lymphocyte count:0.345; Neutrophill count:0.024
28	-0.01	0.001	others:0.247; Spirometry:0.019; Body measurement:0.084; Impedance measures:0.117; Non-cancer illness:0.193; skin Cancer :0.339	others:0.247; Forced expiratory volume in 1-second (FEV1), predicted:0.019; Hip circumference:0.033; Body mass index (BMI):0.029; Trunk predicted mass:0.022; Arm fat percentage (right):0.041; Arm fat percentage (left):0.037; Trunk fat-free mass:0.022; Trunk fat mass:0.018; Childhood sunburn occasions:0.193; skin cancer:0.172; Non-melanoma skin cancer:0.167
23	-0.009	0.02897	others:0.061; Hematological measurement:0.011; Non-cancer illness:0.014; Cancer:0.02; Local environment:0.894	others:0.061; Immature reticulocyte fraction:0.011; Childhood sunburn occasions:0.014; skin cancer:0.01; Non-melanoma skin cancer:0.01; Nitrogen dioxide air pollution; 2010:0.176; Nitrogen dioxide air pollution; 2007:0.171; Nitrogen dioxide air pollution; 2006:0.167; Nitrogen dioxide air pollution; 2005:0.152; Nitrogen oxides air pollution; 2010:0.113; Particulate matter air pollution (pm2.5); 2010:0.091; Particulate matter air pollution (pm2.5) absorbance; 2010:0.024
24	-0.009	0.04296	others:0.056; Body measurement:0.022; Hematological measurement:0.922	others:0.056; Hand grip strength (left):0.012; Hand grip strength (right):0.011; Immature reticulocyte fraction:0.325; Reticulocyte count:0.234; High light scatter reticulocyte percentage:0.142; Reticulocyte percentage:0.121; High light scatter reticulocyte count:0.051; Neutrophill count:0.028; White blood cell (leukocyte) count:0.01; Red blood cell (erythrocyte) count:0.007; Lymphocyte count:0.004

84	-0.015	0.004	others:0.648; Urine assays:0.051; Blood pressure:0.07; Verbal Interview:0.077; Lifestyle and environment:0.154	others:0.648; Creatinine (enzymatic) in urine:0.035; Potassium in urine:0.016; High blood pressure:0.07; Sleep duration:0.022; Fed-up feelings:0.02; Worrier/anxious feelings:0.019; Nervous feelings:0.017; Tea intake:0.072; Time spend outdoors in summer:0.042; Coffee intake:0.022; Time spent outdoors in winter:0.018
----	--------	-------	--	---

5058

5059 **Table S4e.2.** Significant correlations between DG components with PS2 1) randomising the sign of each block of 5Mb in the DG loadings and
5060 bootstrapping the same blocks

5061
5062
5063

1000 GP – FIN vs other European populations gradient.				
DG	Correlation	P-value	contribution_category	contribution_phenotypes
46	-0.011	0.01499	others:0.096; Body measurement:0.045; Non-cancer illness:0.055; Blood pressure:0.181; Spirometry:0.219; Verbal Interview:0.406	others:0.096; Waist circumference:0.045; diabetes:0.02; Diabetes diagnosed by doctor:0.018; Diabetes:0.017; Diastolic blood pressure, automated reading:0.111; Systolic blood pressure, automated reading:0.07; Forced vital capacity (FVC):0.069; Forced vital capacity (FVC), Best measure:0.066; Forced expiratory volume in 1-second (FEV1), Best measure:0.061; Forced expiratory volume in 1-second (FEV1):0.023; Age when periods started (menarche):0.406
27	0.009	0.02498	others:0.093; Hematological measurement:0.907	others:0.093; Mean platelet (thrombocyte) volume:0.408; Neutrophil count:0.11; Red blood cell (erythrocyte) count:0.096; Platelet distribution width:0.052; Haemoglobin concentration:0.05; Red blood cell (erythrocyte) distribution width:0.041; Haematocrit percentage:0.037; Mean corpuscular haemoglobin concentration:0.036; Mean corpuscular haemoglobin:0.034; Monocyte percentage:0.025; Lymphocyte percentage:0.018
20	0.009	0.03996	others:0.189; Verbal Interview:0.016; Non-cancer illness:0.022; Impedance measures:0.051; Hematological measurement:0.723	others:0.189; Statin:0.016; high cholesterol:0.022; Leg fat percentage (right):0.027; Leg fat percentage (left):0.024; Neutrophil count:0.353; White blood cell (leukocyte) count:0.13; Red blood cell (erythrocyte) count:0.068; Monocyte count:0.064; Monocyte percentage:0.038; Mean platelet (thrombocyte) volume:0.037; Immature reticulocyte fraction:0.033
24	-0.009	0.04296	others:0.056; Body measurement:0.022; Hematological measurement:0.922	others:0.056; Hand grip strength (left):0.012; Hand grip strength (right):0.011; Immature reticulocyte fraction:0.325; Reticulocyte count:0.234; High light scatter reticulocyte percentage:0.142; Reticulocyte percentage:0.121; High light scatter reticulocyte count:0.051; Neutrophil count:0.028; White blood cell (leukocyte) count:0.01; Red blood cell (erythrocyte) count:0.007; Lymphocyte count:0.004

5064
5065
5066
5067

Table S4e.3. Significant correlations between DG components with PS1 of European population in 1000 Genomes Project 1) randomising the sign of each block of 5Mb in the DG loadings and bootstrapping the same blocks.

1000 GP – FIN and Southern Europeans vs Northern populations gradient.

D G	Correl ation	P- valu e	contribution_category	contribution_phenotypes
4 6	-0.014	0.00 2	others:0.096; Body measurement:0.045; Non-cancer Illness:0.055; Blood pressure:0.181; Spirometry:0.219; Verbal Interview:0.406	others:0.096; Waist circumference:0.045; diabetes:0.02; Diabetes diagnosed by doctor:0.018; Diabetes:0.017; Diastolic blood pressure, automated reading:0.111; Systolic blood pressure, automated reading:0.07; Forced vital capacity (FVC):0.069; Forced vital capacity (FVC), Best measure:0.066; Forced expiratory volume in 1-second (FEV1), Best measure:0.061; Forced expiratory volume in 1-second (FEV1):0.023; Age when periods started (menarche):0.406
4 5	-0.014	0.00 3	others:0.099; Non-cancer Illness:0.008; Hematological measurement:0.014; Body measurement:0.082; Blood pressure:0.132; Verbal Interview:0.177; Spirometry:0.489	others:0.099; asthma:0.008; Red blood cell (erythrocyte) count:0.014; Waist circumference:0.082; Diastolic blood pressure, automated reading:0.071; Systolic blood pressure, automated reading:0.061; Age when periods started (menarche):0.177; Forced vital capacity (FVC):0.141; Forced vital capacity (FVC), Best measure:0.133; Forced expiratory volume in 1-second (FEV1), Best measure:0.132; Forced expiratory volume in 1-second (FEV1):0.052; Peak expiratory flow (PEF):0.031
9 6	-0.011	0.01 798	others:0.661; Blood pressure:0.019; Eye measurement:0.02; Spirometry:0.025; Sex-specific factors:0.027; Impedance measures:0.034; Non-cancer Illness:0.043; Verbal Interview:0.047; Body measurement:0.123	others:0.661; Pulse rate:0.019; 6mm weak meridian (left):0.02; Forced expiratory volume in 1-second (FEV1), predicted percentage:0.025; Birth weight of first child:0.027; Impedance of arm (left):0.034; aortic dissection:0.022; aortic aneurysm:0.021; Number of treatments/medications taken:0.028; Age at menopause (last menstrual period):0.019; Hand grip strength (right):0.063; Hand grip strength (left):0.06
3 0	-0.01	0.01 099	others:0.341; Eye measurement:0.033; Verbal Interview:0.037; Non-cancer Illness:0.076; Body measurement:0.101; Cancer:0.153; Impedance measures:0.258	others:0.341; Spherical power (left):0.033; Neuroticism score:0.037; Childhood sunburn occasions:0.076; Body mass index (BMI):0.069; Trunk predicted mass:0.031; skin cancer:0.077; Non-melanoma skin cancer:0.076; Arm fat percentage (right):0.081; Arm fat percentage (left):0.07; Trunk fat mass:0.053; Trunk fat mass:0.053
5 2	-0.009	0.03 796	others:0.538; Verbal Interview:0.019; Impedance measures:0.027; DXA assessment:0.171 ; Non-cancer Illness:0.245	others:0.538; Wheeze or whistling in the chest in last year:0.019; Impedance of whole body:0.027; Femur total BMD (bone mineral density) (left):0.027; Femur total BMD (bone mineral density) T-score (left):0.025; Femur total BMD (bone mineral density) (right):0.024; Femur total BMD (bone mineral density) T-score (right):0.022; Femur troch BMD (bone mineral density) (left):0.02; Femur troch BMD (bone mineral density) (right):0.019; Femur shaft BMD (bone mineral density) (left):0.017; Femur troch BMD (bone mineral density) T-score (left):0.016; asthma:0.245
5 9	-0.009	0.03 796	others:0.455; Urine assays:0.011; Non-cancer Illness:0.023; Brain MRI:0.046; Lifestyle and environment:0.218 ; Spirometry:0.247	others:0.455; Sodium in urine:0.011; asthma:0.023; Volume of brain, grey+white:0.016; Volume of grey matter:0.016; Volume of peripheral cortical grey matter:0.014; Water intake:0.131; Tea intake:0.046; Coffee intake:0.041; Peak

				expiratory flow (PEF):0.175; Forced expiratory volume in 1-second (FEV1):0.049; Forced vital capacity (FVC), Best measure:0.022
100	-0.008	0.04795	others:0.746; Impedance measures:0.014; Hematological measurement:0.019; Urine assays:0.024; Sex-specific factors:0.027; Eye measurement:0.036; Verbal Interview:0.039; Lifestyle and environment:0.095	others:0.746; Impedance of leg (right):0.014; Monocyte percentage:0.019; Sodium in urine:0.024; Birth weight of first child:0.027; Corneal hysteresis (left):0.019; Corneal resistance factor (right):0.017; Age started wearing glasses or contact lenses:0.039; Average weekly beer plus cider intake:0.047; Fresh fruit intake:0.02; Tea intake:0.015; Coffee intake:0.014

5068

5069 **Table S4e.4.** Significant correlations between DG components with PS2 of European population in 1000 Genomes Project 1) randomising the
5070 sign of each block of 5Mb in the DG loadings and bootstrapping the same blocks.

5071 References

- 5072 1. Tanigawa, Y. *et al.* Components of genetic associations across 2,138 phenotypes in
5073 the UK Biobank highlight adipocyte biology. *Nat. Commun.* **10**, 4064 (2019).
- 5074 2. Luu, K., Bazin, E. & Blum, M. G. B. pcadapt: an R package to perform genome scans
5075 for selection based on principal component analysis. *Mol. Ecol. Resour.* **17**, 67–77
5076 (2017).
- 5077 3. Diamond, J. & Bellwood, P. Farmers and their languages: the first expansions.
5078 *Science* **300**, 597–603 (2003).
- 5079 4. Bowles, S. Cultivation of cereals by the first farmers was not more productive than
5080 foraging. *Proceedings of the National Academy of Sciences* vol. 108 4760–4765 (2011).

5081
5082

5083 4f) Polygenic prediction for height, eye colour and hair colour in 5084 ancient Danish samples

5085 Anders Rosengren^{1,2}, Vivek Appadurai^{1,2}, Andrew Schork^{1,4}, Andrés Ingason¹⁻³

5086

5087 ¹Institute of Biological Psychiatry, Mental Health Services, Copenhagen University Hospital

5088 ²iPSYCH, The Lundbeck Foundation Initiative for Integrative Psychiatric Research, Aarhus

5089 ³Lundbeck Foundation GeoGenetics Centre, GLOBE Institute, University of Copenhagen,

5090 Copenhagen, Denmark

5091 ⁴Neurogenomics Division, The Translational Genomics Research Institute (TGEN), Phoenix

5092

5093 We predicted relative “genetic” height (i.e. expected increase or decrease in height
5094 compared to the mean of the contemporary Danish population, based on common genetic
5095 variants) as well as eye and hair colour in 100 ancient samples excavated from
5096 archeological sites in Denmark. The estimated age of the ancient samples, all sequenced in
5097 this study, ranged from roughly 10,500 to 3,000 years (See Figure 4 in main text).

5098

5099 The relative genetic height was calculated from summary statistics of a recent GWAS on
5100 adult height in the UK Biobank ¹, using only strand-insensitive autosomal SNPs with robustly
5101 genome-wide significant allelic effects ($P < 1e-15$) with imputation info > 0.8 and minor allele
5102 frequency > 0.05 in our ancient imputation genotype dataset. We excluded long-range
5103 linkage disequilibrium (LD) regions ² and used the “--clump” function in PLINK ³ to prune for

5104 LD ($r^2 < 0.1$ within 10Mb window), rendering a total of 310 effect alleles. Per-sample genetic
5105 height score was then calculated for the 100 ancient samples as well as a subset of 3,467
5106 Danish ancestry male conscripts from a random population subset of the IPSYCH2012 case-
5107 cohort ⁴ by summing allelic effect multiplied with the effect allele imputed dosage (Appadurai
5108 *et al.*, in prep.) across the 310 loci. The genetic height score was moderately correlated with
5109 height in the subset of 3,467 Danish ancestry conscripts ($r^2 = 0.095$, $P = 3.2e-77$), and we
5110 rescaled the score to a unit corresponding to 1 cm change in predicted height with the
5111 median score in the conscript subset as zero. Thus, the PGS in each ancient sample
5112 corresponds to the predicted difference in cm from the average of the present-day Danish
5113 population, assuming that the scores are equally predictive in males and females. The
5114 genetic height score is however limited in two important ways; firstly, the predictive value of
5115 the genetic height score is modest and diminishes with general genetic distance to the
5116 population in which the allelic effects were determined (in this case European ancestry
5117 British); secondly, the score does not take account of important environmental factors such
5118 as health and access to nutrition in childhood.

5119

5120 The genetic predictions of eye and hair colour were done based on the HirisPlex system ⁵.
5121 Out of 24 main effect HirisPlex variants, genotype likelihoods of 18 SNPs were available in
5122 the ancient sample 1000G imputation, and imputed effect allele dosages of these were used
5123 to derive probabilities for brown, blue and grey/intermediate eye colour and blond, brown,
5124 black and red hair colour, following the HirisPlex formulas ⁵.

5125

5126 The results of genetic prediction of height, eye and hair colour are shown alongside results
5127 of other analyses in the composite figure 4 in the main text, titled “Denmark through time”.
5128 When excluding 27 samples with average genomic sequence coverage $< 0.1x$ (indicated in
5129 shaded colour in figure 4 and with “0” in the Cov QC column of the Table S4f.1 below), the
5130 predicted genetic height differed significantly across the three groups defined by the two
5131 major population turnover events from Mesolithic Hunter-Gatherer (ML-HG) to Neolithic
5132 Early Farmer (NL-EF), and later to Neolithic Steppe Pastoralist (NL-SP), and indicated with
5133 thick lines through the panels of figure 4 (ANOVA chi-square test across all three groups,
5134 $P = 5.1 \times 10^{-8}$). A follow-up pairwise testing found significant differences between all three
5135 group pairs (linear regression, $P < 0.05$ in all instances), with the lowest mean (plus/minus
5136 standard error) predicted relative genetic height observed in ML-HG (-1.9 ± 0.3 cm), then
5137 NL-EF (-0.3 ± 0.5 cm), and highest in NL-SP (1.2 ± 0.4 cm). It should be borne in mind that
5138 the population structure of the ancient Danish samples, especially the ML-HG and NL-EF
5139 groups, is different from the current day European ancestry British population (in which
5140 allelic effects for genetic height were estimated), and although we have used only very

5141 robustly associated effect alleles to calculate the genetic height score, it is likely that it will
 5142 not have correlated as well with actual height as it does in the current-day European
 5143 ancestry Danish population (in which the score was rescaled). Therefore, the only
 5144 conclusion that can be drawn from these results is that the common SNP alleles that
 5145 contribute most strongly to increased height in current day European ancestry populations,
 5146 were on average of slightly lower frequency in ML-HG, in similar frequency in NL-EF, and
 5147 slightly higher frequency in NL-SP.

5148

5149 Among the 18 HirisPlex SNPs used to predict eye and hair colour, rs12913832 has the
 5150 strongest overall dark/light pigmentation effect. At the same time rs12913832 has the lowest
 5151 average maximum genotype probability (GPmax) of the HirisPlex SNPs across the ancient
 5152 Danish imputed genotype dataset. To account for this, we applied a second quality filter
 5153 when comparing predicted eye and hair colour probabilities across groups, by requiring a
 5154 GPmax>0.6 for rs12913832 and for at least 15 of the other 17 pigmentation SNPs, which
 5155 removed a further 17 samples in the cross-group comparison (marked with “0” in the Pigm
 5156 QC column in Table S4f.1 below). In this comparison we did not find a significant difference
 5157 in probability of brown eye colour (pEye Brown in Table S4f.1 below) across the three
 5158 groups (ANOVA, P=0.21). On the other hand, the predicted probability of blond hair colour
 5159 differed significantly across groups (ANOVA, P=1.1x10⁻⁹), with the mean likelihood (pHair
 5160 Blond in Table S4f.1 below) increasing over time from ML-HG (0.05 ± 0.01) to NL-EF (0.25 ±
 5161 0.06) and NL-SP (0.43 ± 0.07), although the difference between NL-EF and NL-SP was not
 5162 significant (P=0.07). Although pigmentation traits are polygenic, many of the HirisPlex
 5163 system alleles are so-called main effect alleles and therefore it is likely that the increased
 5164 predicted probabilities for blond hair over time (and corresponding decrease in predicted
 5165 probabilities for black hair) represent a true change in the prevalence of dark and light hair
 5166 colour.

5167

5168 *Supplementary Table S4f.1. Estimated age, genomic coverage, predicted genetic height*
 5169 *difference from current Danish population average, and predicted likelihood of eye and hair*
 5170 *colour of 100 ancient Danish samples sequenced in this study*
 5171

	Age	Age	Genomic	Cov	Pigm	pHeight	pEye	pEye	pEye	pHair	pHair	pHair	pHair
Sample	(ybp)	group	coverage	QC	QC	(cm)	Blue	Inter	Brown	Blond	Red	Brown	Black
NEO254	10463	ML-HG	0.42	1	0	-1.2	0.09	0.13	0.78	0.04	0	0.19	0.77
NEO13	9507	ML-HG	0.01	0	0	-0.7	0.43	0.19	0.38	0.18	0	0.36	0.46
NEO91	9122	ML-HG	1.18	1	1	-1.8	0.44	0.22	0.34	0.03	0	0.26	0.71
NEO759	9028	ML-HG	2.95	1	1	0.2	0.58	0.16	0.26	0.04	0	0.21	0.75

NEO587 8798	ML-HG	1.14	1	1	-3.5	0.73	0.13	0.14	0.01	0	0.18	0.81
NEO123 8182	ML-HG	0.29	1	1	-3	0.09	0.13	0.78	0.11	0	0.26	0.63
NEO19 8163	ML-HG	3.26	1	1	-3.5	0.53	0.27	0.2	0	0	0.28	0.72
NEO122 8146	ML-HG	0.56	1	0	-0.4	0.12	0.22	0.66	0	0	0.15	0.85
NEO600 7817	ML-HG	0.10	0	1	-0.7	0.01	0.09	0.9	0.02	0	0.28	0.7
NEO683 7529	ML-HG	1.82	1	1	-1.8	0.26	0.25	0.49	0.01	0	0.22	0.77
NEO932 7499	ML-HG	2.76	1	1	-4.6	0.64	0.19	0.17	0.01	0	0.15	0.84
NEO589 7478	ML-HG	7.41	1	1	-1	0.46	0.21	0.33	0.02	0	0.24	0.74
NEO748 7129	ML-HG	0.08	0	0	-0.2	0.32	0.17	0.51	0.05	0	0.27	0.68
NEO814 7125	ML-HG	0.06	0	0	-2.6	0.08	0.16	0.76	0.03	0	0.24	0.73
NEO749 7070	ML-HG	1.91	1	1	-1.3	0.57	0.19	0.24	0.03	0	0.26	0.71
NEO791 7048	ML-HG	2.49	1	1	-0.1	0.84	0.11	0.05	0.05	0	0.31	0.64
NEO586 7031	ML-HG	0.20	1	1	-1.4	0.51	0.21	0.28	0.02	0	0.28	0.7
NEO746 6991	ML-HG	0.14	1	0	-3.1	0.16	0.15	0.69	0.11	0	0.29	0.6
NEO583 6981	ML-HG	0.18	1	1	-2.5	0.44	0.15	0.41	0.2	0	0.28	0.52
NEO822 6978	ML-HG	0.06	0	0	0.7	0.47	0.18	0.35	0.13	0	0.32	0.55
NEO930 6888	ML-HG	0.05	0	0	-2.1	0.4	0.18	0.42	0.07	0	0.26	0.67
NEO733 6824	ML-HG	1.32	1	1	-1.1	0.55	0.24	0.21	0.03	0	0.34	0.63
NEO732 6815	ML-HG	0.13	1	1	-4.5	0.34	0.24	0.42	0.03	0	0.36	0.61
NEO745 6790	ML-HG	0.45	1	1	-3	0.03	0.11	0.86	0.01	0	0.16	0.83
NEO856 6777	ML-HG	0.56	1	1	-2.3	0.39	0.22	0.39	0.02	0	0.19	0.79
NEO747 6729	ML-HG	0.25	1	0	-2.1	0.12	0.25	0.63	0	0	0.25	0.75
NEO568 6586	ML-HG	1.98	1	1	-4.4	0.57	0.19	0.24	0.01	0	0.26	0.73
NEO1 6585	ML-HG	0.02	0	0	0.8	0.09	0.15	0.76	0.1	0	0.27	0.63
NEO941 6372	ML-HG	0.14	1	1	-5	0.19	0.21	0.6	0.03	0	0.29	0.68
NEO570 6369	ML-HG	2.86	1	1	-0.6	0.72	0.15	0.13	0.01	0	0.19	0.8

NEO751 6343	ML-HG	0.30	1	1	-0.1	0.62	0.17	0.21	0.11	0	0.33	0.56
NEO852 6308	ML-HG	0.19	1	1	1	0.44	0.24	0.32	0.02	0	0.32	0.66
NEO855 6302	ML-HG	1.38	1	1	-0.4	0.76	0.13	0.11	0.02	0	0.27	0.71
NEO569 6142	ML-HG	0.67	1	1	-2.2	0.64	0.17	0.19	0.01	0	0.2	0.79
NEO598 6075	ML-HG	0.73	1	1	-0.5	0.09	0.15	0.76	0	0	0.17	0.83
NEO853 6047	ML-HG	1.96	1	1	-1.8	0.88	0.08	0.04	0.04	0	0.3	0.66
NEO3 5965	ML-HG	0.03	0	0	-1	0.15	0.16	0.69	0.14	0	0.28	0.58
NEO960 5926	ML-HG	0.15	1	1	-3.9	0.75	0.11	0.14	0.32	0	0.32	0.36
NEO645 5870	ML-HG	0.21	1	1	-0.4	0.85	0.09	0.06	0.09	0	0.33	0.58
NEO962 5786	ML-HG	0.04	0	0	-1.9	0.07	0.11	0.82	0.2	0	0.27	0.53
NEO601 5753	NL-EF	0.08	0	0	1.5	0	0.02	0.98	0.01	0	0.21	0.78
NEO790 5662	NL-EF	0.69	1	0	-4.2	0.28	0.16	0.56	0.25	0	0.29	0.46
NEO891 5661	NL-EF	0.60	1	1	1.1	0.69	0.14	0.17	0.11	0	0.35	0.54
NEO571 5534	NL-EF	0.06	0	0	-1.5	0.11	0.15	0.74	0.29	0	0.28	0.43
NEO23 5533	NL-EF	3.34	1	1	-3.4	0.19	0.27	0.54	0.17	0	0.35	0.48
NEO753 5531	NL-EF	0.16	1	1	-0.4	0	0.03	0.97	0.02	0	0.15	0.83
NEO942 5491	NL-EF	0.89	1	1	0.4	0.7	0.12	0.18	0.08	0	0.33	0.59
NEO29 5489	NL-EF	0.53	1	1	3.9	0.7	0.15	0.15	0.45	0	0.31	0.24
NEO564 5468	NL-EF	0.08	0	0	-2.6	0.35	0.18	0.47	0.13	0	0.32	0.55
NEO41 5462	NL-EF	0.02	0	0	1.4	0.05	0.13	0.82	0.13	0	0.32	0.55
NEO28 5459	NL-EF	0.92	1	0	2.7	0.43	0.2	0.37	0.22	0	0.3	0.48
NEO886 5457	NL-EF	0.27	1	1	-1.9	0.7	0.13	0.17	0.66	0	0.19	0.15
NEO866 5456	NL-EF	1.52	1	1	-0.5	0.84	0.09	0.07	0.72	0	0.19	0.09
NEO595 5452	NL-EF	0.22	1	0	-2.3	0.53	0.15	0.32	0.41	0	0.28	0.31
NEO757 5452	NL-EF	0.13	1	0	-0.3	0.05	0.11	0.84	0.03	0	0.27	0.7
NEO896 5446	NL-EF	0.12	1	0	-1.8	0.01	0.05	0.94	0.03	0	0.17	0.8

NEO945 5445	NL-EF	1.38	1	1	0.3	0.01	0.05	0.94	0.03	0	0.22	0.75
NEO888 5383	NL-EF	0.06	0	0	1.5	0.09	0.13	0.78	0.18	0	0.28	0.54
NEO933 5337	NL-EF	0.52	1	1	2.4	0.65	0.15	0.2	0.36	0	0.31	0.33
NEO744 5333	NL-EF	0.22	1	0	3.6	0.03	0.12	0.85	0.22	0	0.32	0.46
NEO795 5333	NL-EF	0.03	0	0	-1	0.13	0.16	0.71	0.2	0	0.33	0.47
NEO702 5263	NL-EF	0.15	1	1	-0.7	0.52	0.24	0.24	0.02	0	0.31	0.67
NEO7 5242	NL-EF	0.01	0	0	-1.1	0	0.04	0.96	0.03	0	0.23	0.74
NEO597 5210	NL-EF	0.18	1	0	0.5	0.1	0.14	0.76	0.31	0	0.29	0.4
NEO935 5187	NL-EF	5.03	1	1	-1.4	0.09	0.15	0.76	0.44	0	0.27	0.29
NEO865 5179	NL-EF	0.09	0	0	-3.2	0.01	0.06	0.93	0.06	0	0.23	0.71
NEO594 5174	NL-EF	0.05	0	0	-1.8	0.47	0.19	0.34	0.31	0	0.33	0.36
NEO961 5137	NL-EF	0.02	0	0	1.1	0.24	0.19	0.57	0.22	0	0.33	0.45
NEO599 5134	NL-EF	0.19	1	0	-2.6	0.02	0.11	0.86	0.08	0	0.24	0.68
NEO602 5134	NL-EF	0.09	0	0	-1.4	0.25	0.17	0.58	0.47	0	0.25	0.28
NEO566 5130	NL-EF	0.02	0	0	-3.2	0.44	0.17	0.39	0.38	0	0.28	0.34
NEO33 5128	NL-EF	0.05	0	0	-3.1	0.07	0.15	0.78	0.02	0	0.2	0.78
NEO898 5080	NL-EF	3.8	1	1	-4	0.51	0.13	0.36	0.21	0	0.3	0.49
NEO43 5067	NL-EF	0.11	1	0	1.9	0.05	0.14	0.81	0.17	0	0.36	0.47
NEO25 4956	NL-EF	0.36	1	1	1.1	0.09	0.16	0.75	0.08	0	0.33	0.59
NEO925 4947	NL-EF	0.29	1	0	-2.8	0	0.03	0.97	0.05	0	0.25	0.7
NEO943 4614	NL-EF	1.75	1	1	2.1	0.03	0.07	0.9	0.2	0	0.26	0.54
NEO580 4611	NL-EF	0.01	0	0	-1.3	0.03	0.09	0.88	0.09	0	0.28	0.63
NEO792 4493	NL-SP	0.25	1	1	-0.4	0.04	0.09	0.87	0.14	0	0.27	0.59
NEO876 4338	NL-SP	0.05	0	0	-2.4	0.01	0.08	0.91	0.16	0	0.31	0.53
NEO870 4240	NL-SP	0.58	1	1	2.2	0.23	0.26	0.51	0.37	0	0.33	0.3
NEO92 4188	NL-SP	0.61	1	1	-0.9	0.82	0.11	0.07	0.48	0	0.25	0.27

NEO737 4106	NL-SP	0.24	1	1	1	0.77	0.15	0.08	0.27	0	0.41	0.32
NEO738 4103	NL-SP	1.21	1	1	-0.3	0.92	0.05	0.03	0.73	0	0.17	0.1
NEO861 4102	NL-SP	0.38	1	1	0.5	0.79	0.11	0.1	0.74	0	0.17	0.09
NEO878 4026	NL-SP	0.28	1	1	-2.3	0.45	0.15	0.4	0.37	0	0.33	0.3
NEO872 3979	NL-SP	0.08	0	0	-0.1	0.62	0.15	0.23	0.44	0	0.3	0.26
NEO735 3972	NL-SP	0.67	1	1	2.5	0.78	0.11	0.11	0.76	0.01	0.16	0.07
NEO875 3970	NL-SP	0.18	1	0	1.6	0.07	0.11	0.82	0.18	0	0.27	0.55
NEO739 3965	NL-SP	1.88	1	1	0.9	0.01	0.04	0.95	0.02	0	0.15	0.83
NEO934 3809	NL-SP	0.08	0	0	2	0.19	0.15	0.66	0.29	0	0.31	0.4
NEO93 3735	NL-SP	1.98	1	1	-0.3	0.26	0.25	0.49	0.1	0	0.39	0.51
NEO860 3697	NL-SP	0.20	1	0	1.2	0.55	0.2	0.25	0.66	0	0.21	0.13
NEO857 3637	NL-SP	0.04	0	0	-1.9	0.09	0.17	0.74	0.26	0	0.31	0.43
NEO752 3589	NL-SP	2.09	1	1	1.8	0.64	0.12	0.24	0.31	0	0.31	0.38
NEO815 3471	NL-SP	0.11	1	0	4.3	0.48	0.24	0.28	0.33	0.04	0.42	0.21
NEO563 3350	NL-SP	0.81	1	1	3.6	0.69	0.14	0.17	0.59	0.09	0.19	0.13
NEO590 3290	NL-SP	1.01	1	1	2.9	0.85	0.07	0.08	0.84	0.01	0.11	0.04
NEO951 3242	NL-SP	0.45	1	0	1.7	0.01	0.06	0.93	0.2	0	0.34	0.46
NEO946 3094	NL-SP	1.24	1	1	2.2	0.8	0.1	0.1	0.28	0.01	0.37	0.34

5172
5173
5174

5175 References

- 5176 1. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.
5177 *Nature* **562**, 203–209 (2018).
- 5178 2. Price, A. L. *et al.* Long-range LD can confound genome scans in admixed
5179 populations. *American journal of human genetics* vol. 83 132–5; author reply 135–9
5180 (2008).

- 5181 3. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-
5182 based linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 5183 4. Pedersen, C. B. *et al.* The iPSYCH2012 case-cohort sample: new directions for
5184 unravelling genetic and environmental architectures of severe mental disorders. *Mol.*
5185 *Psychiatry* **23**, 6–14 (2018).
- 5186 5. Walsh, S. *et al.* The HirisPlex system for simultaneous prediction of hair and eye
5187 colour from DNA. *Forensic Sci. Int. Genet.* **7**, 98–115 (2013).

5188
5189

5190 4g) Calling chr17q21.31 KANSL1 Duplications in Ancient 5191 Genomes

5192
5193
5194

Alma S. Halgren¹, Andrés Ingason^{2,3}, and Peter H. Sudmant¹

5195
5196
5197
5198
5199

¹Department of Integrative Biology, University of California, Berkeley

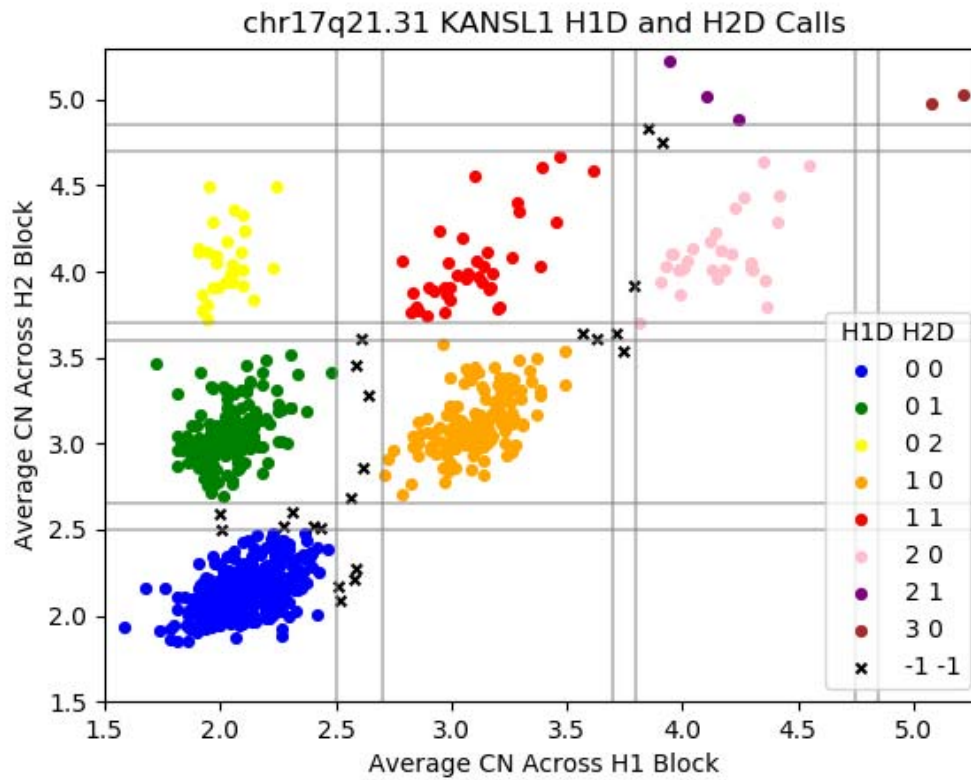
²Institute of Biological Psychiatry, Mental Health Services, Copenhagen University Hospital

³Lundbeck Foundation GeoGenetics Centre, GLOBE Institute, University of Copenhagen,
Copenhagen, Denmark

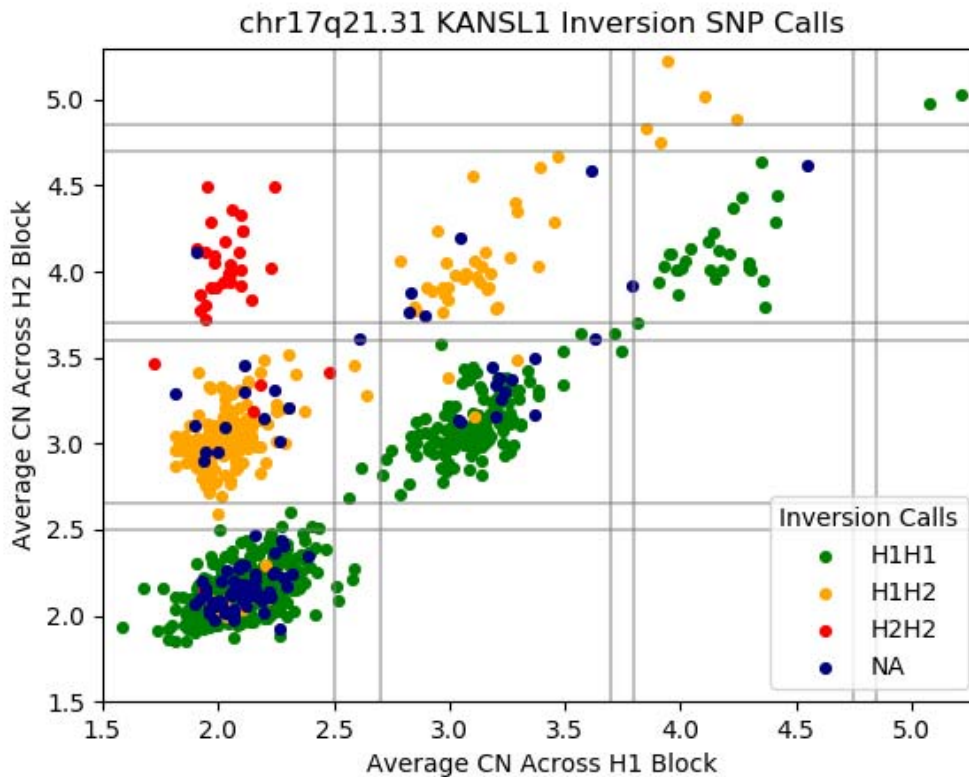
5200
5201
5202
5203
5204
5205
5206
5207
5208
5209
5210
5211
5212
5213
5214

To call the distinct H1 and H2 specific KANSL1 duplications at chr17q21.31 (H1D and H2D respectively), we first removed samples that were too noisy to be accurately genotyped. To do so, we calculated the standard deviation of genome-wide copy number (after removing the top and bottom fifth percentiles of copy number to exclude outliers). We chose standard deviation cutoffs based on a visual inspection of the data – 0.49 for the fastq samples and 0.804 for the bam samples. While bam and fastq samples exhibited genome-wide differences in coverage and noise resulting in these different cutoffs, the selected samples exhibit similar signatures at the KANSL1 locus. Together this resulted in a total of 1143 samples, 427 fastq samples and 716 bam samples. We then set a copy number cutoff of 10 at the KANSL1 locus (copy number signal above 10 is likely noise) and calculated the average copy number in the H1D and H2D coordinate blocks (Fig S4g.1A). We filtered for samples with an average copy number in both the H1D and H2D blocks between 1.5 and 5.25 (values outside of these bounds are likely noise). The samples clustered into duplication genotypes (colored below) with the exception of a handful of samples on the border of genotype groups for which we could not accurately assign a genotype (marked

5215 with 'x'). Figure S4g.1B shows the SNP inversion calls for these samples, which align with
5216 the duplication calls.
5217



5218
5219 **Figure S4g.1A.** *KANSL1* H1D and H2D calls grouped by genotype. The samples marked
5220 with an 'x' are considered ambiguous as they are between groups.



5221
 5222 **Figure S4g.1B.** *KANSL1* H1 and H2 inversion calls mapped atop H1D and H2D average
 5223 copy number values. The inversion and duplication calls align.
 5224

5225 4h) Calculating ancestral contributions to modern complex 5226 phenotypes

5227
 5228 Will Barrie¹ and Dan Lawson²

5229 ¹Zoology Department, University of Cambridge, UK.

5230 ²School of Mathematics and Integrative Epidemiology Unit, University of Bristol, UK.

5231

5232 Introduction

5233 Most studies that look at polygenic risk scores in ancient populations use genotypes of
 5234 ancient individuals, combined with effect sizes from modern GWAS studies, to reconstruct
 5235 risk scores for ancient individuals ¹. This involves exporting effect sizes across space and

5236 time, which is known to dramatically reduce the accuracy of the estimates ². Additionally,
5237 these scores are usually impossible to verify (except with specific phenotypes such as height
5238 where calibration is possible ^{3,4}, and don't necessarily measure what an ancient population
5239 contributed to phenotypic diversity in a modern population(s), especially when there has
5240 been selection or bottleneck events in between.

5241

5242 Here, we aim to use local ancestry information resulting from painting the UK Biobank (UKB)
5243 (Supplementary Note S3h) to estimate ancestral contributions to modern complex
5244 phenotypes, by calculating polygenic risk scores for each ancestry based on local painting
5245 results. This is a well-powered approach due to the large modern sample size, and is a more
5246 direct measure of the variants that a given ancestry contributed to the “white British” genetic
5247 landscape. Thus we can draw conclusions about the differing contributions of each ancestry
5248 to modern genetic risk, whether due to drift or selection. We use bootstrapping to test
5249 whether some ancestries are significantly and systematically over-represented for a
5250 phenotype, indicating selection. Additionally, we look at the ancestral haplotypic background
5251 of a high effect variant, ApoE4, which is implicated in Alzheimer's Disease ^{5,6}.

5252

5253 Methods

5254 We used effect size estimates from the UK Biobank Neale lab GWAS ⁷, and used 1,703 non-
5255 overlapping and approximately independent linkage disequilibrium (LD) blocks ⁸. For each
5256 block, we restricted the SNPs to those with a p-value less than the genome-wide
5257 significance threshold (5e-8), and from these chose the SNP with the lowest p-value. We
5258 then used these SNPs to calculate polygenic risk scores for each ancestry, using ancestry-
5259 specific ‘effect allele frequencies’ derived from the painting.

5260

5261 In order to calculate the effect allele frequency for a given ancestry $f_{\{anc,i\}}$ we used the
5262 formula:

5263

$$5264 \quad f_{\{anc,i\}} = \frac{\sum_j^M \text{Painting certainty}_{\{j,effect\}}}{\sum_j^M \text{Painting certainty}_{\{j,effect\}} + \sum_j^M \text{Painting certainty}_{\{j,alt\}}}$$

5265

5266 Where there are M individuals, and $\sum_j^M \text{Painting certainty}_{\{j,effect\}}$ is the sum of the
5267 painting probabilities for that ancestry of all effect alleles. This calculates an effect allele
5268 frequency for an ancestry which is weighted by the painting probabilities: if a haplotype with
5269 the effect allele was painted with low probability for that ancestry, it will contribute little to the

5270 calculation, and vice versa. One benefit of this approach is that because it only matters how
5271 effect alleles are painted relative to alternate alleles for an ancestry group, and differences in
5272 genome-wide painting averages between ancestries will not cause bias.

5273

5274 To calculate an ancestry-specific PRS we used an additive model, including a transformation
5275 as in Berg & Coop⁹ and in line with (Supplementary Note S4c). We derived standard
5276 deviations for each score by running a block bootstrap (1000 iterations) on (1) loci and (2)
5277 individuals. We calculated polygenic risk scores for 39 traits shown to be significantly over-
5278 dispersed across ancient populations beyond what would be expected under a null model of
5279 genetic drift (Supplementary Note S4c). For computational reasons, we used a random
5280 batch of 48,000 painted individuals to calculate the effect allele frequencies, which is
5281 sufficiently large to approximate the frequencies even for ancestries that are painted less.

5282

5283 Our calculations were limited to the 549,323 SNPs used in the painting of the UKB
5284 (Supplementary Note S3h). This is expected to reduce predictive power compared to using
5285 the full set of imputed SNPs in the UKB, but only slightly¹⁰. There was a ~15% decrease in
5286 the number of SNPs included per phenotype in the PRS calculation compared with the
5287 imputed data.

5288

5289 To test the ancestral background of a single variant, APOE4, we calculated the average
5290 painting score for each ancestry at all sites on the chromosome of haplotypes containing the
5291 effect allele. This makes it clear when there is an excess of a particular ancestry at the site
5292 of interest.

5293

5294 Results

5295 Our results tell us about the ancestral contribution to modern phenotypes in the white British
5296 population (Figure S4h.1, Figure S4h.2), and we stress we are not making claims about the
5297 phenotypes of ancient populations.

5298

5299 We find that Yamnaya, CHG and EHG ancestral contributions (which together form a
5300 'steppe' component) have relatively high scores for height, whereas Farmers and WHG
5301 ancestral contributions have relatively low scores. This accords with most previous studies
5302 ^{3,11,12} but not all ¹³. EHG and Yamnaya both score highly for body mass and basal metabolic
5303 rate.

5304

5305 Hair and skin pigmentation show significant differences between the ancestral contributions,
5306 with risk scores for skin colour for the three hunter-gatherer ancestries being higher (i.e.
5307 darker) than Farmer and Steppe (as in Ju and Mathieson ¹⁴). On the other hand, traits
5308 related to malignant neoplasms of skin show higher scores for the Farmer ancestral
5309 contribution; while Farmer and Yamnaya ancestral contributions have higher scores for
5310 blonde and light brown hair, with the hunter-gatherer ancestries showing higher scores for
5311 dark brown. CHG is the only ancestral contribution which stands out as having a high risk
5312 score for black hair.

5313

5314 Intriguingly, the WHG ancestral component has strikingly high scores for traits related to
5315 cholesterol, blood pressure and diabetes, both when bootstrapping individuals and loci ^{cf. 13}.
5316 In terms of psychiatric traits, the Farmer component scores highest for anxiety, guilty
5317 feelings, and irritability.

5318

5319 Our two bootstrapping methods mean slightly different things. Individuals in the UKB are
5320 related through shared genealogies, and so by bootstrapping over non-independent
5321 individuals (Figure S4h.1) we are testing the consistency of the signal within the population.
5322 From this bootstrapping exercise we can conclude whether a difference in allele frequencies
5323 in ancient populations contributed to phenotypic variation today. Unsurprisingly, with a large
5324 enough sample size most phenotypes will show differences in ancestral contributions for
5325 this, usually due to drift or founder effects. However, this goes further than just reporting risk
5326 scores for ancient populations, because we are looking directly at coalescent tracts in the
5327 British population. We can conclude that “ancestry X contributes higher genetic risk for
5328 phenotype Y in the test population”. On the other hand, because we have used independent
5329 LD blocks to select SNPs to include in the PRS calculation, the requirement for
5330 independence is met when we bootstrap with loci (Figure S4h.2). A positive result here is
5331 therefore much stronger, showing a systematic over/under-representation of an ancestry at
5332 loci affecting a given trait, beyond what is expected given the correlation among individuals.
5333 This points towards selection as an explanation.

5334

5335 The effect/risk allele (rs429358, n=127,760) of ApoE4 is preferentially painted as
5336 WHG/EHG, with a clear depletion of other ancestries (especially Farmer) at this locus
5337 compared to the genome-wide average (Figure S4h.3). This indicates that this allele was
5338 contributed at least in part by hunter-gatherer ancestry into modern (British) populations,
5339 above what we would expect by chance.

5340

5341 Discussion

5342 The methods here directly link genetic contributions from pre-defined ancestries to complex
5343 phenotypes in modern people. For most traits, each ancestry contributed differently to the
5344 modern genetic landscape, with some conveying enhanced or reduced risk either due to drift
5345 (including population bottlenecks/founder events) or selection. Because gradients exist in
5346 these ancestries across the British Isles and further afield ([Supplementary Note S3h](#)), these
5347 differing risk scores indicate how geographically heterogeneous ancestry distributions may
5348 contribute to differing genetic risk profiles, in addition to other factors such as geography,
5349 socio-economic status etc.

5350

5351 A caveat for all studies involving polygenic risk calculation is that they rely on effect size
5352 estimates from an original GWAS which may be affected by population stratification in the
5353 GWAS panel, even when it has apparently been controlled for. This seems to be less of a
5354 problem in the UKB than in previous GWAS studies ¹⁵, but should be kept in mind. One
5355 benefit of our approach is that there is no requirement to export these risk scores across
5356 time and space: we are using effect sizes estimated from the modern population to calculate
5357 ancestral contributions to the same modern population.

5358

5359 ApoE4 is an isoform of the APOE gene, resulting from linkage disequilibrium between two
5360 SNPs, rs429358 and rs7412 ¹⁶, and associated with increased risk for metabolic, vascular
5361 and neurodegenerative diseases in adulthood ¹⁷. It may provide some enhanced cognitive
5362 ability in children and young adults ¹⁸ and other health and immunity benefits, particularly in
5363 highly infected environments ^{e.g. 19}. There are several lines of evidence suggesting a link
5364 between the evolution of diet and the ApoE isoforms: ϵ 2 and ϵ 3 alleles are associated with
5365 lower levels of blood cholesterol ^{20,21}, while ϵ 4 is associated with higher levels, leading some
5366 to speculate that the derived ϵ 3 allele is 'meat-adaptive' ^{22,23}. In a study of South Americans,
5367 there was a five-fold increase in the ApoE4 allele in hunter-gatherers versus horticulturalists
5368 ²⁴, potentially because the immune benefits outweighed the advantages of low blood
5369 cholesterol ²⁵. Generally, ϵ 4 prevalence is higher in indigenous foraging groups such as the
5370 Pygmies, Khoi San, Papuans and some Native Americans, while ϵ 3 is most frequent in
5371 populations with a long-established agricultural economy ²⁶. Finally, ApoE4 is implicated in
5372 higher blood vitamin D levels ²⁷.

5373

5374 The ϵ 4 variant has been shown to be ancestral in humans ²⁸. There is a linear increasing
5375 trend in ϵ 4 prevalence from South to North in Europe, with Sardinians showing the lowest
5376 prevalence ²⁹⁻³¹, while there is a more than two-fold increase in Nordic versus Mediterranean

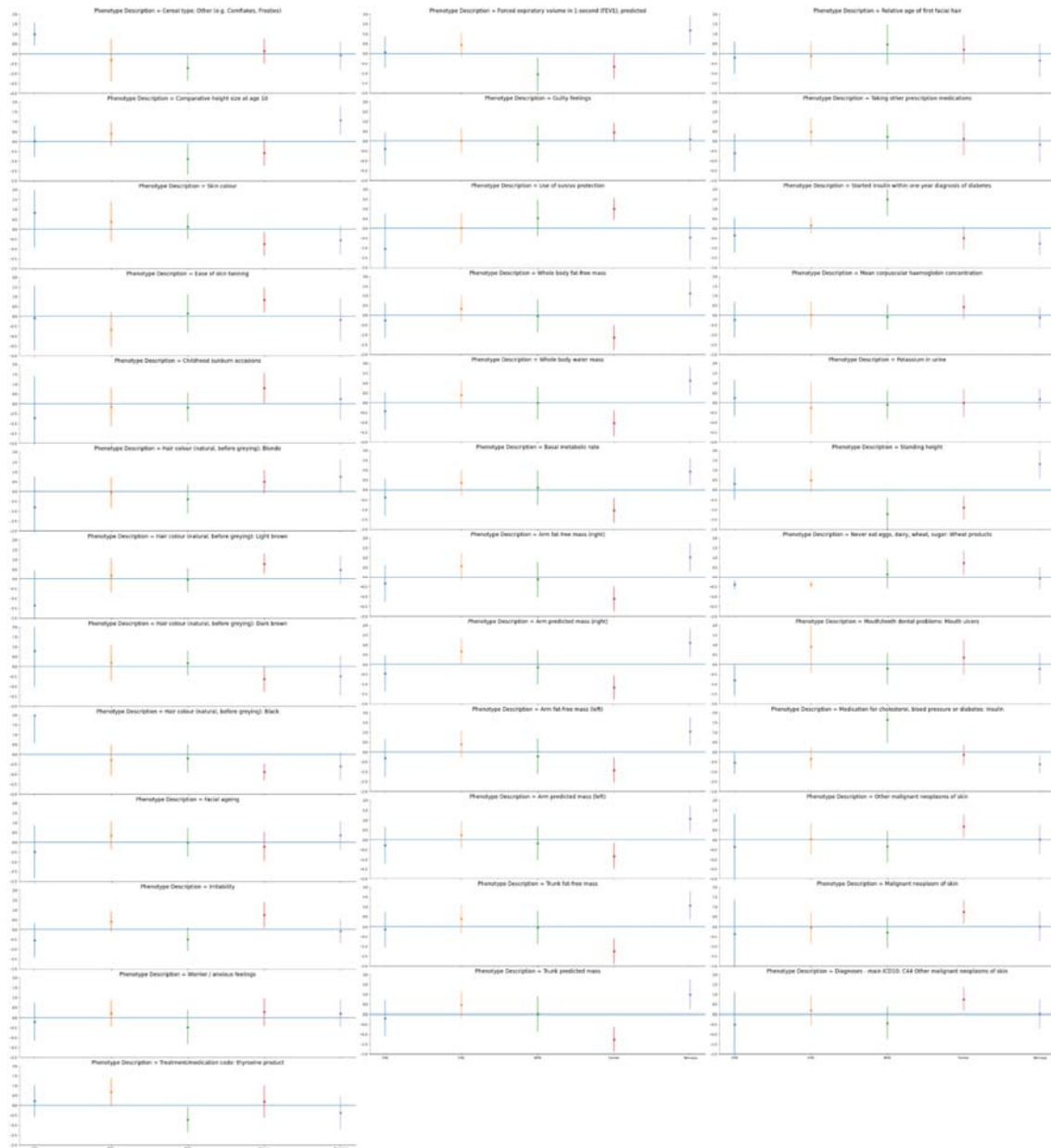
5377 countries³². Sardinians are an unusual population, having the highest level of neolithic
 5378 farmer ancestry of all modern European populations³³. In this light, differences in genome-
 5379 wide ancestry proportions between northern (high WHG/EHG, low Farmer) and southern
 5380 Europe (high Farmer, low WHG/EHG) (Supplementary Note S3h) may explain at least part
 5381 of the differences in frequency of the $\epsilon 4$ variant and subsequent AD genetic risk.

5382 Figures/tables

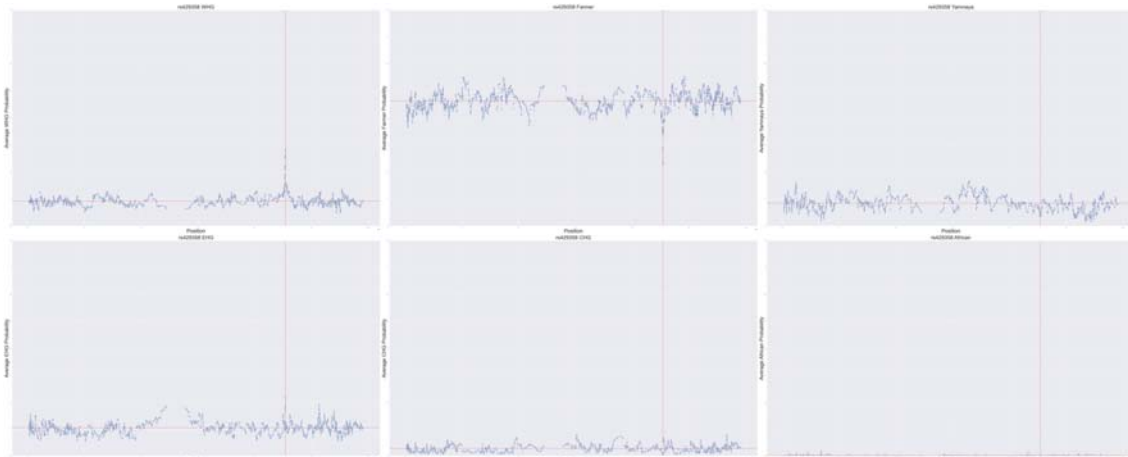


5383 **Figure S4h.1. Ancestry-specific polygenic risk scores with 95% confidence intervals derived**
 5384 **from bootstrapping individuals for phenotypes shown to be significantly over-dispersed**
 5385 **between ancient populations.** Confidence intervals were calculated by re-running PRS calculation
 5386 on random batches of 48,000 individuals, with replacement (1000 iterations), while keeping all other
 5387 annotations intact. Here we show 2 x standard deviation error bars, expected to represent ~95%
 5388

5389 confidence interval under a normal distribution. Bootstrapping individuals tests the extent to which
 5390 ancestry X contributed higher genetic risk for phenotype Y in a given population, either due to drift or
 5391 selection.
 5392



5393
 5394 **Figure S4h.2. Ancestry-specific polygenic risk scores with 95% confidence intervals derived**
 5395 **from bootstrapping loci for phenotypes shown to be significantly over-dispersed between**
 5396 **ancient populations.** Confidence intervals were calculated by bootstrapping independent loci from
 5397 separate LD blocks (1000 iterations), while keeping all other annotations intact. Here we show 2 x
 5398 standard deviation error bars, expected to represent ~95% confidence interval under a normal
 5399 distribution. Bootstrapping loci tests whether there is a systematic bias towards an ancestry for a
 5400 given phenotype across all significant SNPs, possibly indicating selection.
 5401



5402
 5403 **Figure S4h.3. Average painting score for each ancestry at all sites on chromosome 19 of**
 5404 **haplotypes containing the effect allele for ApoE4 (rs429358, n=127,760).** Vertical red line
 5405 indicates the position of the SNP of interest; horizontal red line indicates the average painting score
 5406 for that ancestry for haplotypes containing the effect allele across the entirety of chromosome 19.
 5407 There is a clear excess of WHG/EHG ancestry and a depletion of Farmer ancestry at this locus.

5408 **References**

5409

5410 1. Irving-Pease, E. K., Muktopavela, R., Dannemann, M. & Racimo, F. Quantitative
 5411 Human Paleogenetics: What can Ancient DNA Tell us About Complex Trait Evolution?
 5412 *Front. Genet.* **12**, 703541 (2021).

5413 2. Duncan, L. *et al.* Analysis of polygenic risk score usage and performance in diverse
 5414 human populations. *Nat. Commun.* **10**, 3328 (2019).

5415 3. Cox, S. L., Ruff, C. B., Maier, R. M. & Mathieson, I. Genetic contributions to variation
 5416 in human stature in prehistoric Europe. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 21484–
 5417 21492 (2019).

5418 4. Cox, S. L. *et al.* Predicting skeletal stature using ancient DNA. *American Journal of*
 5419 *Biological Anthropology* **177**, 162–174 (2022).

5420 5. Corder, E. H. *et al.* Gene dose of apolipoprotein E type 4 allele and the risk of
 5421 Alzheimer's disease in late onset families. *Science* **261**, 921–923 (1993).

5422 6. Strittmatter, W. J. *et al.* Binding of human apolipoprotein E to synthetic amyloid beta
 5423 peptide: isoform-specific effects and implications for late-onset Alzheimer disease.
 5424 *Proc. Natl. Acad. Sci. U. S. A.* **90**, 8098–8102 (1993).

- 5425 7. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.
5426 *Nature* **562**, 203–209 (2018).
- 5427 8. Berisa, T. & Pickrell, J. K. Approximately independent linkage disequilibrium blocks in
5428 human populations. *Bioinformatics* **32**, 283–285 (2016).
- 5429 9. Berg, J. J. & Coop, G. A population genetic signal of polygenic adaptation. *PLoS*
5430 *Genet.* **10**, e1004412 (2014).
- 5431 10. Choi, S. W. & O'Reilly, P. F. PRSice-2: Polygenic Risk Score software for
5432 biobank-scale data. *Gigascience* **8**, (2019).
- 5433 11. Mathieson, I. *et al.* Genome-wide patterns of selection in 230 ancient
5434 Eurasians. *Nature* **528**, 499–503 (2015).
- 5435 12. Martiniano, R. *et al.* The population genomics of archaeological transition in
5436 west Iberia: Investigation of ancient substructure using imputation and haplotype-based
5437 methods. *PLoS Genet.* **13**, e1006852 (2017).
- 5438 13. Marnetto, D. *et al.* Ancestral contributions to contemporary European complex
5439 traits. *bioRxiv* 2021.08.03.454888 (2021) doi:10.1101/2021.08.03.454888.
- 5440 14. Ju, D. & Mathieson, I. The evolution of skin pigmentation-associated variation
5441 in West Eurasia. *Proc. Natl. Acad. Sci. U. S. A.* **118**, (2021).
- 5442 15. Berg, J. J. *et al.* Reduced signal for polygenic adaptation of height in UK
5443 Biobank. *Elife* **8**, (2019).
- 5444 16. Rall, S. C., Jr, Weisgraber, K. H. & Mahley, R. W. Human apolipoprotein E.
5445 The complete amino acid sequence. *J. Biol. Chem.* **257**, 4171–4178 (1982).
- 5446 17. de-Almada, B. V. P. *et al.* Protective effect of the APOE-e3 allele in
5447 Alzheimer's disease. *Braz. J. Med. Biol. Res.* **45**, 8–12 (2012).
- 5448 18. Tuminello, E. R. & Han, S. D. The apolipoprotein e antagonistic pleiotropy
5449 hypothesis: review and recommendations. *Int. J. Alzheimers. Dis.* **2011**, 726197 (2011).
- 5450 19. Oriá, R. B., Patrick, P. D., Blackman, J. A., Lima, A. A. M. & Guerrant, R. L.
5451 Role of apolipoprotein E4 in protecting children against early childhood diarrhea
5452 outcomes and implications for later development. *Med. Hypotheses* **68**, 1099–1107

- 5453 (2007).
- 5454 20. Petkeviciene, J. *et al.* Associations between apolipoprotein E genotype, diet,
5455 body mass index, and serum lipids in Lithuanian adult population. *PLoS One* **7**, e41525
5456 (2012).
- 5457 21. Carvalho-Wells, A. L. *et al.* Interactions between age and apoE genotype on
5458 fasting and postprandial triglycerides levels. *Atherosclerosis* **212**, 481–487 (2010).
- 5459 22. Finch, C. E. & Stanford, C. B. Meat-adaptive genes and the evolution of
5460 slower aging in humans. *Q. Rev. Biol.* **79**, 3–50 (2004).
- 5461 23. Allen, J. S., Bruss, J. & Damasio, H. The aging brain: the cognitive reserve
5462 hypothesis and hominid evolution. *Am. J. Hum. Biol.* **17**, 673–689 (2005).
- 5463 24. Reales, G. *et al.* A tale of agriculturalists and hunter-gatherers: Exploring the
5464 thrifty genotype hypothesis in native South Americans. *Am. J. Phys. Anthropol.* **163**,
5465 591–601 (2017).
- 5466 25. Trumble, B. C. *et al.* Apolipoprotein E4 is associated with improved cognitive
5467 function in Amazonian forager-horticulturalists with a high parasite burden. *FASEB J.*
5468 **31**, 1508–1515 (2017).
- 5469 26. Corbo, R. M. & Scacchi, R. Apolipoprotein E (APOE) allele distribution in the
5470 world. Is APOE*4 a 'thrifty' allele? *Ann. Hum. Genet.* **63**, 301–310 (1999).
- 5471 27. Huebbe, P. *et al.* APOE ϵ 4 is associated with higher vitamin D levels in
5472 targeted replacement mice and humans. *FASEB J.* **25**, 3262–3270 (2011).
- 5473 28. Fullerton, S. M. *et al.* Apolipoprotein E variation at the sequence haplotype
5474 level: implications for the origin and maintenance of a major human polymorphism. *Am.*
5475 *J. Hum. Genet.* **67**, 881–900 (2000).
- 5476 29. Corbo, R. M., Scacchi, R., Mureddu, L., Mulas, G. & Alfano, G. Apolipoprotein
5477 E polymorphism in Italy investigated in native plasma by a simple polyacrylamide gel
5478 isoelectric focusing technique. Comparison with frequency data of other European
5479 populations. *Ann. Hum. Genet.* **59**, 197–209 (1995).
- 5480 30. Lucotte, G., Loirat, F. & Hazout, S. Pattern of gradient of apolipoprotein E

- 5481 allele *4 frequencies in western Europe. *Hum. Biol.* **69**, 253–262 (1997).
- 5482 31. Adler, G. *et al.* Bosnian study of APOE distribution (BOSAD): a comparison
5483 with other European populations. *Ann. Hum. Biol.* **44**, 568–573 (2017).
- 5484 32. Trumble, B. C. & Finch, C. E. THE EXPOSOME IN HUMAN EVOLUTION:
5485 FROM DUST TO DIESEL. *Q. Rev. Biol.* **94**, 333–394 (2019).
- 5486 33. Chiang, C. W. K. *et al.* Genomic history of the Sardinian population. *Nat.*
5487 *Genet.* **50**, 1426–1434 (2018).

5488

5489 4i) Pathogenic structural variants in ancient vs. modern-day 5490 humans

5491

5492 Alma S. Halgren¹, Andrés Ingason^{2,3}, and Peter H. Sudmant¹

5493

5494 ¹Department of Integrative Biology, University of California, Berkeley

5495 ²Institute of Biological Psychiatry, Mental Health Services, Copenhagen University Hospital

5496 ³Lundbeck Foundation GeoGenetics Centre, GLOBE Institute, University of Copenhagen,
5497 Copenhagen, Denmark

5498

5499

5500 Rare, recurrent copy-number variants (CNVs) are known to cause neurodevelopmental
5501 disorders and are associated with a range of psychiatric and physical traits with variable
5502 expressivity and incomplete penetrance ^{1,2}. We examined 50 regions susceptible to recurrent
5503 CNV known to be the most prevalent drivers of human developmental pathologies ³ in 1442
5504 ancient Eurasians and 1093 modern human populations (for comparison) to understand the
5505 prevalence of pathogenic structural variants over time.

5506

5507 This analysis examines 1442 ancient humans from primarily West Eurasia and Central Asia
5508 as well as 1093 publicly-available high-coverage modern human genomes encompassing
5509 136 populations worldwide (from the Human Genome Diversity Project ⁴ and the Simons
5510 Genome Diversity Project ⁵). In modern humans and 690 ancient individuals with fastq files
5511 available, paired-end Illumina reads were mapped to the human reference genome GRCh38
5512 with BWA-MEM ⁶. In the remaining 984 samples, only BAM files that had been mapped to
5513 hg19 were available. Of note, in these 984 samples, the filtering of duplicate reads resulted
5514 in the absence of signal over segmental duplications. Nonetheless, we were able to

5515 characterise structural variants intersecting unique sequences in these samples. The large
5516 putatively pathogenic loci which we focused on in this analysis generally consist of unique
5517 sequences flanked by segmental duplications. 232 samples were removed due to low
5518 coverage and genotype yield out of 1674 total ancient samples, leaving 1442 samples for
5519 the final analysis (601 from the fastq hg38 dataset and 841 from the BAM hg19 dataset). In
5520 all samples, average read depth in 1kb sliding genomic windows was extracted from the
5521 subsequent BAM files with pysamstats ⁷. All alternate haplotypes were removed prior to
5522 mapping. To approximate copy number from read depth, we masked tandem repeats
5523 (Tandem Repeat Finder ⁸), corrected read-depth estimates for underlying GC content
5524 (similar method to Sudmant et al. 2013 ⁹), and normalised by median read depth per
5525 individual. We implemented a Gaussian Hidden Markov Model ¹⁰ to call structural variants
5526 from read depth.

5527

5528 We identified CNVs in ancient individuals at ten loci using digital Comparative Genomic
5529 Hybridization ¹¹ (Table S4i.1; Figures S4i.1-S4i.20). Although most of the observed CNVs
5530 (including duplications at 15q11.2 and *CHRNA7*, and CNVs spanning parts of the TAR locus
5531 and 22q11.2 distal) have not been unambiguously associated with disease in large studies,
5532 the identified CNVs include deletions and duplications that have been associated with
5533 developmental delay, dysmorphic features, and neuropsychiatric abnormalities such as
5534 autism (most notably at 1q21.1, 3q29, 16p12.1 and the DiGeorge/VCFS locus, but also
5535 deletions at 15q11.2 and duplications at 16p13.11). However, phenotypes and risk
5536 associated with these structural variants vary widely, and recent population-based studies
5537 ^{12,13} suggest that they may be more common in the general population than previously
5538 thought ^{1,2}. Overall, the carrier frequency in ancient samples is similar to that reported in the
5539 UK Biobank (1.25% vs 1.6% at 15q11.2 and *CHRNA7* combined, and 0.8% vs 1.1% across
5540 the remaining loci combined) ¹². These results suggest that large, recurrent CNVs that can
5541 lead to several pathologies were present in ancient populations at similar frequencies as
5542 modern populations.

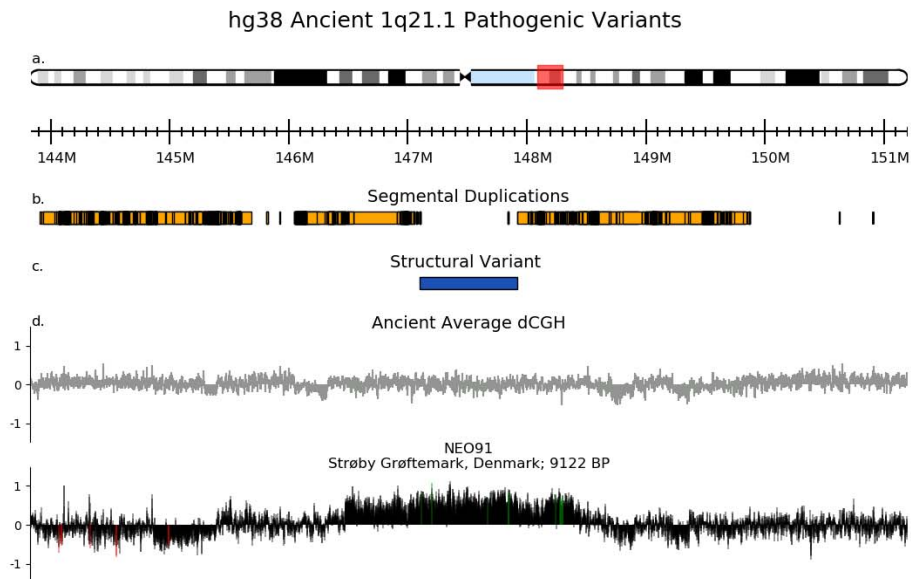
5543

5544

Region	Ancient Deletions	SGDP & HGDP Deletions	UK Biobank Deletions	Ancient Duplications	SGDP & HGDP Duplications	UK Biobank Duplications
1q21.1	0 (0%)	0 (0%)	113 (0.027%)	1 (0.069%)	0 (0%)	177 (0.042%)
3q29	1 (0.069%)	0 (0%)	9 (0.002%)	1 (0.069%)	0 (0%)	5 (0.001%)
15q11.2	4 (0.28%)	2 (0.18%)	1664 (0.39%)	10 (0.69%)	9 (0.82%)	2041 (0.48%)
15q11q13 (BP3-BP4)	1 (0.069%)	0 (0%)	16 (0.004%)	0 (0%)	0 (0%)	53 (0.013%)
15q13.3 (CHRNA7)	0 (0%)	1 (0.09%)	10 (0.002%)	4 (0.28%)	8 (0.73%)	3031 (0.72%)
16p12.1	1 (0.069%)	0 (0%)	246 (0.058%)	1 (0.069%)	0 (0%)	202 (0.048%)
16p13.11	1 (0.069%)	0 (0%)	131 (0.031%)	4 (0.28%)	0 (0%)	828 (0.2%)
22q11.2 (distal)*	4 (0.28%)	0 (0%)	N/A	13 (0.90%)	6 (0.55%)	N/A
DiGeorge-VCFS	0 (0%)	0 (0%)	10 (0.0024%)	1 (0.069%)	0 (0%)	280 (0.066%)
TAR*	1 (0.069%)	0 (0%)	N/A	2 (0.14%)	0 (0%)	N/A

5545
5546
5547
5548
5549
5550
5551
5552
5553
5554
5555
5556
5557

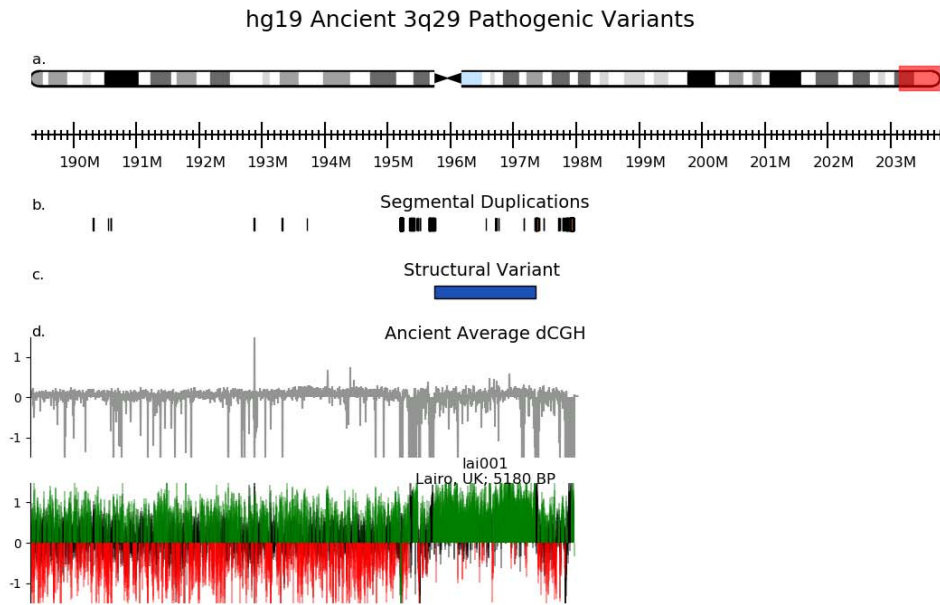
Table S4i.1. Table of 10 pathogenic loci (out of 50 examined) with structural variants identified in the ancient dataset. For each dataset, we report the prevalence of each SV as both the number of individuals identified as well as the percentage in each population (ancient dataset: 1442 samples; modern human Simons Genome Diversity Project (SGDP) ⁵ and Human Genome Diversity Project (HGDP) ⁴ dataset: 1093 samples; UK Biobank dataset: 421268 samples).
*Ancient SVs do not span entire locus, and therefore we cannot compare the UK Biobank frequencies for these loci (which span the entirety of the locus) to the ancient frequencies.



5558
5559
5560
5561
5562
5563
5564

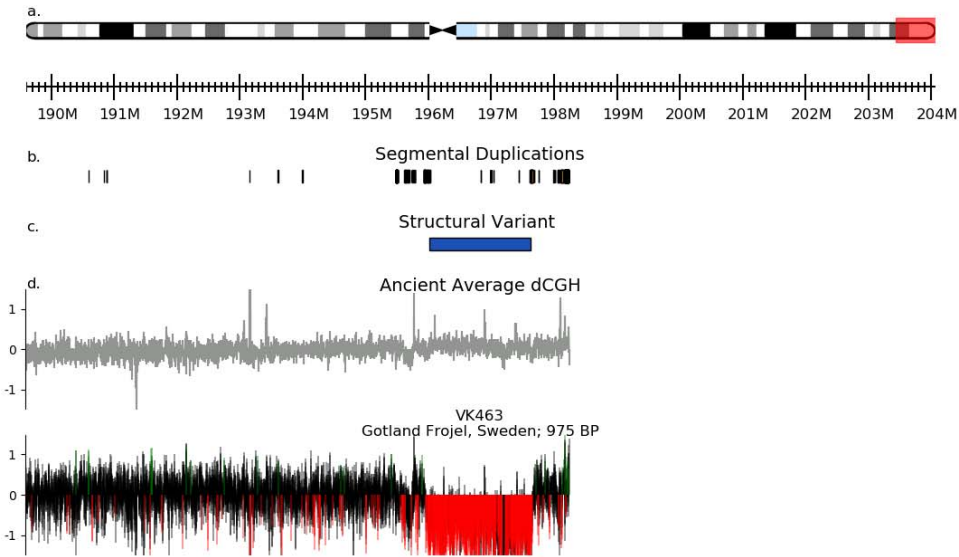
Figure S4i.1. Microduplication of an ancient individual at the 1q21.1 locus from the hg38 fastq dataset. This 2 Mb microduplication is known to be associated with increased risk of neurological and psychiatric problems, delayed development, Tetralogy of Fallot, and micro/macrocephaly ¹⁴⁻¹⁶. a. Chromosome 1 G-bands and axis bar (in Mbp). b. Segmental Duplications of >1000 bases of non-RepeatMasked ¹⁷ sequence from the UCSC Genome Browser ¹⁸. c. Indices of the structural variant from the literature ¹⁷. d. The average dCGH (in

5565 gray) is computed as follows: first, a non-noisy individual for this locus is selected; then, for
5566 every other individual, the \log_2 ratio of its copy number values over the non-noisy individuals'
5567 copy number values is calculated; finally, the average of these ratios is depicted. The \log_2
5568 ratio of NEO91's copy number values over those of the non-noisy individual is shown below
5569 (green = the ratio is at least 1.5 the standard deviation above the average dCGH; red = 1.5
5570 std. dev. below). While there is little deviation from the "norm" copy number at this locus on
5571 average, the dCGH of NEO91 is significantly greater than the average and indicates a
5572 duplication.



5573
 5574 **Figure S4i.2. Microduplication of an ancient individual at the 3q29 locus from the hg19**
 5575 **BAM dataset.** At 3q29, microdeletions and microduplications are associated with speech
 5576 and developmental delay, cleft palate, microcephaly, and increased risk for psychiatric
 5577 disorders^{19,20}. No modern individuals in the SGDP or HGDP datasets present with a
 5578 microdeletion or microduplication at 3q29.
 5579
 5580
 5581
 5582
 5583

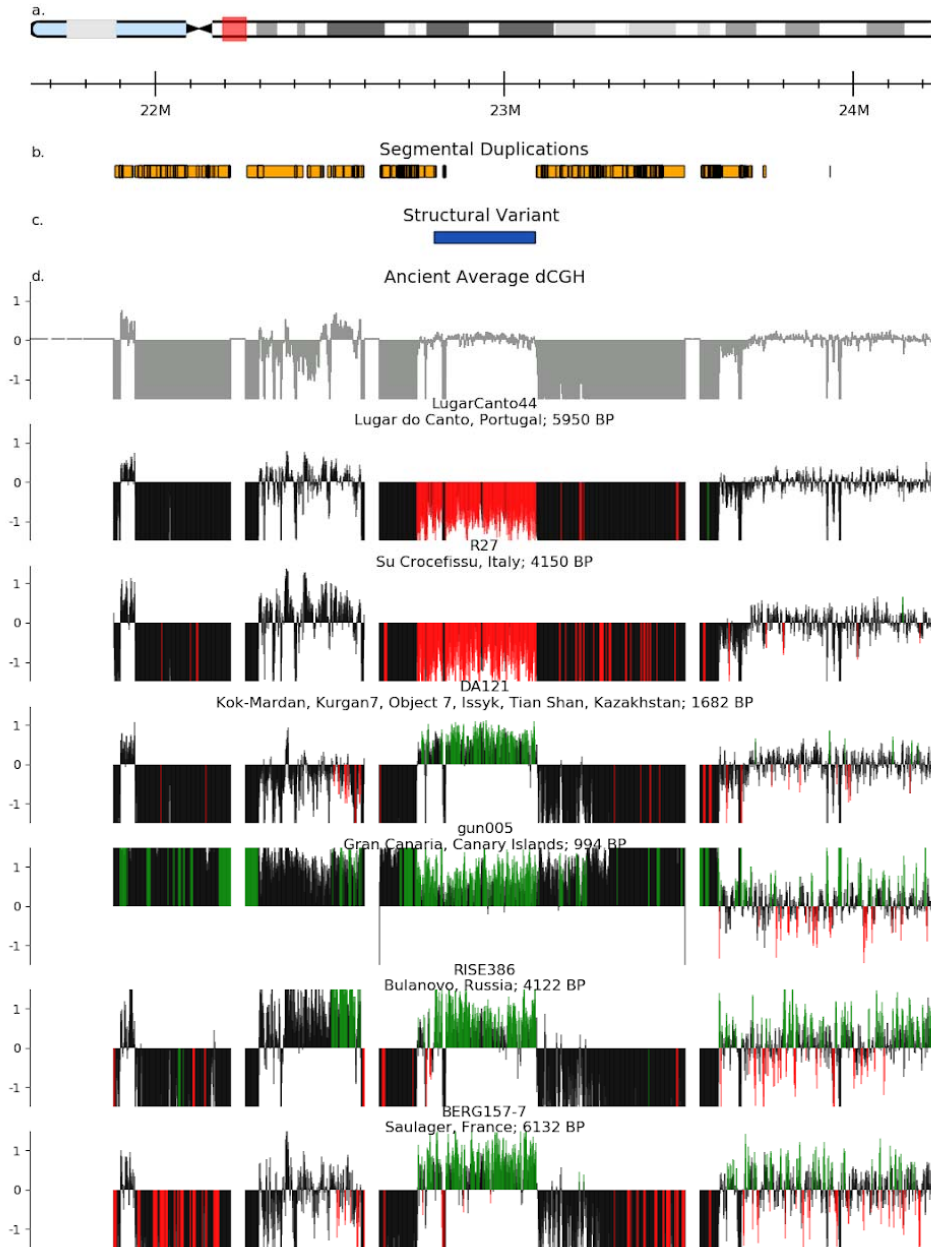
hg38 Ancient 3q29 Pathogenic Variants



5584
5585
5586

Figure S4i.3. Microdeletion of an ancient individual at the 3q29 locus from the hg38 fastq dataset.

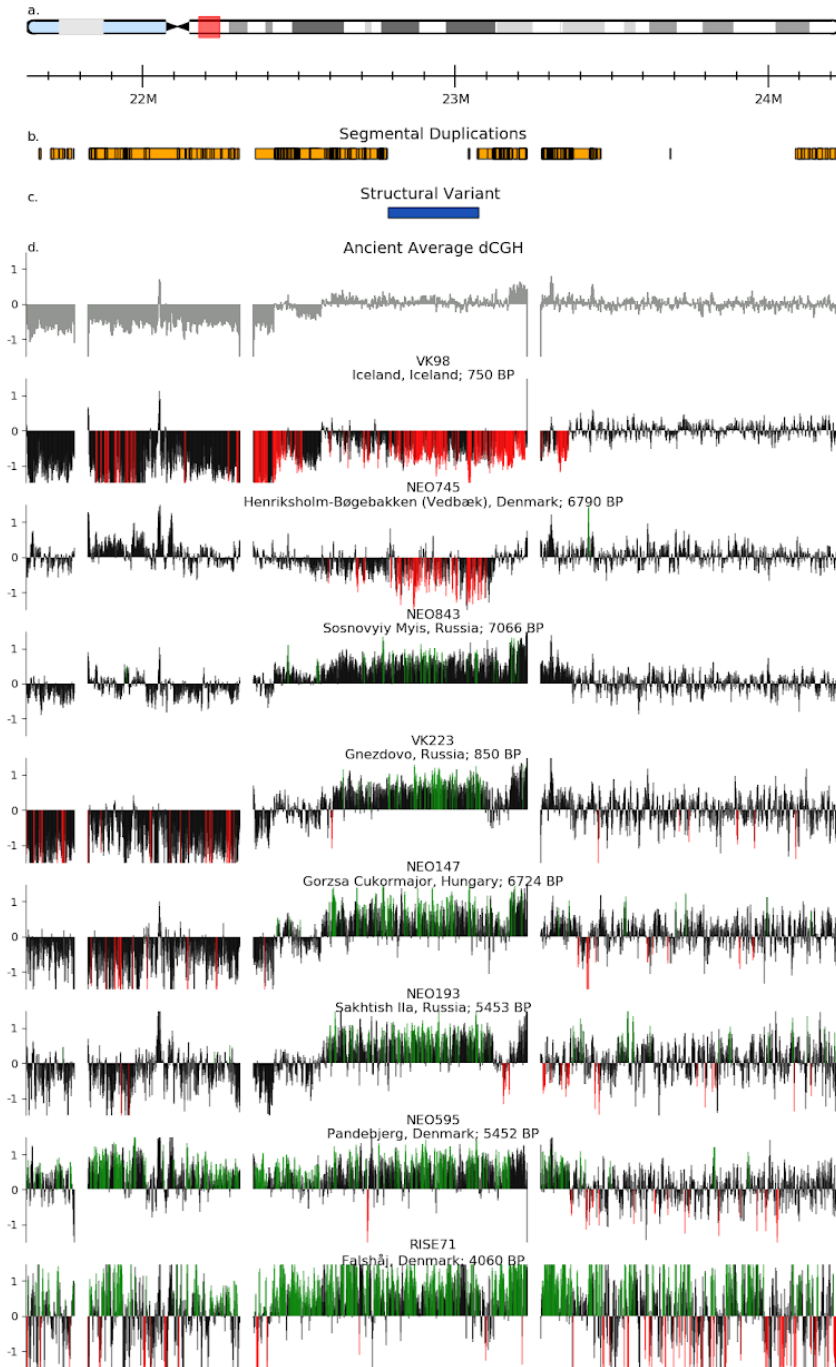
hg19 Ancient 15q11.2 Pathogenic Variants



5587
 5588
 5589
 5590
 5591
 5592
 5593
 5594

Figure S4i.4. Two microdeletions and four microduplications are present at the 15q11.2 locus from the hg19 BAM dataset. Although the microdeletion is considered to confer modest risk of schizophrenia, epilepsy, learning problems, and ADHD, the microduplication is largely considered to be benign²¹⁻²⁴. The blocks flanking the structural variant with copy number 0 are where duplicated reads have been filtered from the dataset.

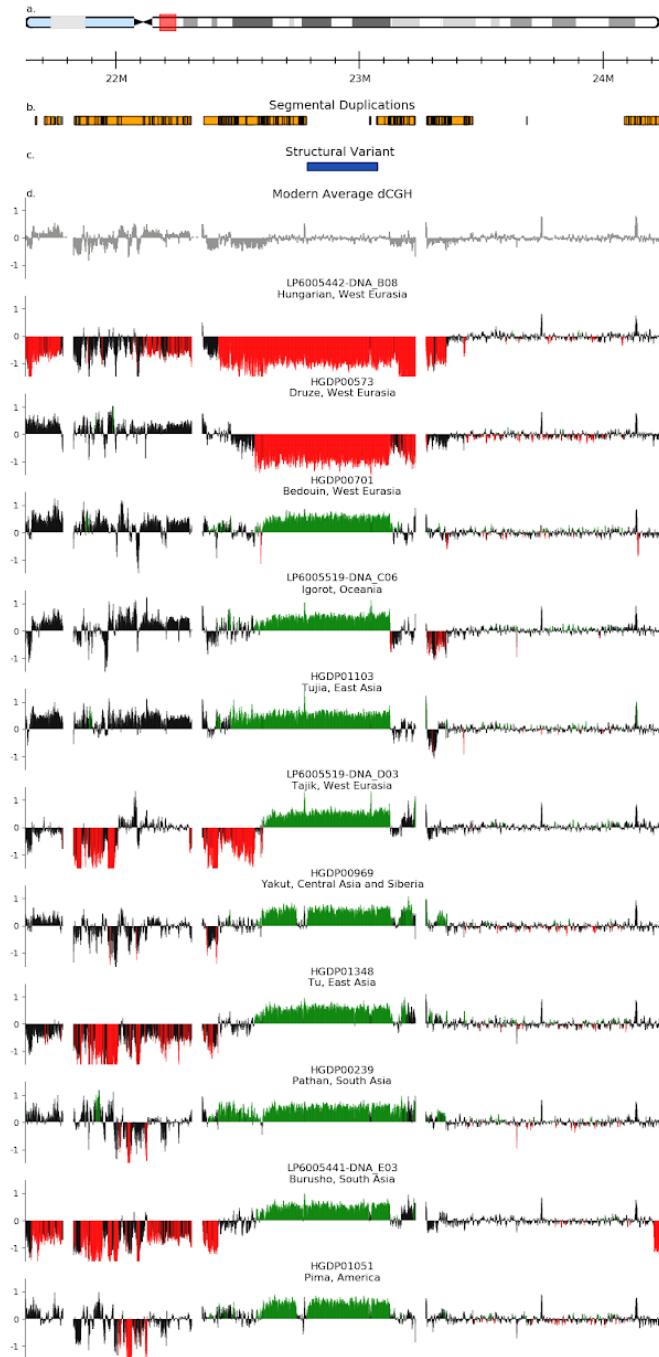
hg38 Ancient 15q11.2 Pathogenic Variants



5595
5596
5597
5598

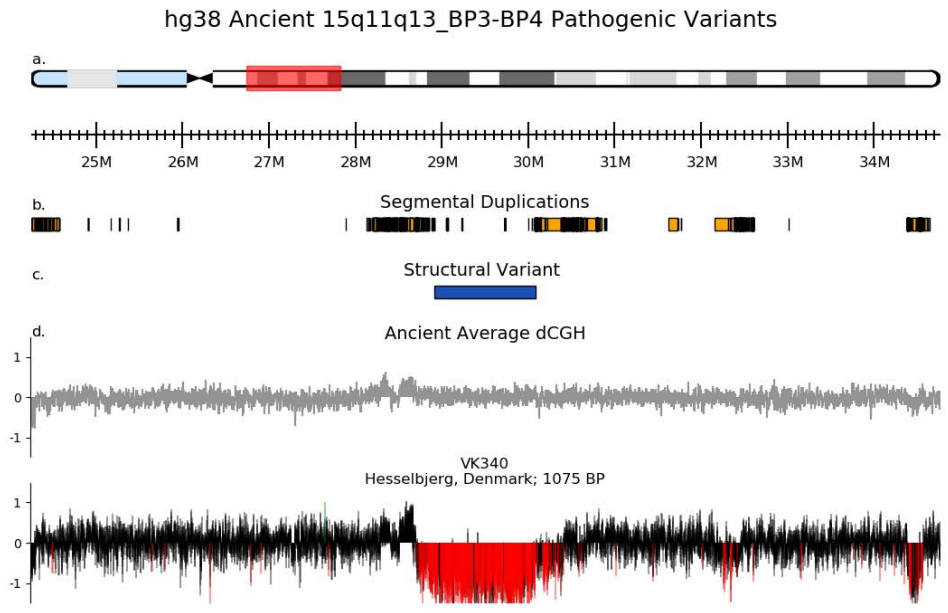
Figure S4i.5. Two microdeletions and six microduplications are present at the 15q11.2 locus from the hg38 fastq dataset.

hg38 Modern 15q11.2 Pathogenic Variants



5599
5600
5601
5602
5603

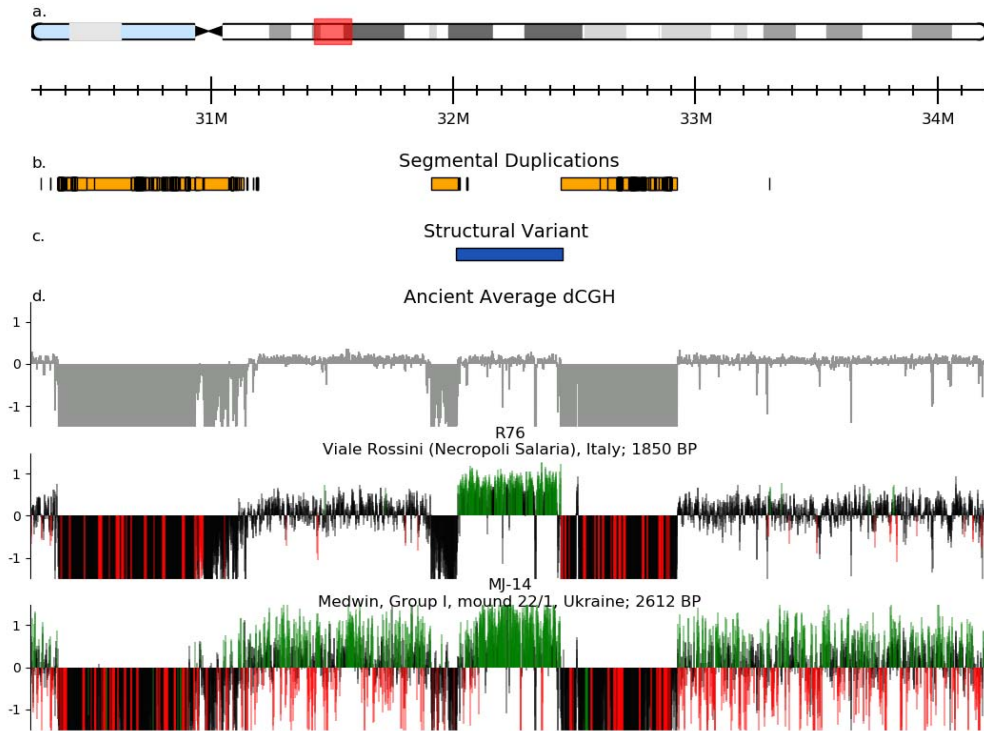
Figure S4i.6. Two microdeletions and nine microduplications are present at the 15q11.2 locus from the modern human dataset. There are perhaps two duplications present – a long (HGDP01103 and HGDP00239) and a short (the other six humans).



5604
 5605 **Figure S4i.7. Microdeletion of an ancient individual at the 15q11q13 (breakpoints BP3-**
 5606 **BP4) locus from the hg38 fastq dataset.** Although the BP3-BP4 microdeletion at
 5607 15q11q13 has not been formally associated with disease, carriers have been reported with
 5608 short stature, microcephaly, hypotonia, and facial dysmorphia ²⁵. None of the studied
 5609 modern individuals have the deletion.

5610
 5611
 5612
 5613
 5614
 5615
 5616
 5617
 5618
 5619
 5620
 5621
 5622
 5623
 5624
 5625
 5626
 5627
 5628
 5629
 5630
 5631

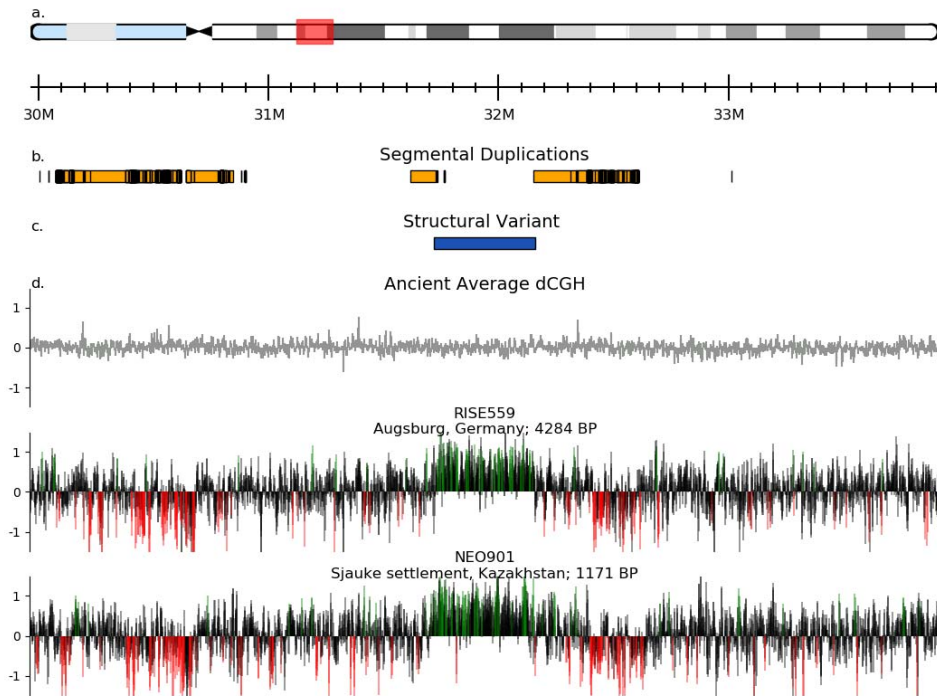
hg19 Ancient 15q13.3_smaller Pathogenic Variants



5632
5633
5634
5635
5636
5637
5638
5639
5640
5641
5642
5643
5644
5645
5646
5647
5648
5649
5650
5651
5652
5653
5654
5655
5656
5657

Figure S4i.8. Two microduplications are present at the 15q13.3 (*CHRNA7*) locus from the hg19 BAM dataset. The *CHRNA7* microdeletion has been associated with developmental delay and psychiatric disorders, but it is less clear which (if any) phenotypes are associated with the *CHRNA7* microduplication and there is evidence that *CHRNA7* duplication carriers have just as good cognitive function as others ²⁶.

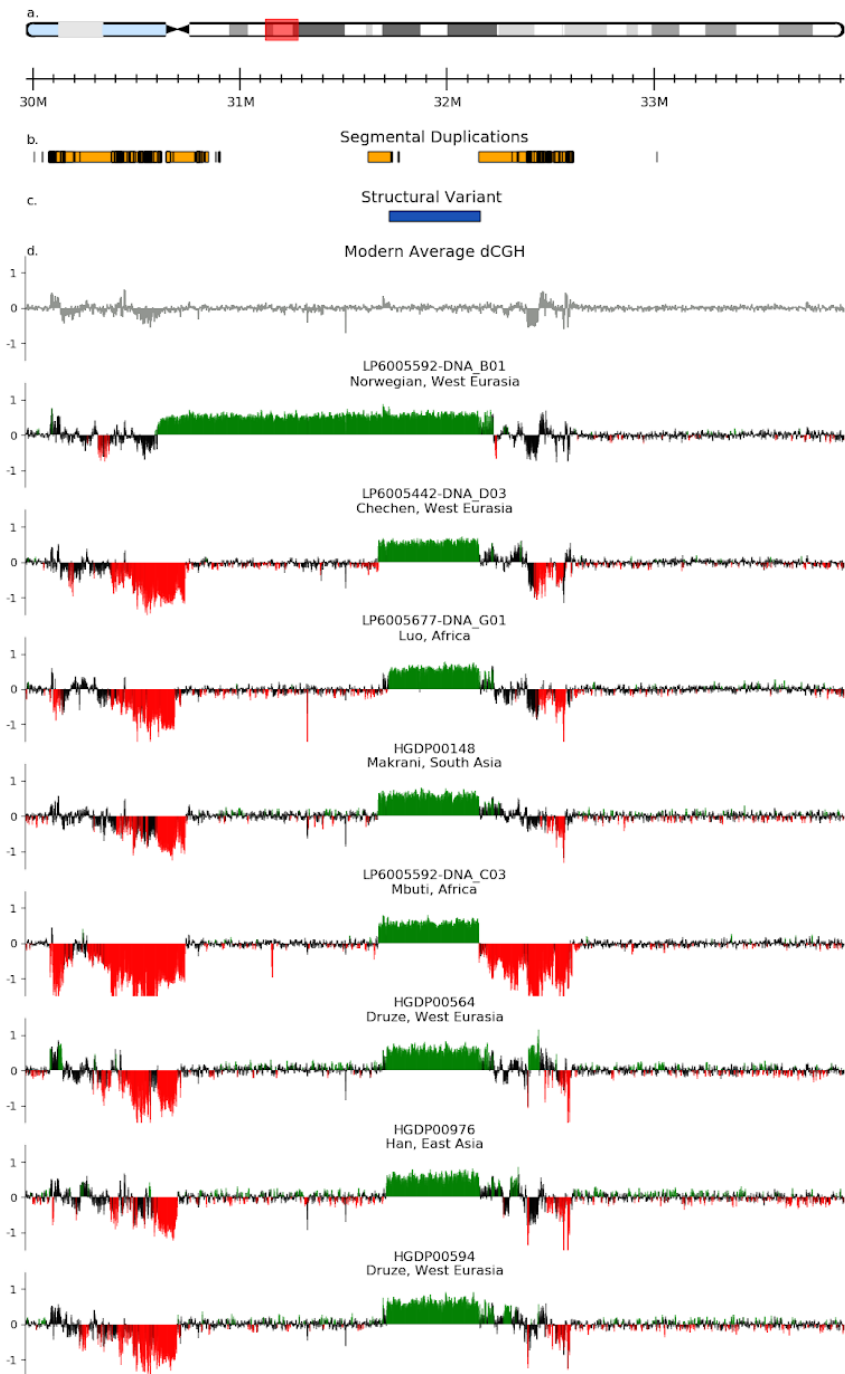
hg38 Ancient 15q13.3_smaller Pathogenic Variants



5658
5659
5660
5661
5662
5663
5664
5665
5666
5667
5668
5669
5670
5671
5672
5673
5674
5675
5676
5677
5678

Figure S4i.9. Two microduplications are present at the 15q13.3 (*CHRNA7*) locus from the hg38 fastq dataset.

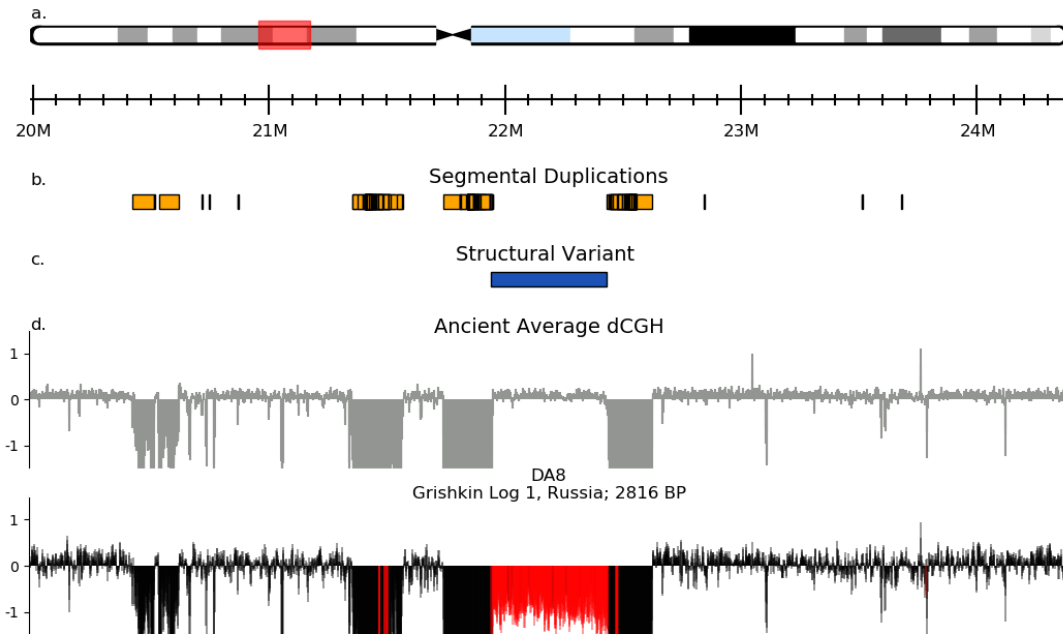
hg38 Modern 15q13.3_smaller Pathogenic Variants



5679
5680
5681
5682

Figure S4i.10. Two different microduplications are present at the 15q13.3 (*CHRNA7*) locus from the modern human dataset.

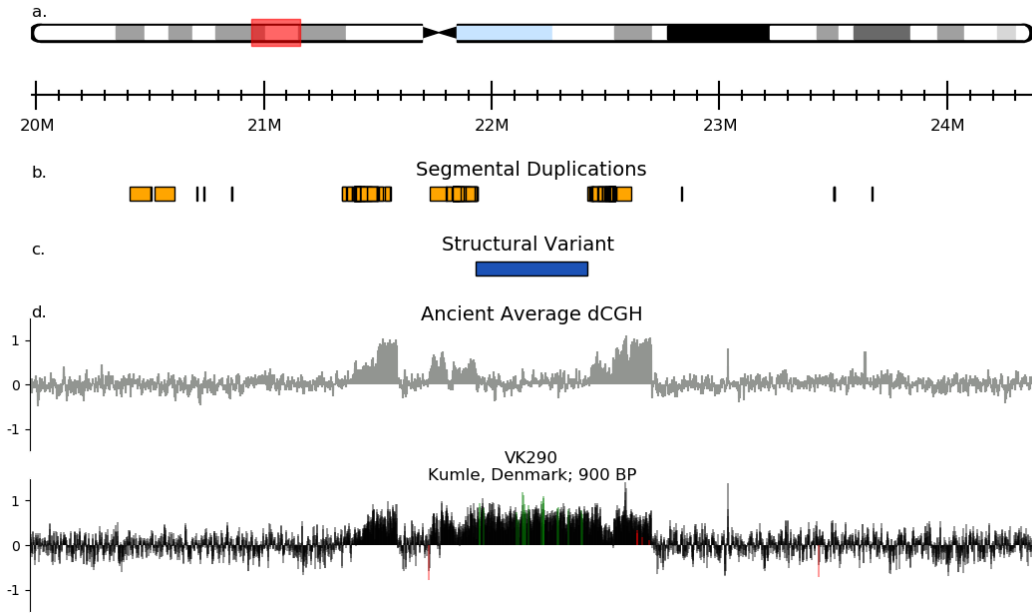
hg19 Ancient 16p12.1 Pathogenic Variants



5683
5684
5685
5686
5687
5688
5689
5690
5691
5692
5693
5694
5695
5696
5697
5698
5699
5700
5701
5702
5703
5704
5705
5706
5707
5708
5709
5710
5711

Figure S4i.11. One ancient individual from the hg19 BAM dataset has a deletion at the **16p12.1 locus**. At this locus, microdeletions and microduplications are associated with developmental delay, cognitive impairment, growth impairment, cardiac malformations, epilepsy, and psychiatric and behavioural problems ^{27,28}.

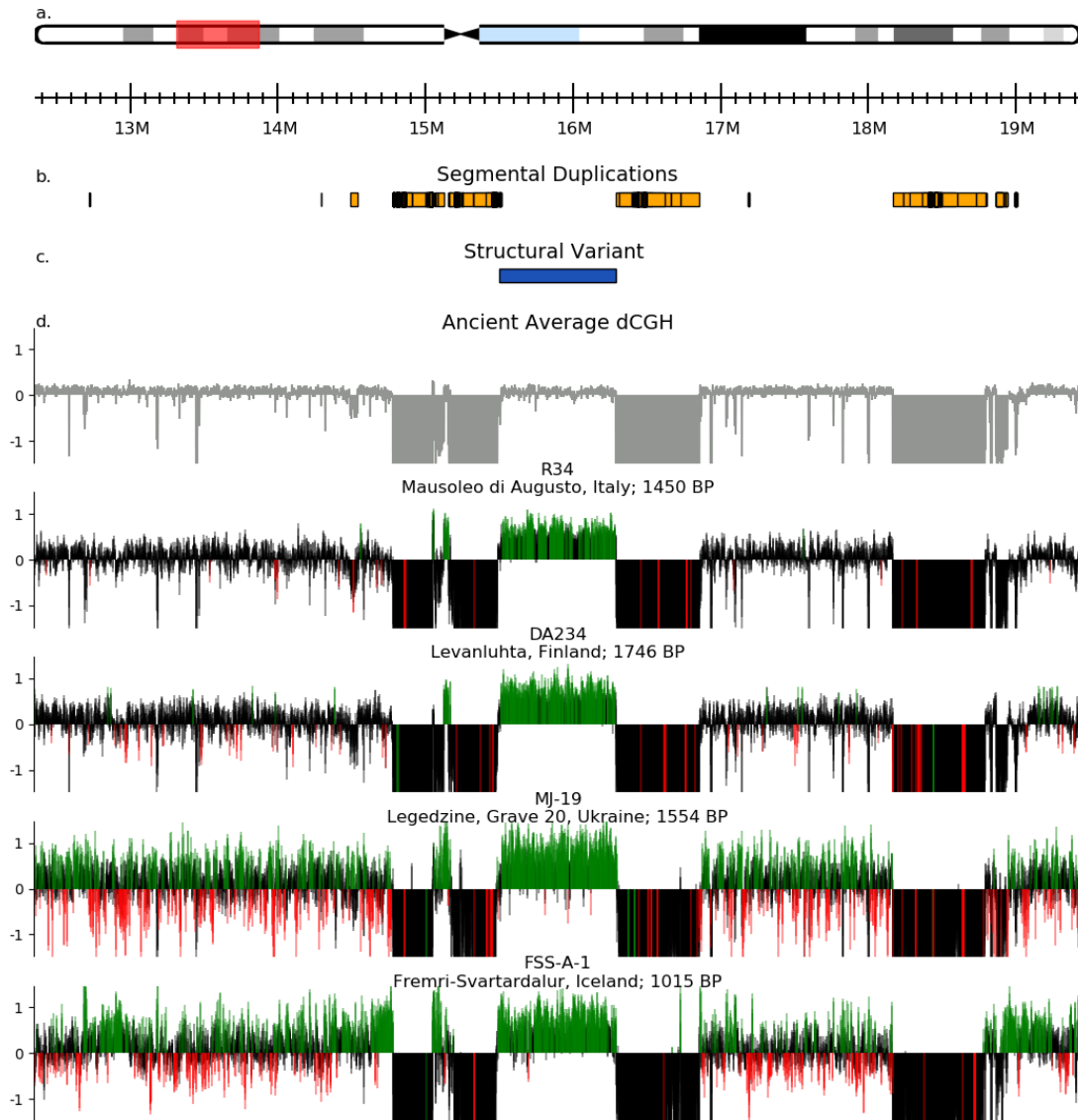
hg38 Ancient 16p12.1 Pathogenic Variants



5712
5713
5714

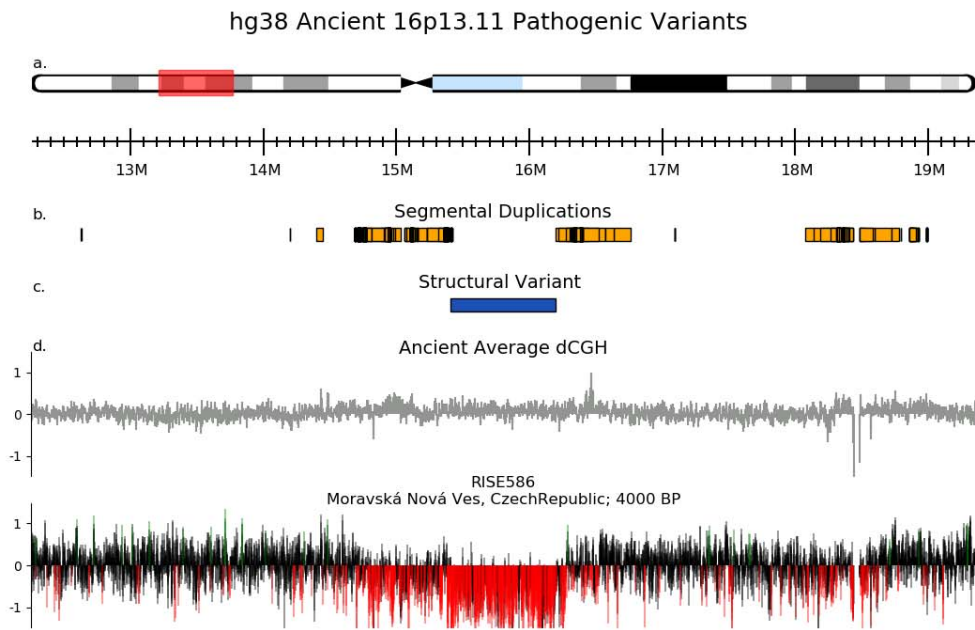
Figure S4i.12. One ancient individual from the hg38 fastq dataset has a duplication at the 16p12.1 locus.

hg19 Ancient 16p13.11 Pathogenic Variants



5715
5716
5717
5718
5719

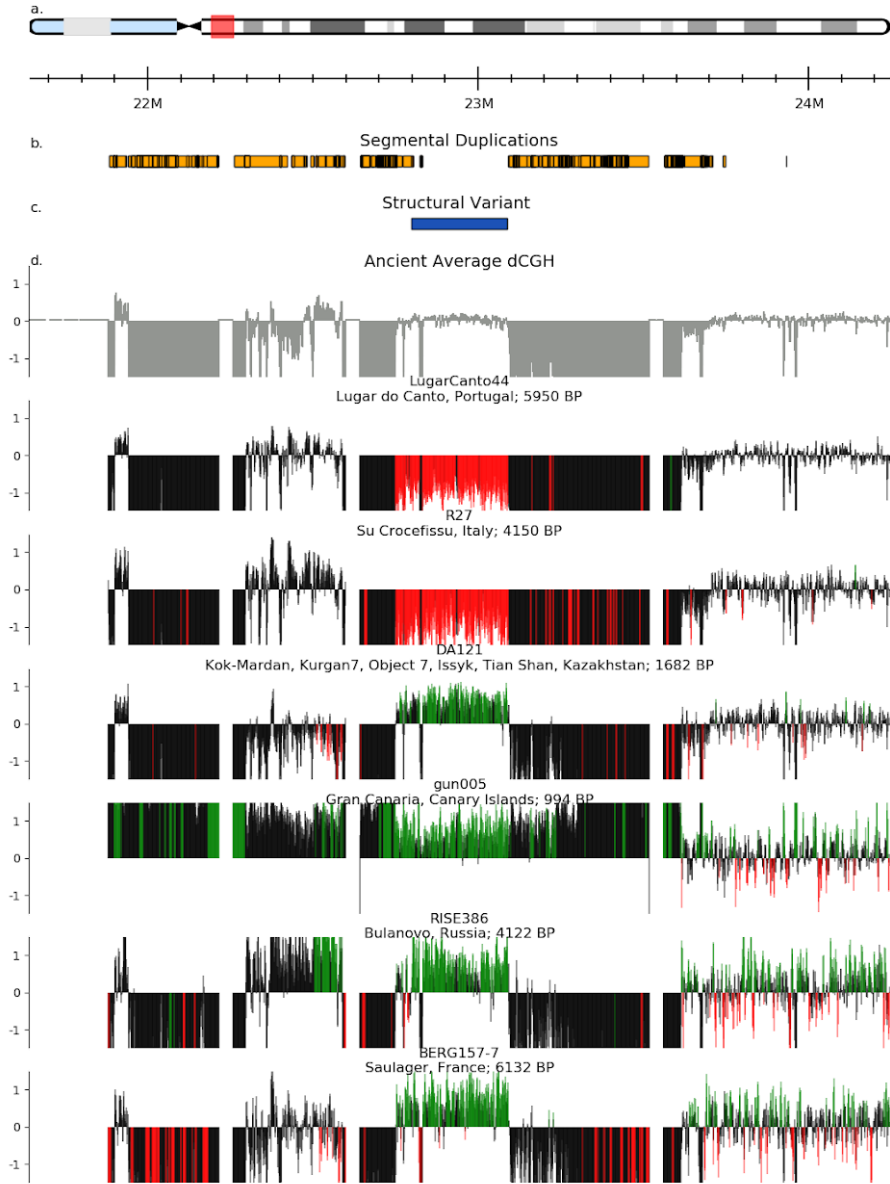
Figure S4i.13. Four microduplications are present at the 16q13.11 locus from the hg19 BAM dataset. At the 16p13.11 locus, microduplications and microdeletions have been associated with several phenotypes including behavioral abnormalities, developmental delay, congenital heart defects, and skeletal abnormalities^{29,30}.



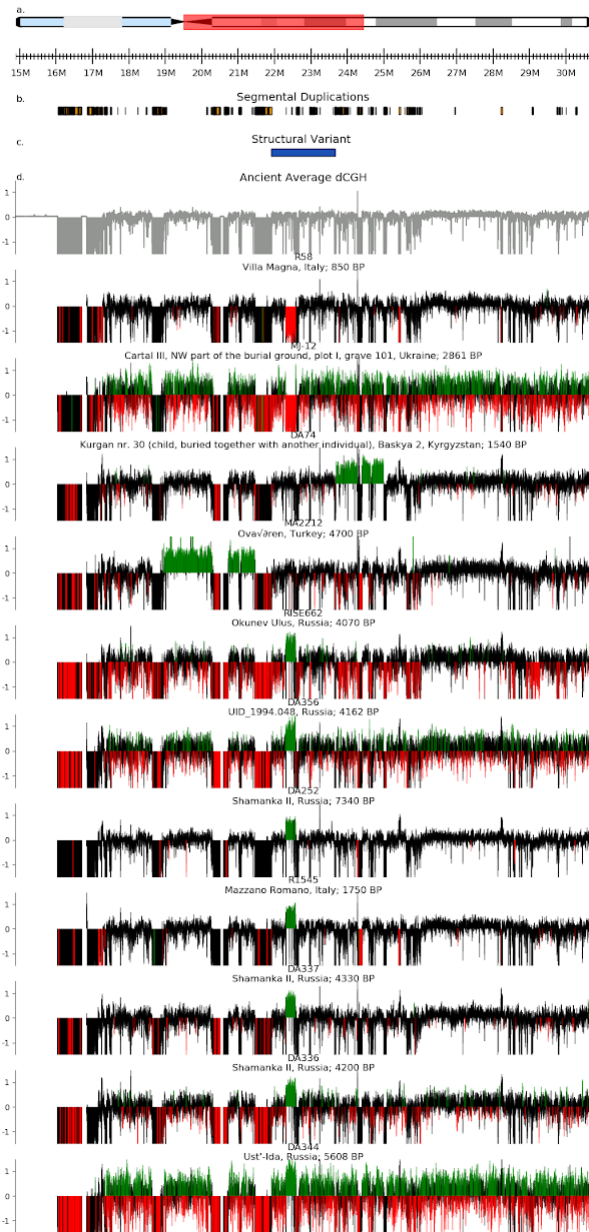
5720
5721
5722

Figure S4i.14. One microdeletion is present at the 16q13.11 locus from the hg38 fastq dataset.

hg19 Ancient 15q11.2 Pathogenic Variants



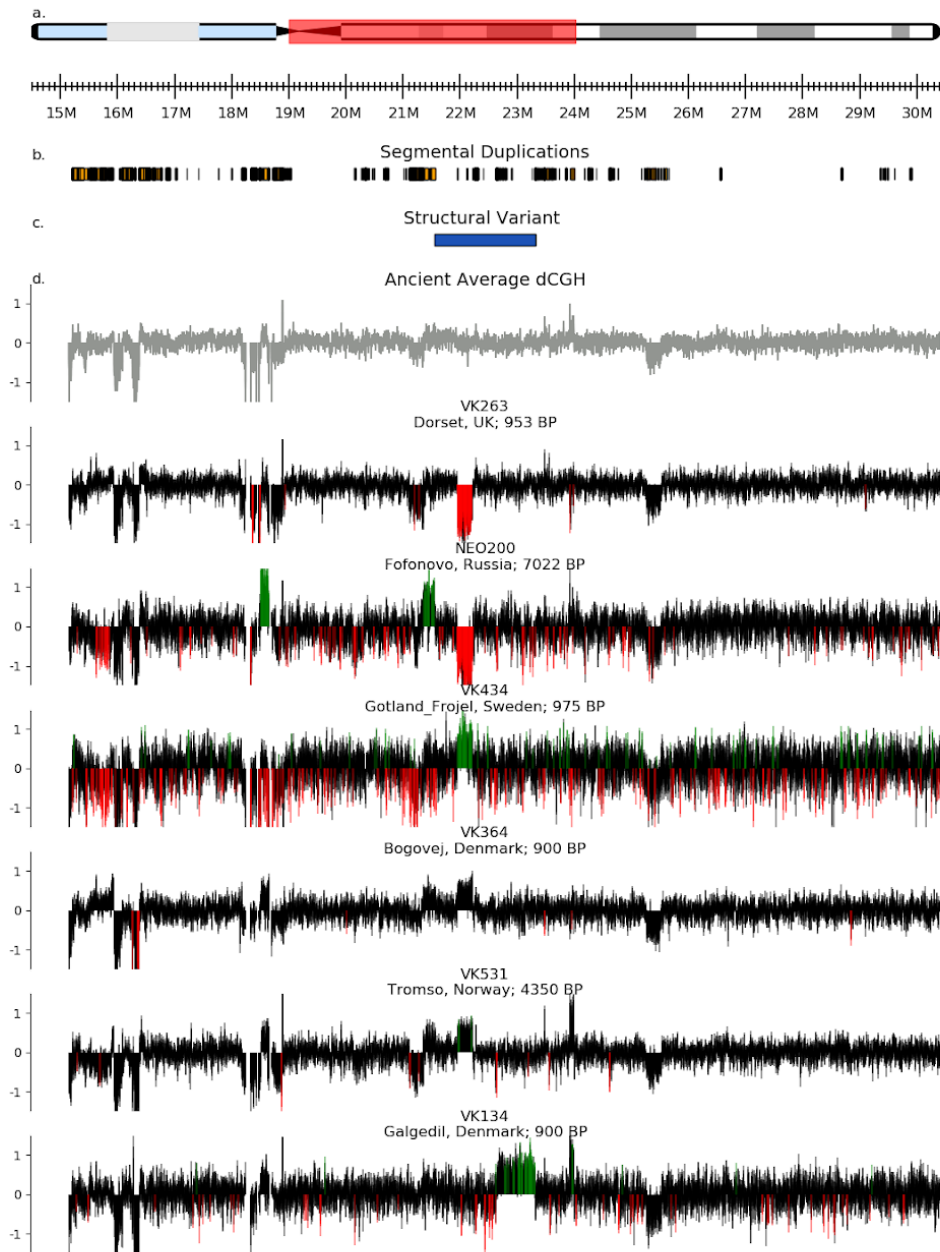
hg19 Ancient 22q11.2_distal Pathogenic Variants



5724
5725
5726
5727
5728
5729
5730
5731
5732
5733

Figure S4i.15. Two microdeletions and nine microduplications are present at the 22q11.2 (distal) locus from the hg19 BAM dataset. At the 22q11.2 locus, the *TOP3B* microdeletion has been reported to be associated with autism, learning disabilities, and dysmorphic features^{31,32}, and *TOP3B* knockout mice exhibit behaviour similar to psychiatric disorders and cognitive impairment³³, but risk estimates from large-scale population-based studies are lacking. All hg19 ancient individuals have this *TOP3B* structural variant with identical breakpoints as a 12-year-old patient with autism, cognitive impairment, and dysmorphic features³⁴ (Figure S4), while two other ancient individuals have other duplication breakpoints (DA74 and MA2212).

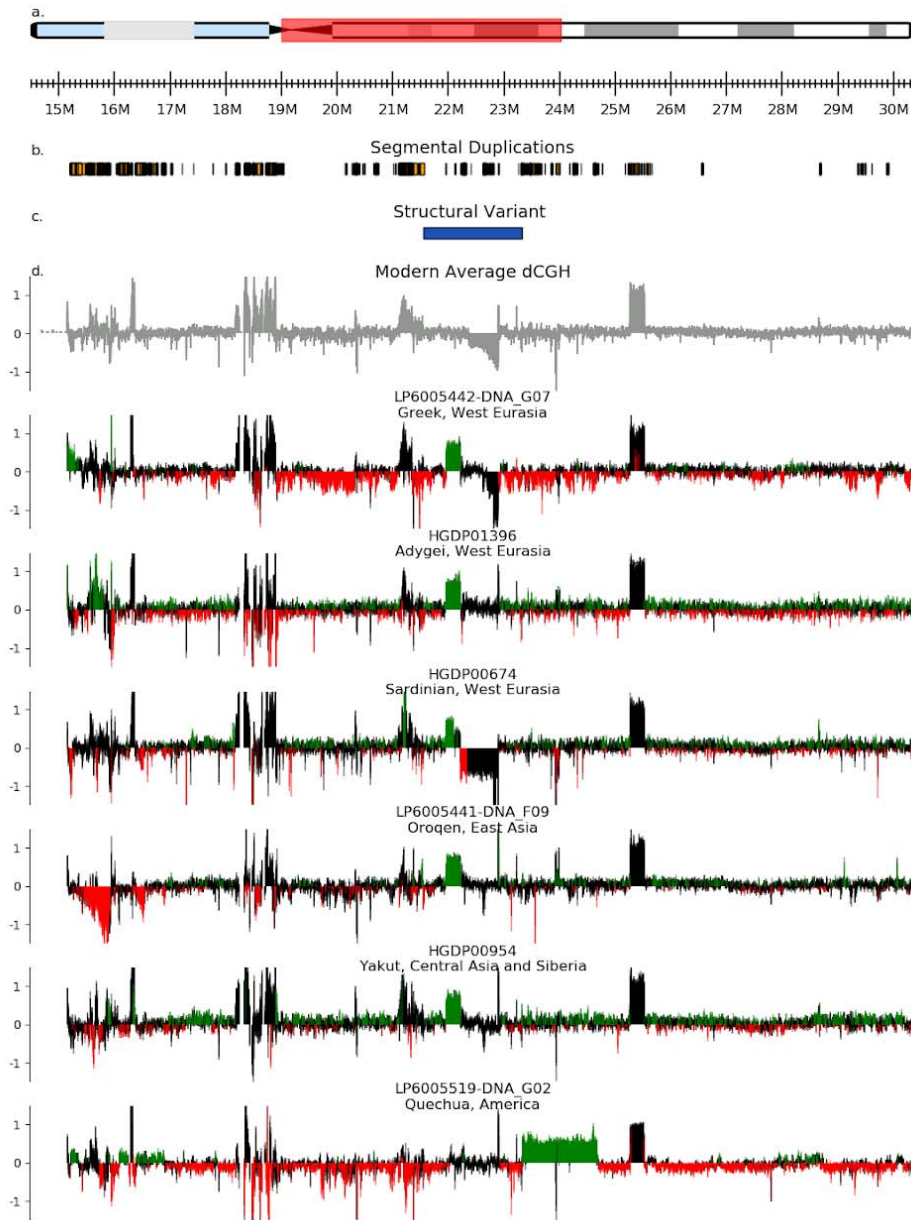
hg38 Ancient 22q11.2_distal Pathogenic Variants



5734
5735
5736
5737
5738

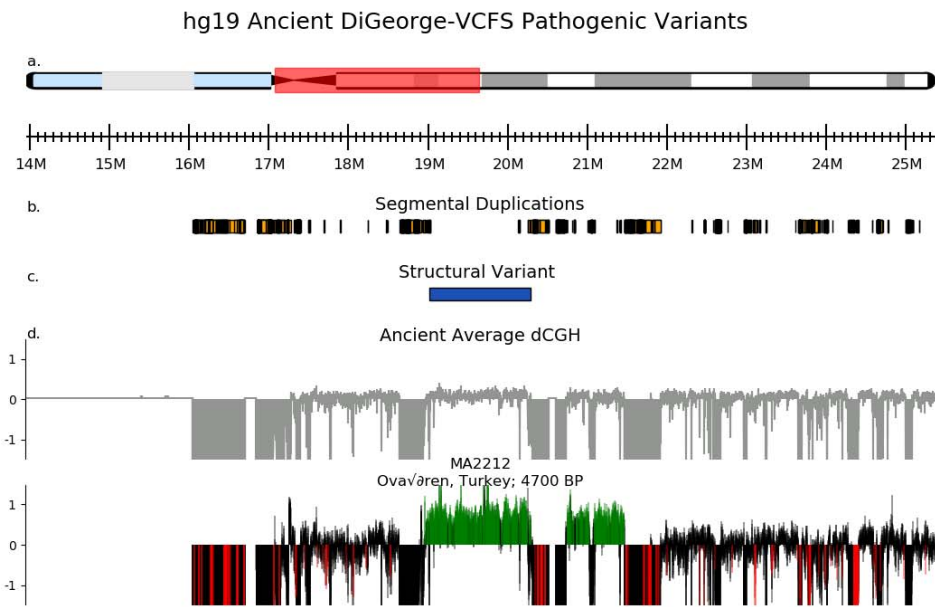
Figure S4i.16. Two microdeletions and four microduplications are present at the 22q11.2 (distal) locus from the hg19 BAM dataset.

hg38 Modern 22q11.2_distal Pathogenic Variants



5739
5740
5741
5742

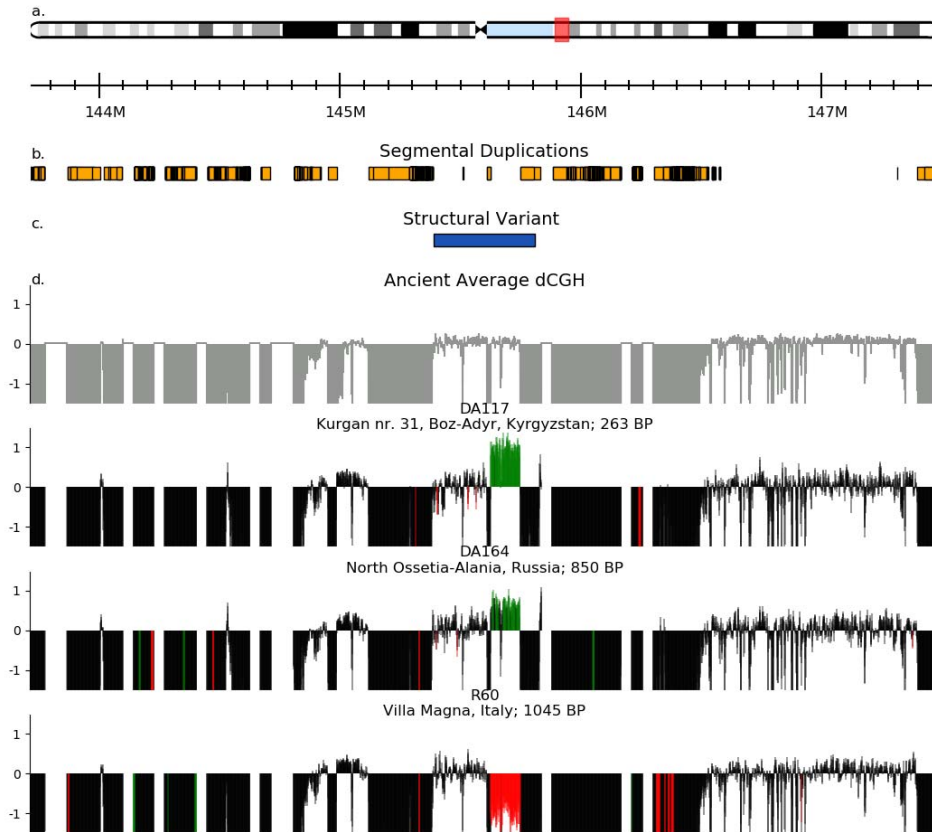
Figure S4i.17. Six microduplications (five of which directly overlap with *TOP3B*) are present at the 22q11.2 (distal) locus from the modern human dataset.



5743
 5744
 5745
 5746
 5747
 5748
 5749
 5750

Figure S4i.18. One ancient individual has a duplication at the DiGeorge-VCFS locus from the hg19 BAM dataset. Duplications at this locus are associated with schizophrenia as well as cardiovascular, parathyroid, thymic, and craniofacial abnormalities³⁵.

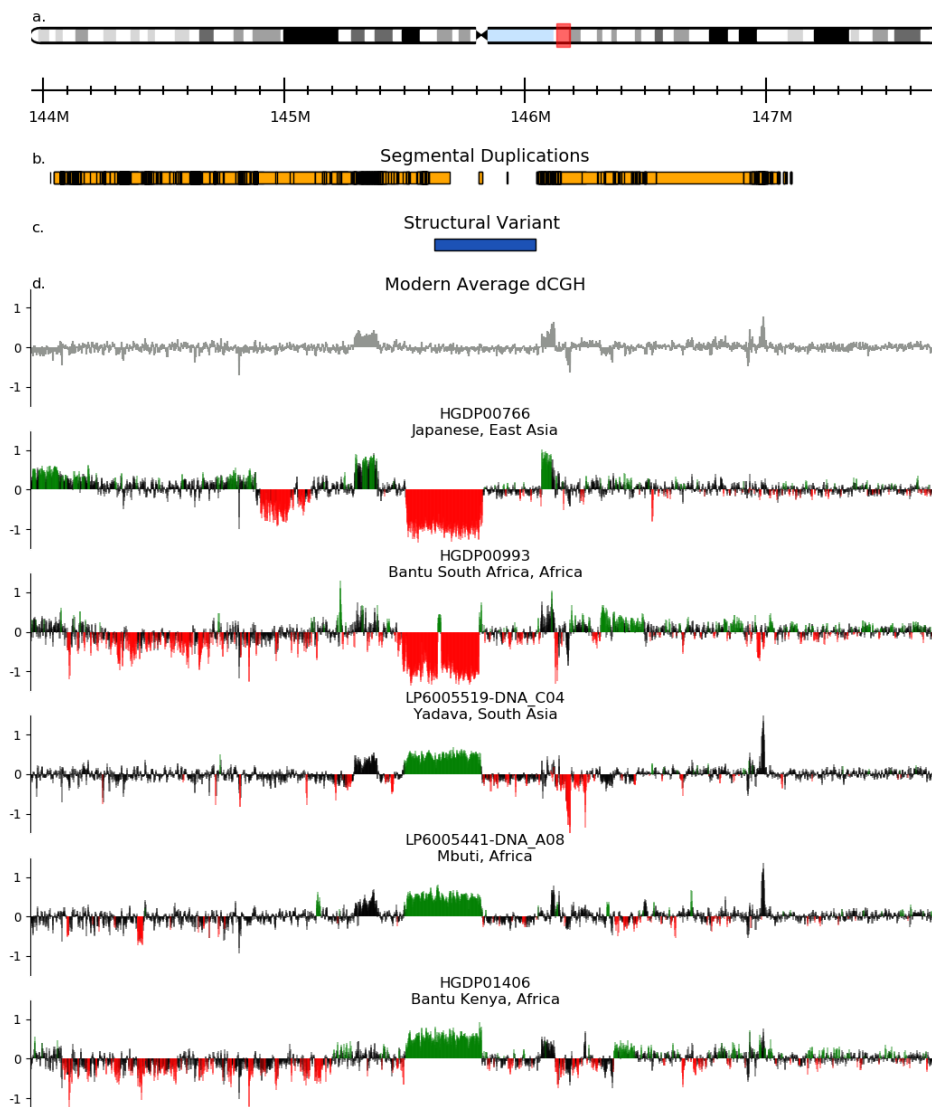
hg19 Ancient TAR Pathogenic Variants



5751
5752
5753
5754
5755
5756
5757

Figure S4i.19. Ancient TAR pathogenic variants. TAR (thrombocytopenia with absent radius) syndrome is a rare genetic disorder featuring the absence of the radius bone in the forearm³⁶; while no ancient individuals have a structural variant across the entire breakpoint of the syndrome, one ancient individual has a deletion across part of the locus and two ancient individuals have a duplication across the same breakpoints. The possible pathological implication of these CNVs is therefore unknown.

hg38 Modern TAR Pathogenic Variants



5758
5759
5760
5761

Figure S4i.20. Two deletions and three duplications are present in the modern human dataset, but with different breakpoints than the ancient human structural variants.

5762 References

- 5763 1. Girirajan, S., Campbell, C. D. & Eichler, E. E. Human copy number variation and
5764 complex genetic disease. *Annu. Rev. Genet.* **45**, 203–226 (2011).
- 5765 2. Weise, A. *et al.* Microdeletion and microduplication syndromes. *J. Histochem.*
5766 *Cytochem.* **60**, 346–358 (2012).
- 5767 3. Girirajan, S. *et al.* Phenotypic heterogeneity of genomic disorders and rare copy-

- 5768 number variants. *N. Engl. J. Med.* **367**, 1321–1331 (2012).
- 5769 4. Bergström, A. *et al.* Insights into human genetic variation and population history from
5770 929 diverse genomes. *Science* **367**, (2020).
- 5771 5. Mallick, S. *et al.* The Simons Genome Diversity Project: 300 genomes from 142
5772 diverse populations. *Nature* **538**, 201–206 (2016).
- 5773 6. Li, H. Burrows-Wheeler Aligner. <http://bio-bwa.sourceforge.net/>.
- 5774 7. Miles, A. *pysamstats: A fast Python and command-line utility for extracting simple*
5775 *statistics against genome positions based on sequence alignments from a SAM or BAM*
5776 *file.* (Github).
- 5777 8. Benson, G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic*
5778 *Acids Res.* **27**, 573–580 (1999).
- 5779 9. Sudmant, P. H. *et al.* Evolution and diversity of copy number variation in the great
5780 ape lineage. *Genome Res.* **23**, 1373–1382 (2013).
- 5781 10. Leffler, E. M. *et al.* Resistance to malaria through structural variation of red
5782 blood cell invasion receptors. *Science* **356**, (2017).
- 5783 11. Sudmant, P. H. *et al.* Diversity of human copy number variation and multicopy
5784 genes. *Science* **330**, 641–646 (2010).
- 5785 12. Crawford, K. *et al.* Medical consequences of pathogenic CNVs in adults:
5786 analysis of the UK Biobank. *J. Med. Genet.* **56**, 131–138 (2019).
- 5787 13. Olsen, L. *et al.* Prevalence of rearrangements in the 22q11. 2 region and
5788 population-based risk of neuropsychiatric and developmental disorders in a Danish
5789 population: a case-cohort study. *The Lancet Psychiatry* **5**, 573–580 (2018).
- 5790 14. Dolcetti, A. *et al.* 1q21.1 Microduplication expression in adults. *Genet. Med.*
5791 **15**, 282–289 (2013).
- 5792 15. Mefford, H. C. *et al.* Recurrent rearrangements of chromosome 1q21.1 and
5793 variable pediatric phenotypes. *N. Engl. J. Med.* **359**, 1685–1699 (2008).
- 5794 16. Brunetti-Pierri, N. *et al.* Recurrent reciprocal 1q21.1 deletions and
5795 duplications associated with microcephaly or macrocephaly and developmental and

- 5796 behavioral abnormalities. *Nat. Genet.* **40**, 1466–1471 (2008).
- 5797 17. Bailey, J. A. *et al.* Recent segmental duplications in the human genome.
5798 *Science* **297**, 1003–1007 (2002).
- 5799 18. Kent, W. J. *et al.* The human genome browser at UCSC. *Genome Res.* **12**,
5800 996–1006 (2002).
- 5801 19. Coyan, A. G. & Dyer, L. M. 3q29 microduplication syndrome: Clinical and
5802 molecular description of eleven new cases. *Eur. J. Med. Genet.* **63**, 104083 (2020).
- 5803 20. Glassford, M. R., Rosenfeld, J. A., Freedman, A. A., Zwick, M. E. & Mulle, J.
5804 G. Novel features of 3q29 deletion syndrome: Results from the 3q29 registry. *Am. J.*
5805 *Med. Genet. A* **170**, 999–1006 (2016).
- 5806 21. Stefansson, H. *et al.* CNVs conferring risk of autism or schizophrenia affect
5807 cognition in controls. *Nature* **505**, 361–366 (2014).
- 5808 22. Kirov, G. *et al.* The penetrance of copy number variations for schizophrenia
5809 and developmental delay. *Biol. Psychiatry* **75**, 378–385 (2014).
- 5810 23. Gudmundsson, O. O. *et al.* Attention-deficit hyperactivity disorder shares copy
5811 number variant risk with schizophrenia and autism spectrum disorder. *Transl.*
5812 *Psychiatry* **9**, 258 (2019).
- 5813 24. Jønych, A. E. *et al.* Estimating the effect size of the 15Q11.2 BP1–BP2
5814 deletion and its contribution to neurodevelopmental symptoms: recommendations for
5815 practice. *J. Med. Genet.* **56**, 701–710 (2019).
- 5816 25. Rosenfeld, J. A. *et al.* Deletions flanked by breakpoints 3 and 4 on 15q13 may
5817 contribute to abnormal phenotypes. *Eur. J. Hum. Genet.* **19**, 547–554 (2011).
- 5818 26. Kendall, K. M. *et al.* Cognitive performance and functional outcomes of
5819 carriers of pathogenic copy number variants: analysis of the UK Biobank. *Br. J.*
5820 *Psychiatry* **214**, 297–304 (2019).
- 5821 27. Girirajan, S., Pizzo, L., Moeschler, J. & Rosenfeld, J. *16p12.2 Recurrent*
5822 *Deletion*. (University of Washington, 1993).
- 5823 28. D’Angelo, D. *et al.* Defining the Effect of the 16p11.2 Duplication on

- 5824 Cognition, Behavior, and Medical Comorbidities. *JAMA Psychiatry* **73**, 20–30 (2016).
- 5825 29. Nagamani, S. C. S. *et al.* Phenotypic manifestations of copy number variation
5826 in chromosome 16p13.11. *Eur. J. Hum. Genet.* **19**, 280–286 (2011).
- 5827 30. Allach El Khattabi, L. *et al.* 16p13.11 microduplication in 45 new patients:
5828 refined clinical significance and genotype-phenotype correlations. *J. Med. Genet.* **57**,
5829 301–307 (2020).
- 5830 31. Stoll, G. *et al.* Deletion of TOP3 β , a component of FMRP-containing mRNPs,
5831 contributes to neurodevelopmental disorders. *Nat. Neurosci.* **16**, 1228–1237 (2013).
- 5832 32. Ahmad, M. *et al.* Topoisomerase 3 β is the major topoisomerase for mRNAs
5833 and linked to neurodevelopment and mental dysfunction. *Nucleic Acids Res.* **45**, 2704–
5834 2713 (2017).
- 5835 33. Joo, Y. *et al.* Topoisomerase 3 β knockout mice show transcriptional and
5836 behavioural impairments associated with neurogenesis and synaptic plasticity. *Nat.*
5837 *Commun.* **11**, 3143 (2020).
- 5838 34. Kaufman, C. S., Genovese, A. & Butler, M. G. Deletion of TOP3B Is
5839 Associated with Cognitive Impairment and Facial Dysmorphism. *Cytogenet. Genome*
5840 *Res.* **150**, 106–111 (2016).
- 5841 35. Agatsuma, S. & Hiroi, N. Chromosome 22q11 and schizophrenia. *Nihon*
5842 *Shinkei Seishin Yakurigaku Zasshi* **25**, 79–84 (2005).
- 5843 36. Boussion, S. *et al.* TAR syndrome: Clinical and molecular characterization of
5844 a cohort of 26 patients and description of novel noncoding variants of RBM8A. *Hum.*
5845 *Mutat.* **41**, 1220–1225 (2020).

5846
5847

5848
5849
5850
5851

5852

5853

5854