

# 1 **Patterns of population structure and genetic variation within the Saudi** 2 **Arabian population**

3 D.K. Malomane<sup>1\*</sup>, M. P. Williams<sup>2</sup>, C.D. Huber<sup>2</sup>, S. Mangu<sup>3,‡</sup>, M. Abedalthagafi<sup>4,5,‡</sup>, C. W. K.  
4 Chiang<sup>1,6,‡\*\*</sup>

5 <sup>1</sup>Center for Genetic Epidemiology, Department of Population and Public Health Sciences, Keck  
6 School of Medicine, University of Southern California, Los Angeles, CA.

7 <sup>2</sup>Department of Biology, Pennsylvania State University, University Park, PA.

8 <sup>3</sup>Titus Department of Clinical Pharmacy, Alfred E. Mann School of Pharmacy and Pharmaceutical  
9 Sciences, University of Southern California, Los Angeles, CA.

10 <sup>4</sup>Department of Pathology and Laboratory Medicine, Emory University Hospital, Atlanta, GA.

11 <sup>5</sup>Genomics Research Department, King Fahad Medical City, Riyadh, Saudi Arabia.

12 <sup>6</sup>Department of Quantitative and Computational Biology, University of Southern California, Los  
13 Angeles, CA.

14 <sup>‡</sup>These authors contributed equally to supervise this work.

15 \*Correspondence: [kholo642@gmail.com](mailto:kholo642@gmail.com).

16 \*\*Correspondence: [charleston.chiang@med.usc.edu](mailto:charleston.chiang@med.usc.edu).

17

## 18 **ABSTRACT**

19 The Arabian Peninsula is considered the initial site of historic human migration out of Africa.

20 The modern-day indigenous Arabians are believed to be the descendants who remained from

21 the ancient split of the migrants into Eurasia. Here, we investigated how the population history

22 and cultural practices such as endogamy have shaped the genetic variation of the Saudi

23 Arabians. We genotyped 3,352 individuals and identified twelve genetic sub-clusters that

24 corresponded to the geographical distribution of different tribal regions, differentiated by

25 distinct components of ancestry based on comparisons to modern and ancient DNA references.

26 These sub-clusters also showed variation across ranges of the genome covered in runs of

27 homozygosity, as well as differences in population size changes over time. Using 25,488,981

28 variants found in whole genome sequencing data (WGS) from 302 individuals, we found that the

29 Saudi tend to show proportionally more deleterious alleles than neutral alleles when compared

30 to Africans/African Americans from gnomAD (e.g. a 13% increase of deleterious alleles  
31 annotated by AlphaMissense between 0.5 - 5% frequency in Saudi, compared to 7% decrease of  
32 the benign alleles;  $P < 0.001$ ). Saudi sub-clusters with greater inbreeding and lower effective  
33 population sizes showed greater enrichment of deleterious alleles as well. Additionally, we  
34 found that approximately 10% of the variants discovered in our WGS data are not observed in  
35 gnomAD; these variants are also enriched with deleterious annotations. To accelerate studying  
36 the population-enriched deleterious alleles and their health consequences in this population,  
37 we made available the allele frequency estimates of 25,488,981 variants discovered in our  
38 samples. Taken together, our results suggest that Saudi's population history impacts its pattern  
39 of genetic variation with potential consequences to the population health. It further highlights  
40 the need to sequence diverse and unique populations so to provide a foundation on which to  
41 interpret medical- and pharmaco- genomic findings from these populations.

## 42 INTRODUCTION

43 Saudi Arabia is the largest country in the Arabian Peninsula (AP), the central hub of the world  
44 that connects Africa, Asia and Europe. The AP is considered one of the initial sites of historic  
45 human migration out of Africa (OOA), with presence of human footprints reported at least since  
46 50 – 60 thousand years ago (kya) and as early as 85 – 120 kya<sup>1-6</sup>. The contribution of the  
47 earliest expansion in present-day Arabians or other modern non-Africans has not been fully  
48 explored. However, genetic evidence suggests that all present-day Middle Eastern populations  
49 predominately descend from the same ancestral OOA population, as is the case for the other  
50 non-Africans<sup>6</sup>.

51 The genetic diversity of today's Arabians is shaped by a complexity of ancestries from historic  
52 and recent splits and admixture events. An early divergence of Arabian ancestors from other  
53 non-Africans is estimated to have happened shortly after the OOA event<sup>3</sup>. Among the non-  
54 Africans, Arabians carry a higher proportion of a deeply diverged 'ghost' ancestry, labeled 'Basal  
55 Eurasian,' and lower levels of Neanderthal admixture<sup>7,8</sup>. It has been hypothesized that the  
56 Arabians descended from the Basal Eurasians which diverged from other non-Africans before  
57 the major Neanderthal admixture<sup>7-10</sup>. Alternatively, the Basal Eurasians diverged from the non-  
58 Africans shortly after the OOA and was isolated until experiencing a later admixture in the  
59 Middle East around 25kya, which diluted the Neanderthal ancestry<sup>11</sup>. Since the OOA, Arabians  
60 have experienced series of admixtures, and the present-day Arabians have shared ancestries  
61 with various groups including Africans, South Asians, Europeans, Levantines, and Iranians<sup>6,12,13</sup>.

62 Despite the rich history of ancestries and being in the center of the world, for centuries the  
63 genetic pool of the Arab countries and the Greater Middle East (GME) have been greatly  
64 influenced and refined by mating practices. Arab countries have a high rate of endogamous and  
65 consanguineous marriages<sup>14,15</sup>, especially in Saudi Arabia with rates as high as 58%<sup>16,17</sup>. These  
66 endogamous marriages are meant to preserve family structure and strengthen bonds, as well as  
67 to ensure cultural, religious, financial and social stability<sup>17-19</sup>. Many of the consanguineous  
68 marriages are found among close relatives (e.g. 28.4% among first cousins<sup>16</sup>), but can also

69 extend to members of the same or related tribal groups. Endogamy leads to regional genetic  
70 isolation and population substructure. A recent study analyzing the population structure of  
71 Saudi Arabia based on less than a thousand indigenous genotyped samples showed a signature  
72 of tribal stratification within the population <sup>20</sup>. Furthermore, because many deleterious  
73 mutations are recessive-acting, consanguineous unions have the potential of increasing the  
74 burden of deleterious alleles in a population as these deleterious recessive alleles are co-  
75 inherited in offsprings <sup>14,21</sup>. This could increase the prevalence of genetic disorders, some of  
76 which have indeed been observed in Saudi Arabia <sup>22,23</sup>. While in the long run these deleterious  
77 recessive alleles are likely exposed to purifying selection due to increased homozygosity <sup>24,25</sup>,  
78 previous studies in the GME region have found no evidence of genetic purging of deleterious  
79 alleles due to the long-term practices of endogamy and consanguinity. Instead, intense  
80 inbreeding and/or reproductive compensation have been suggested to counteract the  
81 effectiveness of purifying selection in consanguineous populations <sup>26–29</sup>. Moreover, with small  
82 effective population sizes and inbreeding, variants acting additively tend to accumulate at a  
83 much higher rate and negative selection is less effective in removing weakly deleterious alleles  
84 <sup>30–33</sup>. Overall, small effective population size and intense inbreeding through consanguinity may  
85 result in an abundance of deleterious alleles due to its negative impact on the effectiveness of  
86 negative selection.

87 In the present study, we genotyped 3,352 individuals with high-density SNP array and whole  
88 genome sequenced (WGS) 302 individuals to investigate how the population history and  
89 cultural practices have shaped the genetic structure of the Saudi population. We investigated  
90 the pattern of admixture in Saudi sub-populations through the lens of both modern and  
91 available ancient DNA samples and inferred the population size trajectories over time. Finally,  
92 we leveraged the 302 whole genome sequenced individuals to further explore the impact of the  
93 population history on the distribution of genetic variation within social structure and potential  
94 consequences to today's population health.

95

## 96 RESULTS

### 97 Genetic substructure and admixture patterns of Saudi Arabians

98 We merged 3,352 genotyped Saudi individuals after quality control (see **Methods**) with 302  
99 whole genome sequencing (WGS) samples, based on 603,833 shared segregating sites, to  
100 explore the population structure. We performed principal component analysis (PCA) on the  
101 combined set and projected the first 10 principal components (PCs; **Figure S1**) down to 2  
102 dimensions using Uniform Manifold Approximation and Projection (UMAP). Average Silhouette  
103 Width (ASW) clustering on the UMAP results suggested that twelve genetic sub-clusters within  
104 the Saudi population best fit the data, although visually 6 to 8 sub-clusters may also be sensible  
105 (**Methods; Figure 1A**). The distribution of individuals by ASW clusters are presented in **Figure**  
106 **S2A**. To aid in the geographical interpretation of these sub-populations, we intersected the  
107 clustering results with self-reported or predicted tribal geographic labels from the cohort  
108 (**Methods**). Due to privacy protection and ethical restrictions, we did not have access to specific  
109 tribal name of each individual but rather the geographic regions of the tribes. We found that the  
110 12 clusters corresponded to geographical structure of the tribes within Saudi Arabia, with each  
111 cluster generally consisting of a majority of its members from a single geographical region  
112 (Central, West, North, South, or East) whether using harmonized tribal labels or self-reported  
113 labels when available, except for clusters 11 and 12 (**Table S1, Figure 1A** and **Figure S2B**). Both  
114 clusters 11 and 12 had multiple dominating tribal regions. We note that there were multiple  
115 separate genetic clusters affiliated to the same geographic regions (e.g. clusters 2, 3, and 9 from  
116 Central region; 4, 5, 7, and 10 from the Western region, etc.). This observation is unlikely due to  
117 errors in inferring tribal regional labels, since previous studies using completely self-identified  
118 indigenous tribal information also showed limited inter-tribal marriages within a region <sup>20</sup>.  
119 Among the clusters we inferred, cluster5 from the Western region appeared to be most  
120 differentiated from the rest of the cohort, in both UMAP (**Figure 1A**) and PCA (PCs 6 and 7;  
121 **Figure S1**).

122 For a global comparison, we compared the Saudi clusters to the populations from the Human  
123 Genome Diversity Panel (HGDP)<sup>34</sup>. Consistent with previous reports<sup>20,35</sup>, the Saudi individuals  
124 clustered between Africans, Central & South Asians and Europeans and were the most distant to  
125 East Asians (**Figure 1B**). Cluster12 showed the strongest affinity towards the African reference  
126 individuals, followed by cluster3, while cluster11 showed affinity towards Europeans, Africans  
127 and Central & South Asians. The remaining clusters co-localized mainly with the Middle Eastern  
128 reference individuals from the HGDP panel.

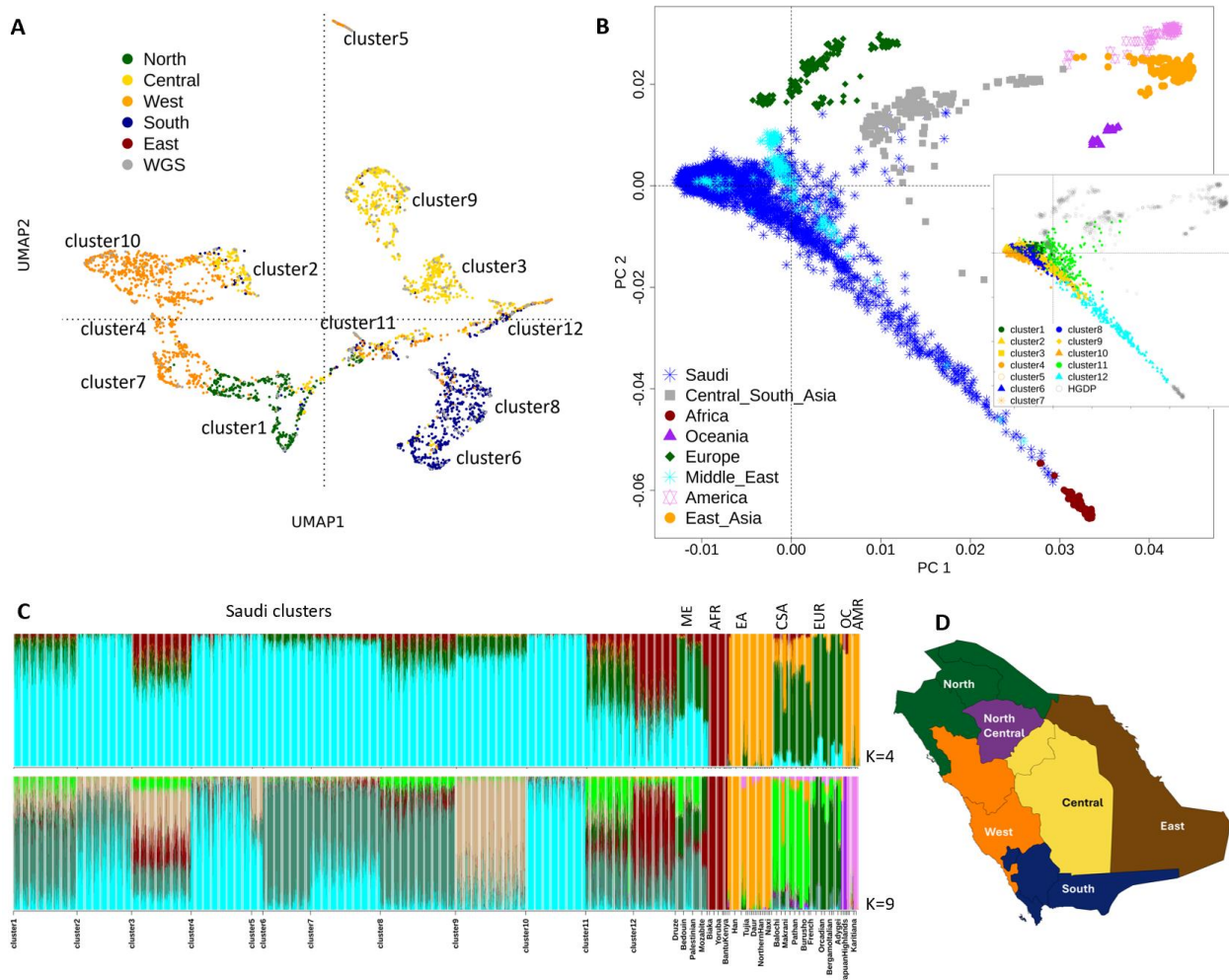
129 Our observation of population structure from PCA is also supported by unsupervised  
130 ADMIXTURE analysis combining Saudi with HGDP populations. For instance, at K = 4, clusters 11  
131 and 12 also exhibited the highest levels of admixture (**Figure 1C** and **Table S2**). We labeled  
132 ancestry components by the HGDP population with dominating or highest admixture  
133 proportions (**Table S2**). The most dominating ancestry in the Saudi clusters was one largely  
134 shared with the HGDP Middle Eastern populations (Druze, Bedouin, Mozabite, and Palestinian  
135 (**Figure 1C**, top), which we termed Middle Eastern-like (ME-like; cyan) ancestry, with Bedouin  
136 showing the largest amount of this ancestry among the HGDP Middle Eastern individuals (**Table**  
137 **S2**). On the other hand, clusters 12, 11, and 3 had on average less than two-thirds of this ME-  
138 like ancestry component and were enriched with African (AFR)- (red) and/or European (EUR)-  
139 like (green) ancestries.

140 One of the lowest cross validation errors occurred at K = 9 (**Figure S3C**), which also introduced a  
141 new ancestry component distinguishing the CSA-like ancestry from the EUR-like ancestry (For K  
142 = 5 to 8 (**Figure S3A**) the ADMIXTURE algorithm were mostly distinguishing ancestries within  
143 Saudi Arabian clusters themselves). At K = 9 (**Figure 1C** bottom, and **Table S3**), the relationship  
144 between cluster11 and the Central & South Asians (CSA) that we observed on PCA can also be  
145 observed, where cluster11 carried more (average proportion = 0.223) of such CSA-like ancestry  
146 compared to other clusters (average proportions less than 0.1). We also observed three ME-like  
147 ancestries (**Table S3** and **Figure 1C**, bottom). One of the Middle Eastern-like ancestry (ME-2)  
148 that is dominant in several (> 10% in 10 out of the 12) clusters, particularly in clusters 6, 8, 1,  
149 and 7, is also found in Sardinians and other Italians & Adygei but completely missing in the

150 Russians. The other Middle Eastern-like ancestries (ME-1 and ME-3) are found distributed in a  
 151 subset of the clusters. In particular, the ME-3 ancestry is found in high proportions in clusters  
 152 10, 4, 5, and 2 (average proportions = 0.45 - 0.92), and was also found in HGDP-Bedouin (but  
 153 absent from HGDP-Druze, HGDP-Mozabite and HGDP-Palestinian). This ancestry appears to be  
 154 enriched in Qataris Bedouins and Saudi Arabians but not other Middle Eastern populations, and  
 155 was suggested to reflect an indigenous Arab ancestry<sup>3</sup>.

156

157



158 **Figure 1. The genetic structure of Saudi Arabians and its relation to global populations.** (A) A  
 159 two-dimensional UMAP of Saudi Arabians based on the top 10 principal components. Each  
 160 individual is colored based on the affiliated tribal region (see (D)). WGS samples did not have  
 161 self-reported or harmonized tribal affiliation and are assigned their own color. (B) PCA of Saudi  
 162 Arabian clusters and HGDP populations. Saudi Arabians are grouped in a single group. Inset

164 shows clusters 1 - 10 colored according to the most prevalent tribal region represented in the  
165 cluster (see **Table S1**). Because clusters 11-12 has no single dominating tribal region, they were  
166 assigned distinct separate colors. (C) Admixture analysis of Saudi Arabian clusters and HGDP  
167 populations for  $K = 4$  (top) and  $K = 9$  (bottom). ME – Middle Eastern, AFR – African, EA – East  
168 Asian, CSA – Central & South Asian, EUR – European, OC – Oceania, AMR – American. The  
169 names of Saudi clusters and HGDP populations are shown on the bottom X-axis. However, due  
170 to limited space some of the labels for smaller populations from HGDP are omitted. Grouped  
171 regional labels are shown on the top X-axis of plots. We show the admixture results of the Saudi  
172 clusters alone in **Figure S3B**. (D) A regional map of Saudi Arabia with matching colors to the  
173 regional labels in (A) and (B).

174

175 The evidence for admixture for clusters 12, 11, and 3, together with clusters 8 and 1 were also  
176 corroborated by admixture  $f_3$ -statistics. Using all possible pairs of HGDP populations as potential  
177 proxies of ancestral sources (possibly through shared ancestry), only these five Saudi sub-  
178 clusters showed any significantly negative  $f_3$ -statistics indicative of admixture (**Tables S4 - 8**). We  
179 further investigated the degree of shared drift between the Saudi clusters and the HGDP  
180 populations using the outgroup  $f_3$ -statistics. We used HGDP-Han population as the outgroup  
181 over the typical choices of the San, Mbuti, or Yoruba populations since we expect African being  
182 a plausible admixing source in Saudi and indeed found positive gene flow between the HGDP  
183 African populations and the Saudi clusters (**Methods; Figure 1C**). We found the pattern of the  
184 outgroup  $f_3$ -statistics to be similar among Saudi clusters 1 - 10, in contrast to clusters 11 and 12  
185 (**Figure S4**). Clusters 1 - 10 showed highest shared drift to other Saudi clusters followed by the  
186 Middle Eastern and European populations (**Figure S4**). On the other hand, cluster12 showed  
187 most shared drift to HGDP-African populations, even more so than it is to other Saudi clusters,  
188 further corroborating the results in **Figure 1B** and **C**. In relation to other Middle Eastern  
189 populations from HGDP, clusters 1 - 10 showed greater shared drift to Bedouin while cluster12  
190 was most related to the Mozabite. Given the high African ancestry in cluster12, this relationship  
191 with Mozabite population from North Africa is not surprising, and could have either arose from  
192 shared sub-Saharan ancestry or an Arabic admixture event during the Islamic expansion into  
193 North Africa about 1,200 – 1,400 kya<sup>36</sup>, or a combination of both. Interestingly, in relation to  
194 the European populations, all Saudi clusters shared the greatest drift with the Sardinians, and



195 the least with the Russians. The relationship to the Sardinians is in concordance with Charati et  
196 al.,<sup>37</sup>. Sardinians heavily harbor early farmer Neolithic ancestry, which expanded into Europe  
197 from the Near East and Anatolia<sup>38–40</sup>. Saudi Arabians are estimated to have split from Sardinians  
198 around 20 kya<sup>6</sup>. Their relationship might be reflecting the ancient Neolithic ancestry, a product  
199 of the Arabian migration into Italian islands or a result of continuous and recent admixtures  
200 <sup>37,40–42</sup>.

201

## 202 **Genetic legacy of ancient ancestry in modern day Saudi Arabians**

203 We expanded our understanding of Saudi genetic history by integrating 302 WGS Saudi  
204 individuals with ancient DNA (aDNA) datasets (1240k AADR v54.1<sup>43</sup> and ancient Bahrain  
205 individuals<sup>12</sup>; **Figure S5; Methods**). Our integrated analysis with aDNA data corroborated many  
206 of our findings above on population structure and admixture history using only modern DNA.  
207 Previous studies have shown that Eastern Arabian Peninsula (AP) populations have higher  
208 ancient Zagros mountains/Caucasus mountains hunter-gatherer (CHG)-related ancestry than  
209 Western AP populations<sup>7,11,44</sup>. We observed a similar geographical divide among Saudi  
210 Arabians. Using an Epipaleolithic sample from the Natufian culture (Natufian EpiP) to represent  
211 Levantine ancestries and a sample from the Ganj Dareh Neolithic settlement in western Iran  
212 (Ganj Dareh N) along with other CHG to represent Zagros ancestries, we found that Saudi  
213 clusters mainly from the West and South region of Saudi Arabia (clusters 6, 10, 5, 8, 7, and 4)  
214 showed significant ( $f_4$  Z-score < -3) excess shared drift with Levantine ancestries relative to  
215 Zagros ancestries (**Figure S6**, bottom, and **Table S9**). We also evaluated the spatiotemporal  
216 distribution of African ancestry amongst the Saudi clusters using the outgroup  $f_3$ -statistic of  
217 form  $f_3(\text{Han.DG}; \text{ancient and present-day African population, Saudi cluster})$  along with 108  
218 African populations spanning from the present-day to ~15 kya. Temporal analysis across four  
219 time bins ([0,0], (0, 1,000 kya], (1,000 kya, 4,000 kya], and (4,000 kya, 15,500 kya]) revealed  
220 that Saudi clusters 12 and 3 consistently exhibited elevated African ancestry (**Figure S7 and**  
221 **Table S10**), consistent with **Figures 1B, 1C, and S4**. Importantly, this increased genetic affinity

222 appears associated with African populations south of approximately 20°N latitude, whereas  
223 genetic affinity to Northern African ancestry remains relatively uniform across all time bins and  
224 Saudi clusters. This finding is also corroborated by  $f_4$ -statistic confirming that Saudi cluster12  
225 shows the highest level of African ancestry (**Figure S6**, top and **Table S11**).

226 Previous research has identified genetic continuums amongst present-day AP populations with  
227 respect to ‘Basal Eurasian’ ancestry—a hypothesized ghost lineage that diverged from the  
228 primary out-of-Africa lineage prior to Neanderthal introgression<sup>10,45</sup>. We used the  $f_4$ -statistic of  
229 form  $f_4(\text{Saudi cluster, Han.DG; Ust-Ishim, ancient African population})$  to estimate the relative  
230 amount of Basal Eurasian ancestry (i.e. the drift basal to the shared drift between Han.DG and  
231 Ust-Ishim, a 45 ky sample from western Siberia) across the Saudi cluster cohort. Strongly  
232 negative  $f_4$  values of this form is consistent with elevated shifted drift with Basal Eurasian  
233 ancestry. However, recent African admixture could confound this signal, and different African  
234 aDNA samples are likely to represent Basal Eurasian ancestry to different degrees, thereby  
235 prompting us to test a range of ancient African populations. We found that most Saudi clusters,  
236 regardless of geographical locations, showed similar levels of Basal Eurasian ancestry (**Figure 2B**  
237 and **Table S12**). We found that Saudi cluster12, and to a lesser extent cluster3, indeed showed  
238 stronger negative  $f_4$  values when the African aDNA source is from Central, Southern, or Eastern  
239 parts of Africa, consistent again with their recent African admixture from these regions (**Figure**  
240 **2B** and **Table S12**). The estimates of Basal Eurasian ancestries across Saudi clusters are much  
241 more comparable to each other when using the North African references (**Figure 2B** and **Table**  
242 **S12**), consistent with the recent proposal that the Epipaleolithic Iberomaurusian Taforalt  
243 population is the best proxy for Basal Eurasian ancestry<sup>7</sup>. Moreover, these results confirm that  
244 the North African aDNA references are not closely related to the admixing African source in  
245 Saudi clusters 12, and 3.

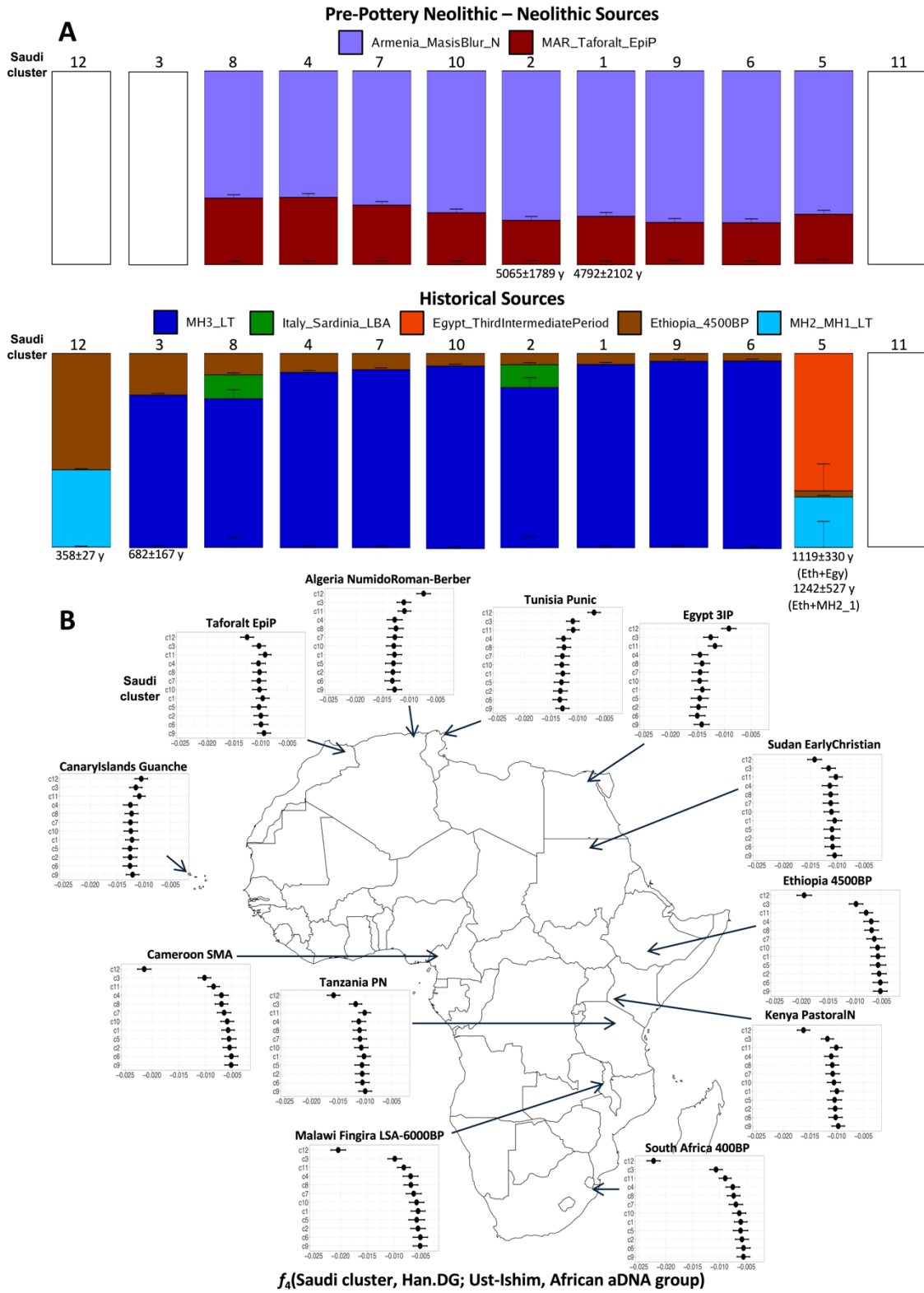
246 To investigate the recourse of Saudi cluster ancestries through time, we implemented qpAdm  
247 modeling with chronologically stratified sources across four temporal bins: Paleolithic-Neolithic  
248 (P-N), Chalcolithic (C), Bronze Age (BA), and Historical (H), maintaining a fixed set of 11 right  
249 groups throughout (**Methods** and **Tables S13 - S16**). Saudi clusters were most successfully

250 modeled using sources from the temporal bookends, the earliest (P-N) and most recent (H)  
251 periods (**Figure 2A**). The poor performance of Chalcolithic and Bronze Age sources in modeling  
252 Saudi cluster ancestry (**Table S14 - S15**) likely reflects the current absence of well-represented  
253 sources, such as an absence of ancient Arabian Peninsula genetic data from these intermediate  
254 periods. In the P-N period a two-source model combining similar proportions of North African  
255 (Taforalt EpiP) and Neolithic Armenian (MasisBlur N) components appeared to be the most  
256 plausible model for most of the Saudi clusters, corroborating  $f_4$ -statistic evidence of  
257 approximately equal ancient North African genetic affinity (**Figure 2B**). Clusters 3, 11, and 12  
258 were rejected in this relatively simple model, reflecting their complicated admixture history as  
259 described above. Re-analysis of clusters 3, 11, and 12 by removing Mbuti.DG from the qpAdm  
260 construct (owing to their greater African-related ancestry) did not result in a better fitted model  
261 either (data not shown). We used DATES<sup>46</sup> to estimate admixture timing and only Saudi  
262 clusters 1 ( $4792 \pm 2102$  years) and 4 ( $5065 \pm 1789$  years) returned well-fitting admixture timing  
263 estimates (normalized root mean standard deviation (NRMSD)  $< 0.7$ , Z-score  $> 2$  and **Table S17**).

264 Analysis of Saudi cluster ancestry through Historical sources revealed consistent ancestral  
265 contributions from two Bahraini groups (MH3 LT and MH1-MH2 LT; MH3 LT has greater  
266 Levantine-related ancestry<sup>12</sup>). The MH3 LT + Ethiopia 4500BP model was plausible for seven  
267 Saudi clusters: clusters 1, 3, 4, 6, 7, and 9 with only Saudi cluster3 returning a well-fitting  
268 admixture date estimate ( $682 \pm 167$  years). For Saudi clusters 2 and 8, the inclusion of Sardinia  
269 LBA as a third source on top of Ethiopia 4500BP + MH3 LT appeared to improve the model.  
270 However, admixture timing analysis revealed high-uncertainty estimates for both two-source  
271 and three-source qpAdm models, whereby MH3 LT + Sardinia LBA was well-fitting for only  
272 cluster8 ( $612 \pm 835$  years) and the Ethiopia 4500BP + Sardinia LBA model well-fitting for Saudi  
273 clusters 2 ( $474 \pm 625$  years) and 8 ( $820 \pm 952$  years). Interestingly the Arabian ancestry  
274 component for both Saudi clusters 12 and 5 is best modeled by the MH1-MH2 LT group,  
275 characterized by reduced Levantine affinity<sup>12</sup>. Consistent with their high African ancestry, Saudi  
276 cluster12 is plausibly modeled possessing  $0.60 \pm 0.005$  Ethiopia 4500BP ancestry, with the  
277 estimated timing of their Ethiopia 4500BP + MH1-MH2 LT admixture model at  $358 \pm 27$  years.  
278 Ancestry modeling of Saudi cluster5 required a third source component from Egypt's Third

279 Intermediate Period (Egypt 3IP:  $0.71 \pm 0.14$ ; MH1-MH2 LT:  $0.26 \pm 0.13$ ). DATES analysis of Saudi  
280 cluster5's ancestry formation revealed well-fitting estimates for Ethiopia\_4500BP + Egypt 3IP  
281 ( $1120 \pm 330$  years) and Ethiopia\_4500BP + MH2\_MH1\_LT ( $1242 \pm 528$  years) models, while  
282 Egypt + MH2\_MH1\_LT modeling failed statistical fitting criteria. Finally, the unique ancestry  
283 configuration of Saudi cluster11 demonstrated above (**Figures 1C, S3, and S4**), is also manifest  
284 in qpAdm modeling whereby all models across all time periods fit the data poorly. Taken  
285 together, these findings suggest the present-day Saudi Arabian ancestry component was formed  
286 through multiple ancient non-local ancestry contributions to a predominant local Arabian  
287 background – represented by Bahraini groups MH3 LT and MH1-MH2 LT. The majority of  
288 clusters showed compatibility with MH3 LT and Ethiopian ancestry, while specific clusters  
289 exhibited unique patterns: clusters 2 and 8 incorporated Sardinian ancestry, clusters 12 and 5  
290 showed strong African components with MH1-MH2 LT base, and cluster11 displayed a  
291 distinctive genetic profile unable to be plausibly modeled with the current sources, further  
292 revealing the importance of future ancient DNA research in this region.

293



295 **Figure 2. Ancestry compositions in Saudi Arabians estimated with aDNA data as reference.** (A)  
296 Barplots for plausible ( $p$ -value  $\geq 0.01$  and admixture weights between 0 and 1) qpAdm models  
297 grouped by age brackets of source populations (top and bottom; **Methods**). For Pre-Pottery  
298 Neolithic – Neolithic sources (top), three clusters were rejected under the Armenia\_MasisBlur N  
299 + MAR\_Taforalt\_EpiP qpAdm model at the statistical threshold cut-off: cluster12, 3, and 11. We  
300 display under the corresponding qpAdm barplot well-fitting ( $nrmsd < 0.7$  and  $Z > 2$ ) estimates of  
301 admixture timing in years. (B) ‘Basal Eurasian’ ancestry estimated from  $f_4$ -statistic of form  
302  $f_4(\text{Saudi cluster, Han.DG; Ust-Ishim, African aDNA group})$  with varying ancient African groups<sup>47–</sup>  
303<sup>55</sup>. We plotted three standard errors for each  $f_4$ -statistic. The Saudi cluster (y axis) order in each  
304 plot is retained throughout (c12, c3, c11, c4, c8, c7, c10, c1, c5, c2, c6, and c9) following  
305 decreasing value for the statistic  $f_4(\text{Saudi cluster, Han.DG; Ust-Ishim, Ethiopia 4500BP})$ .  
306 Significant (absolute Z-score  $> 3$ ) negative  $f_4$ -statistic values indicate the Saudi cluster possesses  
307 excess shared drift basal to the shared drift between the groups (Han.DG and Ust-Ishim),  
308 commonly interpreted as deriving from a population basal to the OOA event (i.e. the Basal  
309 Eurasian).

310

## 311 **Genetic variation within Saudi is shaped by the social structure**

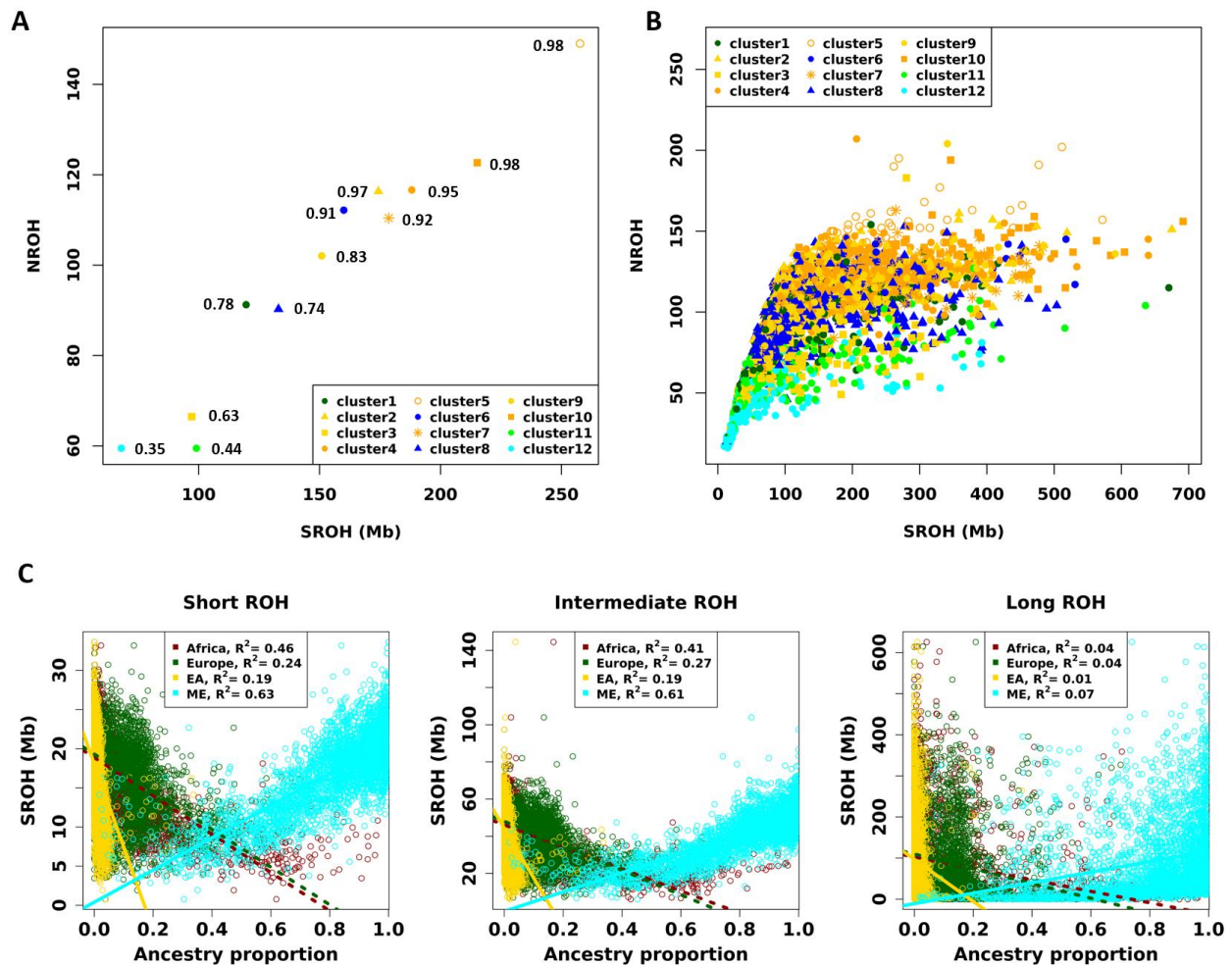
312 The sub-clusters of Saudi Arabians identified in this study exhibits a diverse distribution of runs  
313 of homozygosity (ROH) between the individuals. The general pattern of number of ROH (NROH)  
314 vs. the sum total length of ROH (SROH) showed a distinct relationship with the proportion of  
315 ME-like ancestries (**Figure 3A**). Across Saudi clusters, the median total sum of runs of  
316 homozygosity ranged from 38.12 Mb to 232.6 Mb, while the median number of ROH ranged  
317 from 42 to 150 ROHs. In terms of the mean, cluster12 had the shortest mean length and  
318 smallest mean number of ROHs, followed by cluster11 and cluster3 (**Figure 3A**), which is  
319 characteristic of larger effective population size and consistent with greater admixture from  
320 more diverse African ancestral populations (**Figure 1C**). Cluster5 had the highest burden of ROH,  
321 with highest average number and total length of ROH, followed by cluster10, reflecting a  
322 consequence of both long-term small effective size and/or consanguinity<sup>56</sup>.

323 We also followed a previous approach<sup>57</sup> and divided the ROHs into three classes based on  
324 length:  $< 635$  kb for short ROHs, between 635kb and 1671kb for intermediate ROHs, and  $>$   
325 1671kb for long ROHs (**Methods**). Short ROHs indicate homozygosity from ancient or distant

326 ancestry, i.e. background relatedness. Intermediate ROHs likely arise from background  
327 relatedness with moderate level of inbreeding from past few generations, often due to reduced  
328 population sizes or reproductive isolation (e.g. due to geographic or cultural preferences), or  
329 from recent bottlenecks followed by recovery. Long ROHs indicate recent inbreeding and are  
330 common in populations with high levels of consanguinity<sup>56–58</sup>. When classified by the sizes, we  
331 can observe that the overall pattern of NROH vs. SROH (**Figure 3B**) are driven by the long ROHs  
332 (**Figure S8**). For both the short and intermediate ROHs, there are clear linear relationships  
333 between NROH and SROH (**Figure S8**). In contrast, for the long ROHs, as SROH increase per  
334 individual genome, the NROHs are not increasing at the similar linear pattern as observed for  
335 short and intermediate ROHs. That is, for individuals with greater SROH due to the long ROHs,  
336 they do not have proportionally greater NROH compared to those with less SROH, suggesting  
337 that the contributions of SROHs are driven by fewer but longer ROHs in this length class due to  
338 recent consanguinity. Therefore, the consequences of consanguinity in not only increasing SROH  
339 but also increasing the variance of SROH in a population<sup>56,59</sup>. We also observed the impact of  
340 this when considering each Saudi sub-clusters. In general, clusters 12, 11, and 3 have the fewest  
341 NROH and the smallest SROH across the ROH classes while clusters 5, 2, and 10 tended to have  
342 the most NROH and longest SROHs (**Figure S9A** and **S9B**). The ranked order by both NROH and  
343 SROH across the 12 Saudi clusters were very similar for both the short and intermediate length  
344 ROHs (**Figure S9A** and **S9B**), but varied for the long ROH class, implying a different pattern of  
345 recent inbreeding that differed from ancient demographic events.

346 To further support the relationship between ROH and ancestry components, we modeled the  
347 ROH by ancestry proportion, based on admixture analysis at  $K = 4$  with the HGDP populations.  
348 We found that SROH increase with the increase in ME ancestry proportions, while they are  
349 negatively correlated with the proportion of African, European and East Asian ancestries (**Figure**  
350 **S9C**). This observation is seen across length classes of ROHs, though more attenuated for long  
351 ROHs (**Figure 3C**). We reasoned that this ancestry effect across length classes is likely reflecting  
352 the long-term endogamous marriages and recent consanguinity associated with the ME-like  
353 ancestry relative to admixture component of other ancestries.

354



355

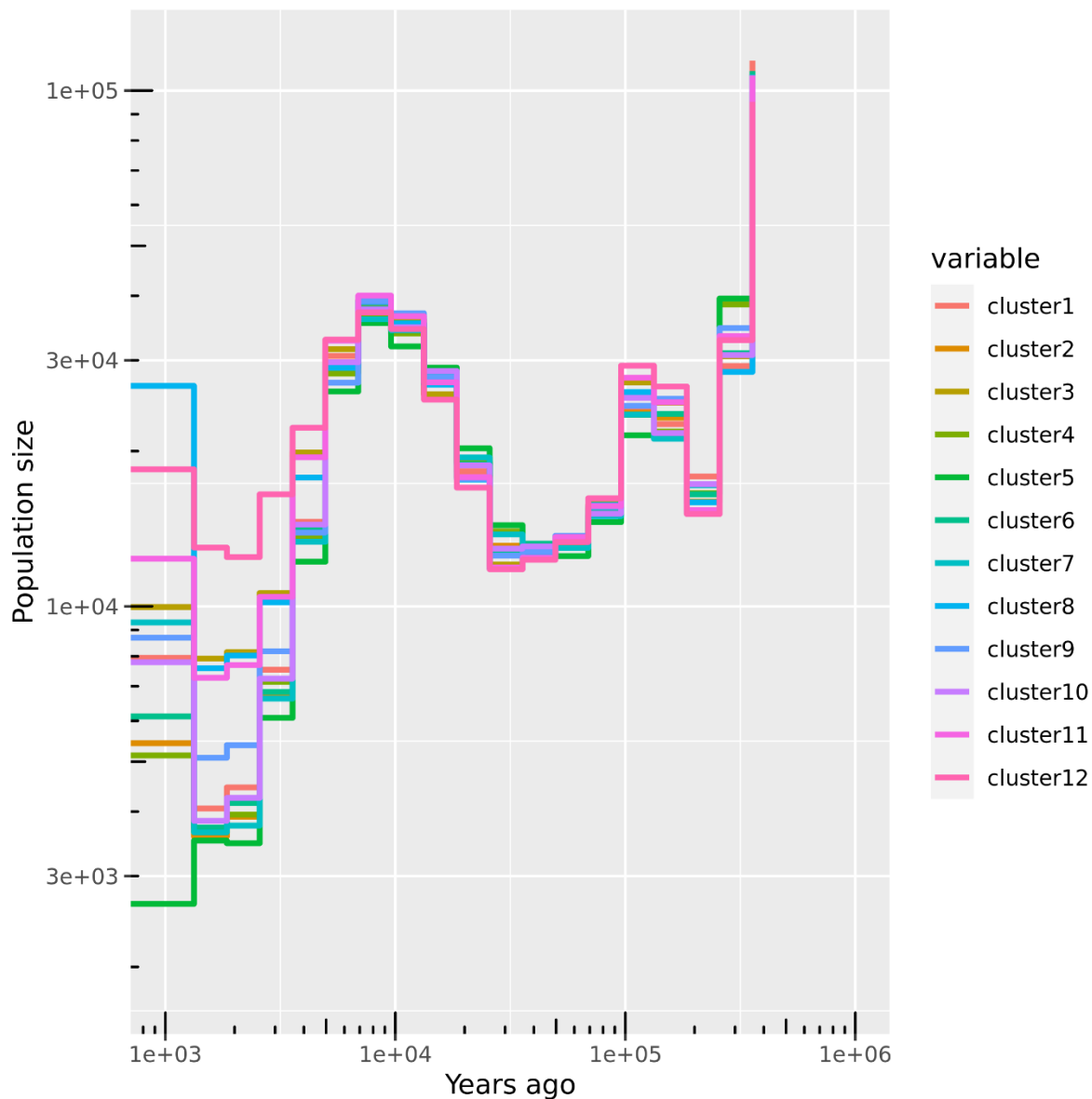
356 **Figure 3. Runs of homozygosity in Saudi Arabians.** (A) Average total length and number of ROH  
 357 per cluster. The numbers next to the symbol represents the mean ME-like ancestry proportion.  
 358 (B) Total length and number of ROH per individual across the Saudi Arabian cohort. For (A) and  
 359 (B), symbols are colored by the geographical region associated with each cluster (**Figure 1D**). (C)  
 360 Total length of ROH vs ancestry proportion per individual stratified by three length classes of  
 361 ROHs. ROH – Runs of homozygosity, ME - Middle Eastern, EA – East Asia.

362

363 We leveraged the dense marker information from the 302 WGS individuals to reconstruct  
 364 genome-wide genealogies and infer the population size trajectories within the Saudi social  
 365 substructure (**Figure 4**). Consistent with the out-of-Africa event, all clusters experienced and  
 366 subsequently recovered from a decline in the effective population size ( $N_e$ ) about 100 kya. All



367 Saudi clusters reached a local maximum  $N_e$  around 9 -10 kya, a period consistent with the early  
368 Holocene Wet Phase / Holocene Humid Period, characterized by wet conditions which resulted  
369 in expansion of lakes and rivers and extensive grasslands<sup>60</sup>. Following the Holocene period,  
370 populations in the Arabian Peninsula experienced another bottleneck dating around 6 - 7 kya,  
371 along with divergence among sub-clusters. This period coincided with the Arabian aridification,  
372 which is responsible for the desert conditions in most of the Arabia as we know it today<sup>6,60</sup>.  
373 Clusters with less ME-like ancestries and stronger signature of admixture (such as clusters 12,  
374 11, 8, & 3), showed less severe decline in  $N_e$  compared to those that have high ME-like  
375 ancestry. Cluster5 in particular, showed the most severe bottleneck and remained low in  $N_e$  in  
376 the recent times, consistent with long-term isolation. Cluster5 appears to resemble the pattern  
377 of the tribe labelled as T25 in a previous study<sup>20</sup>: both originated from the Western region  
378 showing the highest level of inbreeding within the respective study. T25 is said to have been  
379 subjected to strict intratribal marriages. Such social practices can indeed result in persistent  
380 small  $N_e$  as observed here, as well as our observed pattern in ROH (**Figure 3**).



381

382 **Figure 4. Population size trajectories between the Saudi Arabian sub-clusters.** Effective  
383 population sizes were computed from genealogical trees using RELATE (see **Methods**). The  
384 number of samples per cluster used for the estimates can be found in **Table S1**.

385

386 **Saudi's social structure does not impact imputation accuracy but lack of**  
387 **reference representation does**

388 Using the Trans-Omics for Precision Medicine (TOPMed) reference panel<sup>61,62</sup>, we imputed the  
389 genotypes based on Saudi array data and compared the imputation accuracy of the Saudi to the  
390 Europeans from the United Kingdom 10,000 Genomes (UK10K) project matched by sample size  
391 and SNP content to further evaluate the impact population history may have on haplotypic  
392 pattern of variation and implications for genetic epidemiology studies in Saudi today.  
393 Unsurprisingly, the imputation accuracy was lower for Saudi compared to Europeans across all  
394 MAF bins (**Figure S10A**). This is consistent with what was reported in Cahoon et al.,<sup>63</sup> which  
395 showed Saudi Arabia among the populations with the lowest imputation accuracy when  
396 compared to Europeans and populations within North America. Across Saudi sub-clusters, the  
397 imputation accuracies were quite similar, except for a slight difference in cluster12 which had  
398 lower imputation accuracy of common variants (**Figure S10B**). This observation is consistent  
399 with Cluster12 showing elevated shared drift with Africans from the Southern and/or Eastern  
400 region (**Figures 2 and S7**), which may not be well-represented in the TOPMed.

401

## 402 **Allelic architecture of Saudi Arabians**

### 403 **Genetic variation in Saudi Arabian WGS data**

404 Having investigated extensively the population structure, admixture and demographic history  
405 and their impact on ROHs in the Saudi, we then leveraged the WGS data from 302 Saudi  
406 individuals to investigate the consequence of population history on the pattern of genetic  
407 variation in Saudi. A total of 25,488,981 autosomal variants were called and retained after  
408 quality control (QC) (**Methods**), of which 2,459,950 (9.7%) variants were not previously  
409 identified in gnomAD v4.1<sup>64,65</sup> and thus are potentially novel or Saudi-specific. We refer to  
410 these variants as the “previously unknown variants”. As expected, the previously unknown  
411 variants are highly enriched with rare alleles (83% of them are singletons in our dataset,  
412 compared to 32% singletons among the known variants; **Figure S11**). Of all variants, 63% have  
413 MAF < 1%. Although there are some Middle Eastern individuals in gnomAD v4, they are

414 proportionally underrepresented in this global dataset (2,884 exomes among 730,947 total, 147  
415 genomes among 76,215 total), resulting in a large number of Saudi Arabian variants missing  
416 from the database (although pipeline differences may explain some of the missing variants).  
417 Previously, Almarri et al.,<sup>6</sup> found that of 23.1 million single nucleotide variants identified in 147  
418 Arabian and Levantine individuals, 4.8 million (20.8%) were not found in the Human Genome  
419 Diversity Project (HGDP-CEPH) global dataset. Taken together, both our and Almarri et al.,  
420 studies showed that variation from the Arabian Peninsula are not yet well captured. Thus,  
421 genetic association studies will be limited to only the common variation, which we have also  
422 shown to be currently sub-optimally imputed (**Figure S10**). Here we provide all the 25,488,981  
423 variants that remained after QC and their allele frequencies, see the “**Data availability**” section  
424 for access.

425 We compared allelic frequency spectra and allelic homozygosity in the Saudi Arabians (all WGS  
426 individuals) to the Middle Eastern population in gnomAD (gnomAD-MID). The two have  
427 relatively similar patterns in the genome-wide alternative allele frequency spectra though Saudi  
428 had proportionally slightly fewer common variants (**Figure S12A**). The allele frequencies are  
429 highly concordant ( $r = 0.98$ ) between the two populations (**Figure S12B**), but Saudi Arabians  
430 have approximately 2x more homozygous genotypes than gnomAD-MID (e.g. an average of 20%  
431 vs. 10% of the genotypes are homozygous for variants with alternative allele frequency > 5% in  
432 Saudi and gnomAD-MID, respectively; **Figure S12A** and **S12C**). The higher proportion of  
433 inbreeding suggests that the Saudi and the gnomAD-MID population are not reflective of the  
434 same underlying populations. However, because the frequency spectra and correlation of allele  
435 frequencies are highly similar (**Figure S12**), we thus utilize both samples to compare the pattern  
436 of variation with gnomAD African/African Americans (gnomAD-AFR) and non-Finnish Europeans  
437 (gnomAD-EUR) below to better understand the impact of the unique history in the Arabian  
438 Peninsula on its current pattern of variation.

439

440 **Distribution of functionally deleterious variants**

441 We annotated the variants using three different annotation tools: VEP (v.110) (McLaren et al.,  
442 2016), AlphaMissense (Cheng et al., 2023), and Genomic Pre-trained Network (GPN) (Benegas  
443 et al., 2023). AlphaMissense predicts the pathogenicity of missense variants<sup>66</sup> while GPN  
444 predicts the deleteriousness for both coding and non-coding variants. The distribution of the  
445 variants by functional classes are shown in **Table S18**. We first examined the set of previously  
446 unknown variants, relying on AlphaMissense and VEP only as GPN is precomputed only for  
447 variants found in gnomAD<sup>67</sup>. In addition to being enriched with rarer alleles, proportionally  
448 more unknown variants (19.63%) were annotated to be deleterious than known ones found in  
449 gnomAD (7.2%). This implies that the previously unknown variants are not just sequencing  
450 errors distributed randomly across the genome, but are enriched for rare variants of functional  
451 relevance in the Saudi that are maintained or have not been purged from the population.

452 Non-Africans are expected to have more deleterious alleles due to the relaxation of purifying  
453 selection during the Out-of-Africa (OOA) bottleneck as well as the introduction of new  
454 deleterious mutations during population expansion<sup>30,68,69</sup>. Despite some opposing reports on  
455 this hypothesis<sup>70,71</sup>, there has been empirical evidence in isolated populations having an excess  
456 of functionally deleterious alleles<sup>72-74</sup>. In addition to the OOA bottleneck, Saudi has a deep  
457 culture of endogamy and consanguinity, and these demographic factors are known to  
458 potentially increase the burden of deleterious alleles in a population due to decreased efficacy  
459 of purifying selection<sup>33</sup>. Here we investigated the allelic architecture of functionally deleterious  
460 alleles in the Saudi population, compared to other continental populations from gnomAD.

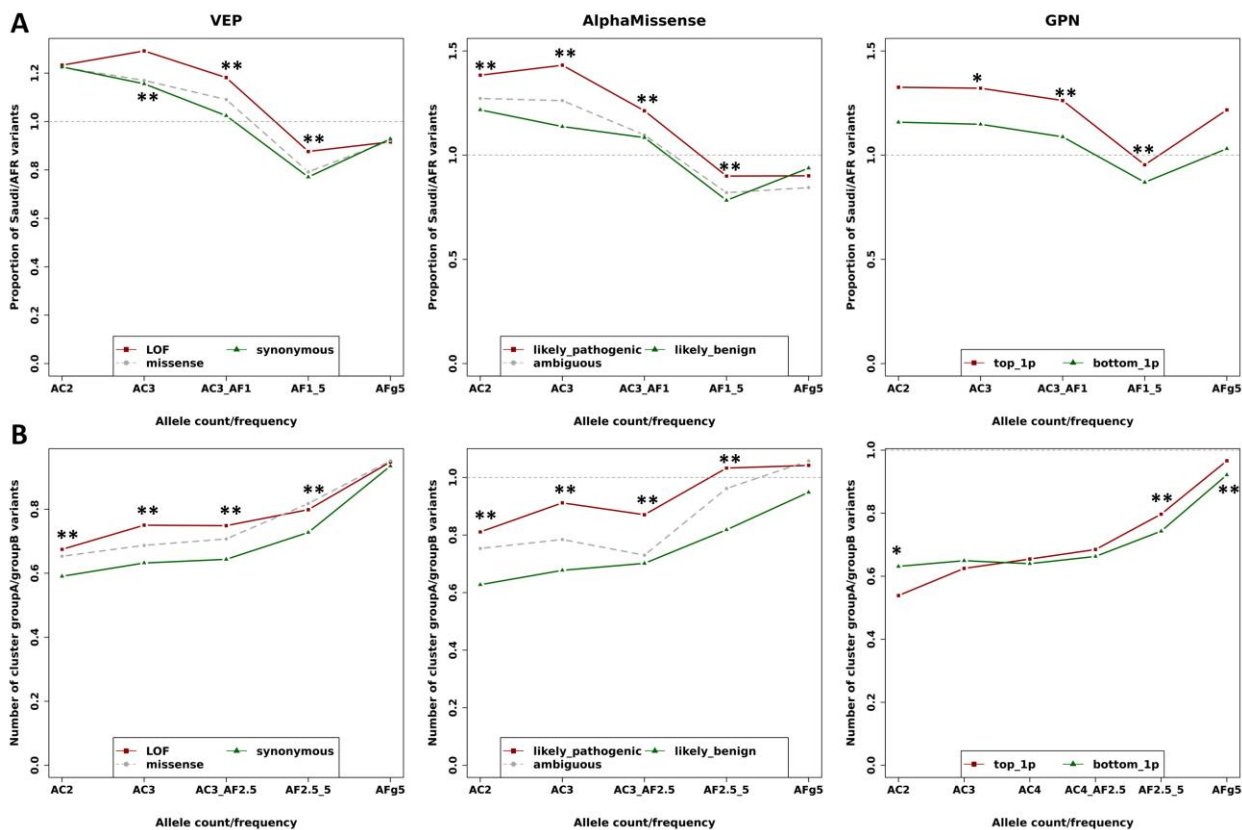
461 Compared to gnomAD-AFR individuals, the Saudi tend to show proportionally more deleterious  
462 alleles than those annotated to be benign or neutral across algorithms (**Figure 5A**), particularly  
463 for variants up to ~5% frequency. Overall, relative to gnomAD-AFR, between the 0.5 - 5%  
464 frequency, we found a 13% proportional increase of deleterious (likely pathogenic) alleles  
465 annotated by AlphaMissense in the Saudi Arabians compared to 7% proportional decrease of  
466 the benign alleles ( $P < 0.01$ ; **Figure 5A**). When annotated by VEP and GPN, at the same  
467 frequency range, we observed a consistent pattern i.e. a 3% proportional increase in loss of  
468 function variants in the Saudi Arabians compared to 10% proportional decrease in neutral  
469 (synonymous) ones ( $P < 0.01$ ) by VEP, and an 11% proportional increase in the first percentile of

470 alleles by deleteriousness compared to 3% proportional decrease in the 99<sup>th</sup> percentile (e.g. the  
471 most likely neutral) of alleles when annotated by GPN (**Figure 5A**).

472 This pattern of enrichment for functionally deleterious alleles is also qualitatively observed  
473 when comparing the exome samples from gnomAD-MID to gnomAD-AFR (**Figure S13**), taking  
474 advantage of the larger sample size for gnomAD-MID exomes and the same data processing  
475 pipeline in gnomAD. The pattern is also qualitatively observed when comparing Saudi to  
476 gnomAD-EUR, though the difference may be more attenuated in some allele frequency bins  
477 (e.g. for AlphaMissense annotation; **Figure S14**). The less significant finding when comparing  
478 Saudi to gnomAD-EUR is probably because Europeans also showed proportionally more  
479 deleterious than neutral alleles across all frequency bins, as previously reported<sup>68,75</sup> and  
480 replicated here (**Figure S15**).

481 We also compared the enrichment of deleterious alleles between Saudi sub-clusters. Because of  
482 the smaller number of individuals within each cluster having WGS data (**Table S1**), we grouped  
483 the clusters into two groups: groupA which contained clusters with greater inbreeding and  
484 lower effective population sizes (clusters 2, 4, 5, 6, 9, and 10), and groupB which has less  
485 inbreeding and higher effective population sizes (clusters 12, 11, 3, and 8). We left out cluster1  
486 from this analysis as it tends to fall in the middle of the two groups. GroupB had generally  
487 greater number of variants compared to groupA (**Figure S16**), consistent with its higher genetic  
488 diversity and less inbreeding. GroupA, with greater inbreeding and lower effective population  
489 sizes, showed greater enrichment of deleterious alleles (**Figure 5B**).

490



491  
 492 **Figure 5. Distribution of minor allele frequency across functional classes.** (A) Ratio of Saudi to  
 493 gnomAD-AFR variants. The sample size of gnomAD-AFR is based on downsampling to Saudi  
 494 sample size,  $n = 302$ . (B) Ratio of Saudi cluster groupA to cluster groupB variants. The sample size  
 495 of cluster groupB is based on downsampling to groupA sample size,  $n = 124$ . Variant functional  
 496 consequences were annotated based on VEP (loss-of-function, missense, or synonymous  
 497 variants), AlphaMissense (likely pathogenic, likely benign, and ambiguous), and GPN. AC and AF  
 498 refer to allele count and allele frequency, respectively. AFG5 refers to allele frequency greater than  
 499 5%. Top\_1p refers to variants with the top 1% of GPN scores (more deleterious) and Bottom\_1p  
 500 refers to variants with the bottom 1% of GPN scores (more neutral). AFR denotes the gnomAD-  
 501 AFR sample. LOF refers to Loss of function. \*\* and \* denote frequency bins with significant  
 502 difference between the most deleterious (red) and most neutral (green) through bootstrapping  
 503 at  $p < 0.01$  and  $< 0.05$ , respectively.

504

## 505 DISCUSSION

506 Scholars have noted the complexity of diverse histories in shaping the genetic architecture of  
 507 Arabian Peninsula populations, and called for better characterization of each population for  
 508 better understanding of their genetics and health<sup>35,76</sup>. On one hand, being situated

509 geographically at the crossroads between Africa and Eurasia, which facilitates intercontinental  
510 interactions, is expected to increase heterozygosity and genetic diversity in the AP populations.  
511 On the other hand, Saudi Arabian culture is rooted in endogamous practices which increases  
512 homozygosity which can result in health consequences <sup>22,25</sup>. Here we elucidated the fine-scale  
513 population structure of Saudi Arabian population using 3,252 genotyped and 302 WGS  
514 individuals from various geographic regions within the country and investigated the genetic  
515 variation within the social structure and impact of the demographic histories on the population  
516 health.

517 Our analyses concur with previous findings about the presence of sub-population structure  
518 within Saudi Arabia <sup>20,77</sup>, with twelve distinct genetic sub-clusters being identified in our data.  
519 We chose to infer the sub-clusters based on genetic similarities and use the resulting sub-  
520 clusters as units of analysis throughout the study. We note that clustering by genetic similarity  
521 in an admixed population could misrepresent the population structure in the dataset, such as  
522 when multiple clusters of Arabian origin are combined into a single cluster due to sharing a  
523 common admixing source (e.g. African). However, we elected for this approach in part because  
524 of limited tribal affiliation information at the individual level due to privacy concerns, thus we  
525 could not rely on self-reported tribes as units of analysis. Moreover, ancestry-specific  
526 approaches to infer population structure <sup>78,79</sup> are limited to situations where the admixing  
527 ancestries are divergent and that references for the ancestral populations are available, both of  
528 which are understudied or unavailable for Saudi Arabians. While our observation of population  
529 structure should be re-examined in the future with more self-reported demographic  
530 information or with improved methodologies to exclude the impact of admixture, we also note  
531 that our clustering is not driven solely by non-Arabian admixture. For instance, both cluster12  
532 and cluster3 exhibit strong levels of African admixture (**Figure 1C; Figure 2A - B**), but they also  
533 showed affinity to different Bahraini aDNA references with varying levels of Levantine-related  
534 ancestry (**Figure 2A**). Therefore, in this case, level of African admixture is not the only reason  
535 that separated clusters 12 and 3 from the rest of the Saudi Arabian with less admixture.  
536 Furthermore, the fine-scale structure in the Saudi Arabian population, brought on at least in  
537 part by the social practices, has previously been reported when studying self-identified



538 indigenous tribes<sup>20</sup> and we do observe some similarities of such structure with our clusters.  
539 Proceeding with the genetic sub-clusters defined, we then continue to elucidate several key  
540 features in the pattern of genetic variation in Saudi Arabia.

541 While geographic proximities increase the chances for gene flow between individuals around  
542 same geographic areas, the strong consanguineous and endogamous culture in Saudi Arabia is  
543 expected to limit such interactions, resulting in distinct sub-clusters residing in relatively close  
544 proximities that we and others have observed<sup>20</sup>. We also note that, while recent and/or  
545 ongoing admixture may be taking place between the Saudi Arabians and modern-day Africans,  
546 Europeans and Central & South Asia, the spread of such genetic pool is likely restricted by the  
547 endogamous and consanguineous practices as previously suggested for the Emiratis<sup>35</sup>. Both  
548 endogamy and consanguinity increase the burden of ROH (SROH) which increases health risks.  
549 In El-Mouzan et al.,<sup>80</sup> the highest rate of consanguineous marriages was found in Madinah  
550 from the Western region of the country. Similarly, our results show footprints of inbreeding in  
551 all the sub-clusters from the Western region (clusters 4, 5, 7, and 10) as demonstrated by the  
552 longest ROHs, lowered  $N_e$ , and less evidence of admixture. On the other hand, the sub-groups  
553 that intermarry benefit from the increased effective population sizes and genetic diversity,  
554 especially those with elevated African admixture. The sub-clusters with elevated admixture in  
555 our study also show reduced ROH numbers and sizes.

556 Consistent with the Arab slave trade and Islamic expansion in the 7<sup>th</sup> (1,300 – 1,400 years ago)  
557 century, almost all of the Saudi clusters show a recovery in effective population sizes from the  
558 Arabian aridification (about 6 - 7 kya) bottleneck, at varying ranges, except for the isolated sub-  
559 cluster from the Western region (cluster5). Whilst the admixture dates for Saudi cluster5 qpAdm  
560 model of Ethiopia\_4500BP + Egypt (1120±330 years ago) and Ethiopia\_4500BP + MH2\_MH1\_LT  
561 (1242 ± 528 years ago) align with the 7th century Islamic expansion we note that historical  
562 documentation of significant North African to Arabian migrations during this period remains  
563 limited and thus, the possible conduit for Northern/Northeastern African ancestry to Saudi  
564 cluster5 remains inconclusive. The increase in the effective population sizes from the recent  
565 times is more for sub-clusters with higher African admixtures. Previous studies have reported

566 dominating African admixture source in Arabia originating from Bantu speakers from the East  
567 (Kenya) or South Africa, dating from 400 to 1754 years ago<sup>6,13,81</sup>, which is consistent with the  
568 Arab trade slave expansion in the 7<sup>th</sup> (1,300 – 1,400 years ago) century. In a recent study<sup>20</sup>, the  
569 Saudi tribes with highest African admixture were results of recent admixture events, as recently  
570 as 11 generations ago. This may be a sign of continuous admixture beyond the Arab trade slave  
571 expansion. Indeed, we observe a recent ( $358 \pm 27$  years ago) estimated timing of Saudi  
572 cluster12 ancestry formation of Ethiopia 4500BP + MH1-MH2 LT. This recent timing, coupled  
573 with distinct Central, Eastern, and Southern African genetic affinity, suggests a possible  
574 connection to 17<sup>th</sup> - 20<sup>th</sup> century Red Sea and trans-Saharan slave trades, as opposed to the  
575 older North African Islamic expansion, during which East African regions, including Ethiopia and  
576 Eritrea (historically referred to as Abyssinia), Somalia, and Sudan, were significant sources of  
577 enslaved individuals<sup>82</sup>. Additionally, Central African regions extending into present-day Chad  
578 and the Congo Basin contributed to this population due to extensive trade networks.

579 The impact of the OOA bottleneck and population size histories on the allelic architecture of  
580 deleterious alleles in non-African populations has been a matter of debate<sup>30–32,68–71,73</sup>. We  
581 observed an abundance of rare and low frequency (AF = 0.5 – 5%) deleterious alleles in Saudi  
582 Arabians when compared to gnomAD-AFR. This enrichment in deleterious alleles can be  
583 explained to some extent by the demographic history of the Saudi Arabians beyond the OOA  
584 bottleneck event. First, the high percentage of previously unknown deleterious mutations that  
585 we found in sequencing could be driven by consanguinity/inbreeding. Second, although  
586 consanguinity and endogamy could in principle expose deleterious alleles to negative selection  
587 through genetic purging, purging is less effective when effective population sizes are small  
588 thereby the deleterious alleles may drift to high frequencies and even become fixed<sup>33,75</sup>,  
589 particularly if population isolation followed a bottleneck<sup>31,72,73,83</sup>. Saudi Arabia is not a  
590 completely isolated population, but there is reproductive isolation due to their social practices.  
591 Taken separately, the subgroups with high prevalence of endogamy are indeed enriched for  
592 deleterious alleles than those that have high levels of admixture (**Figure 5B**). Our results are  
593 consistent with other empirical studies that show enrichment of deleterious in populations with

594 different demographic histories following population bottleneck <sup>72,75</sup> and further confirms no  
595 evidence of genetic purging in the Saudi Arabians as previously reported <sup>26</sup>.

596 Overall, our results shows that Saudi's population history impacts its pattern of genetic variation  
597 with potential consequences to the population health. The legacy of endogamy and  
598 consanguinity in the population poses health risks and the frequency of this practice has not  
599 shown signs of decline <sup>84,85</sup>. There have been initiatives to raise public awareness on the health  
600 risks of close relative marriages in Saudi Arabia and other Middle East countries which include  
601 mandatory premarital health screening for recessive illnesses and genetic counselling <sup>22,23,86,87</sup>.  
602 Even though the mandatory premarital screening in Saudi Arabia has so far not been very  
603 successful in discouraging or preventing at-risk marriages, it has fostered a more informed pre-  
604 birth decisions, reducing the prevalence of children born with the health complications through  
605 altering some of the cultural behaviors including adoption of prenatal detection and therapeutic  
606 abortion <sup>14,85,87</sup>. With increasing public education and awareness on risk factors associated with  
607 endogamous and consanguineous unions, its prevalence may change in the future, especially  
608 because the consanguineous marriages are generally most prevalent in poor, rural and least  
609 educated societies of Arabia compared to urbanized and more educated counterparts <sup>88</sup>.

610 We note a common issue with regards to the under-representation of Arabian countries in  
611 global cohorts which is also raised by others <sup>6,26,89</sup>. For example, for the genomes in gnomAD  
612 database, only 0.2% of samples are of Middle Eastern origins, compared to 44.6% and 27.3% as  
613 Europeans and Africans, respectively. Even with the newly increased exome data, the Middle  
614 Easterners only make up 0.38% of gnomAD, compared to 4.65% Africans and 77.07%  
615 Europeans. As a result, the human genetics field in general is missing variants that are enriched  
616 to Saudi Arabia and Middle East. It has also been suggested that GME populations tend to  
617 harbor more variants unique to the region <sup>26</sup> e.g. 28% in the Qatar <sup>90</sup>. The under-representation  
618 accentuates the understudying of regionally enriched alleles, and contributes to reduced  
619 accuracy of polygenic prediction models when applied to Arabic populations <sup>91</sup> and lower  
620 imputation accuracy when using state-of-the-art reference panel like TOPMed <sup>63</sup> (**Figure S10**).

621 This further highlights the need to sequence diverse and unique populations and include them  
622 in large global genetic data platforms such as gnomAD, TOPMed and others.

623

## 624 **METHODS**

### 625 **Data collection, processing and quality control**

626 For all studied samples, written informed consent was obtained from each participant. The  
627 studies were approved under the Saudi Genome Project by the Institutional Review Board at  
628 King Abdulaziz City for Science and Technology and King Fahad Medical City. In compliance with  
629 Saudi privacy legislation and the protection of human subject confidentiality, the sharing of raw  
630 genotyping and clinical data is restricted. Access to this data requires prior approval from the  
631 Saudi National Bioethics Committee.

### 632 **Array data**

633 **Sample collection, genotyping and quality control.** A total of 3,752 samples were collected in  
634 Saudi Arabia between the years 2017 - 2020 as control individuals for various projects, such as  
635 the GenOMICC International project and covid19 host genetics consortium <sup>92</sup> studies.

636 Individuals were genotyped on the Axiom Genome-wide CEU 1 Array including customized  
637 variants following the manufacturer's specifications for sample preparation, including whole  
638 genome amplification, fragmentation, denaturation, and hybridization. Genome-wide SNP  
639 genotyping was performed using the automated, high-throughput GeneTitan system from  
640 Affymetrix.

641 We filtered individuals with sample call rates  $< 0.9$  using PLINK v1.9 <sup>93,94</sup> on each plate  
642 individually before merging the autosomal SNPs across the different plates, resulting in a  
643 merged set of 757,790 SNPs. We removed duplicates and non-biallelic variants, retaining  
644 703,986 SNPs. We then filtered SNPs with greater than 10% missing rate and SNPs that did not

645 pass Hardy Weinberg Equilibrium (HWE) test ( $P < 10^{-6}$ ) using PLINK, resulting in a total of  
646 606,349 SNPs for analysis. We lifted over the genomic coordinates from human reference  
647 genome hg19 to hg38. We phased the data using Beagle v5.2<sup>95</sup>.

648 **Removing close relatives and filtering out outliers.** Using the 3,752 Saudi samples and 606,349  
649 SNPs, we pruned the dataset by linkage disequilibrium (LD) (using the command `--indep—`  
650 `pairwise 50 5 0.8` in PLINK), resulting in 547,307 SNPs to estimate individuals' relatedness using  
651 King v2.2.5<sup>96</sup>. We removed twins (or duplicated individuals) as well as first degree relatives,  
652 retaining 3,403 samples. Furthermore, we performed PCA and performed two iterations of  
653 outlier (defined as being  $> 6$  standard deviation (SD) away from the mean in any of the first 10  
654 PCs), resulting in 3,352 samples left for further analyses.

655 **Defining samples' tribal affiliation and imputing missing tribal information.** We aimed to use  
656 available demographic information, i.e. tribal affiliations, in validating and interpreting the  
657 results of clustering based on genetic data. However, 82% of the individuals in our data (2,740  
658 of the 3,352) did not have self-reported tribal information. We thus imputed such information  
659 using the software HARE (harmonized ancestry and race/ethnicity) package<sup>97</sup> based on the  
660 available Self-identified Race/Ethnicity (SIRE) tribal information of 612 individuals. SIRE in our  
661 data were derived from either self-report or individual's family name that is presumed to reflect  
662 their tribal affiliation (**Table S1**). The HARE package combines genetically inferred structure  
663 based on PCA with available SIRE information to train a support vector machine (SVM) classifier  
664 that could correct for potentially mislabeled SIRE and predict the race/ethnicity, in this case  
665 tribal label, for those individuals missing SIRE. We used the HARE to impute tribal information of  
666 the samples missing a SIRE label in our dataset using the first 30 PCs as the input data. We used  
667 the highest predicted membership probability ( $L_1$ , see Fang et.,<sup>97</sup> for more details) labels to aid  
668 in the interpretation of the population sub-clusters that we infer from genetic data.

669

670 **Whole genome sequencing (WGS) data**

671 **Sequencing information and processing.** In addition to the genotyped samples, 349 samples  
672 were whole-genome sequenced (WGS) to a targeted depth of 30x. The samples were prepared  
673 following the Illumina's TruSeq Nano sample preparation protocol and sequenced on an Illumina  
674 HiSeq X-ten machine. The raw sequences were aligned against the human reference genome  
675 GRCh38 using the Burrows-Wheeler Aligner (BWA) version 0.7.10<sup>98</sup>. Picard tools version 1.117  
676 was used to mark duplicates<sup>99</sup>. All sample preparation, sequencing, sequence alignment, pre-  
677 processing, quality control before calling of variants and BAM file augmentation were  
678 performed by deCODE genetics (<https://www.decode.com>), and a more detailed information on  
679 these steps is documented in Jónsson et al.,<sup>100</sup>.

680 **Variant calling and filtering.** We merged the gVCFs of the 349 samples using CombineGVCFs in  
681 GATK<sup>99</sup> and subsequently performed a joint genotyping calling using GenotypeGVCFs. We  
682 performed variant quality score calibration (VQRS) on the combined samples using  
683 VariantRecalibrator and ApplyVQRS in GATK<sup>101</sup>. We supplied the homo sapiens reference  
684 assembly 38 (Homo\_sapiens\_assembly38.fasta) and used the following resources: HapMap III  
685 variants were used as training and truth sets with prior priority of 15, 1000G omni2.5 sites were  
686 used as training set with prior priority of 12, 1000G phase1 high confidence SNPs was used as  
687 training set with prior priority of 10 and the dbSNP138 as known SNPs with prior probability of  
688 2. For the annotations, we included the QD, MQ, MQRankSum, ReadPosRankSum FS and SOR.  
689 We used 99% sensitivity level to filter the SNPs.

690 **Quality control on samples and SNPs.** All 349 samples had missing genotyping rate < 10%. We  
691 excluded 302,640 SNPs with missing rate > 10% and 53,981 SNPs based on HWE threshold ( $P <$   
692  $10^{-6}$ ), leaving 26,781,476 SNPs. We removed non-biallelic sites which left 26,408,559 variants.  
693 Further filtering was applied on specific downstream analyses when appropriate. To exclude  
694 outliers in our samples, we merged the 349 WGS samples with our array data and the HGDP  
695 dataset at segregating SNPs shared across all datasets. A principal component analysis (PCA)  
696 was performed using PLINK and we used HARE to impute missing self-reported individual  
697 nationalities (e.g. self-identified nationality as Saudi or not). We excluded 8 samples which were

698 not imputed as a Saudi. We then removed monomorphic sites which were introduced by calling  
699 the variants including these potentially non-Saudi samples, leaving 25,488,981 variants.

700 We filtered samples based on relatedness using King software v2.2.5<sup>96</sup>. For estimating the  
701 relatedness, we randomly sampled 550,000 SNPs with minor allele frequency > 1% after LD  
702 pruning (*--indep-pairwise 50 5 0.5* using PLINK) to estimate the relatedness. We removed 37  
703 twins/duplicates and first-degree relatives. Using the PCA, we further removed 2 samples that  
704 appeared as extreme outliers (> 6 SD on any of the first 10 PCs), leaving 302 samples. Haplotype  
705 phasing was performed on the remaining 302 samples and 25,488,981 variants using Eagle  
706 v2.4.1<sup>102</sup>.

707

## 708 **Annotation of variants**

709 We annotated the variants using the popular VEP (v.110)<sup>103</sup> as well as two recently published  
710 annotation tools, AlphaMissense<sup>104</sup> and Genomic Pre-trained Network (GPN)<sup>67</sup>. The  
711 AlphaMissense only annotates missense variants and has three functional classes, “likely  
712 pathogenic”, “ambiguous” and “likely benign”. The GPN annotates all genomic variants and  
713 assign a deleteriousness score to each variant in gnomAD (v3). We downloaded the pre-  
714 computed scores from  
715 <https://huggingface.co/datasets/songlab/gnomad/resolve/main/test.parquet>, accessed  
716 2/9/2024.

717

## 718 **Merging of Saudi whole-genome-sequence data with ancient genomes**

719 We downloaded the Allen Ancient DNA Resource (AADR) v.54.1 Eigenstrat files which are  
720 genotyped according to hg19 coordinates. We were kindly provided the Bahrain aDNA<sup>12</sup>  
721 genome bed files by Rui Leite Portela Martiniano, which were originally mapped to GRCh38. To  
722 have all our genotype files on a consistent reference genome, we mapped the Saudi whole-

723 genome-sequenced data and ancient Bahrain samples back to human reference genome hg37  
724 using liftOver. We then filtered the variants through PLINK 1.9-beta7 using parameters *--geno 0 --*  
725 *snps-only --make-bed --allow-no-sex*. Prior to merging the Saudi WGS, AADR, and Bahrain  
726 datasets we filtered mistyped SNPs where the rsIDs are shared between the datasets, but the  
727 AADR reference allele does not match either the reference or alternate allele (n=520). The Saudi  
728 WGS and Bahrain aDNA datasets were then converted from plink to packedancestrymap format  
729 through Eigensoft convert function with parameter *familynames: NO*. Finally, we merged the  
730 Saudi WGS and AADR datasets with the Eigensoft mergeit function with parameters  
731 *strandcheck: YES*. The mergeit program merges two data sets into a third, which has the union  
732 of the individuals and the intersection of the SNPs in the first two. We first merged the AADR  
733 and Saudi WGS datasets. The merging of the AADR and Saudi datasets resulted in the filtering of  
734 14,239 SNPs due to A/T or C/G strand checks and 91 SNPs due to allele mismatch. In addition,  
735 there were 770,115 genotype strand flips with the final dataset consisting of 1,032,250 retained  
736 SNPs. We then merged the Bahrain aDNA resulting in a final packedancestrymap genotype file  
737 of 1,030,352 SNPs.

738

## 739 **Data analyses**

740 We used the larger collection of Saudi genotyped samples to investigate the genetic  
741 substructure and historic admixtures of the population. We then utilized the high-density  
742 genome-wide marker information from the WGS data to investigate differences in genetic  
743 ancestries with aDNA, population size trajectories, and allelic architecture of functional variants  
744 within the social structure of the Saudi population.

## 745 **Evaluation of population structure**

746 We merged the fully filtered array and WGS datasets, based on segregating markers. We  
747 performed PCA followed by UMAP<sup>105</sup> to combine the first 10 PCs and reduced them into two-  
748 dimensions in order to explore the population structure. Based on the UMAP results, we



749 assigned individuals to subpopulations using K-means clustering from the R package *stats*<sup>106</sup>. To  
750 determine the optimal number of K clusters, we used the Average Silhouette Width (ASW)  
751 (**Figure S17**;<sup>107</sup>) which is a popular and trusted method to produce quality clustering<sup>108</sup>. The  
752 ASW uses values between -1 and 1 to measure how similar/dissimilar is an object to others  
753 within its cluster as well as objects in different clusters, with higher numbers representing a  
754 better fit and appropriateness of clustering. Likewise, a high ASW value corresponds to an  
755 optimal number of K clusters for partitioning a particular set of objects<sup>107</sup>. We validated these  
756 clustering by evaluating the concurrence between the clusters and the tribal region  
757 assignments. We used these clusters as representative of the social structure and also used  
758 them in the whole genome sequencing samples to evaluate patterns of genetic diversity within  
759 the Saudi population.

760 **Analysis of ancestry components.** We conducted the unsupervised admixture analysis using  
761 ADMIXTURE software v1.3<sup>109</sup>. We conducted 10 independent runs of admixture analysis for  
762 each K and retained the run with maximum likelihood. We used the cross-validation procedure,  
763 implemented in the program, to identify the best number of ancestral populations K which fits  
764 our data.

765 **Evaluating patterns of admixture.** To further test for the presence of admixture within the  
766 identified clusters, we performed supervised admixture analysis using the  $f_3$ -statistics from the  
767 ADMIXTOOLS 2 package v2.0.4<sup>110</sup>. We computed the  $f_3$ -statistics using the Saudi clusters as  
768 targets and using all pairs of populations in the HDGP data<sup>34</sup> as source populations i.e.  $f_3$ (Saudi  
769 cluster; HGDP population 1, HGDP population 2).

770 We also used the outgroup  $f_3$ -statistics to investigate the degree of shared drift between Saudi  
771 clusters and the HGDP populations. For this statistic, the African San, Mbuti or Yoruba are often  
772 considered as outgroups for investigation of non-African populations. However, HGDP African  
773 populations have showed to be highly admixed with some of the Saudi clusters, and the  
774 outgroup should be close enough but should not be part of the ingroups<sup>111</sup>. We first used  $f_4$ -  
775 statistics of the form  $f_4$ (HGDP population, HGDP population; Saudi cluster, Saudi cluster) to

776 determine a suitable outgroup for the Saudi clusters from the HGDP reference populations, as  
777 recommended by Pattersons et al.,<sup>111</sup>. We found that nearly all HGDP population combinations  
778 showed positive gene flow between the HGDP population and either one or both of the Saudi  
779 clusters, thereby violating the outgroup assumption. The HGDP populations that did not violate  
780 the outgroup assumption in this test were eight East Asian populations: Han, Miao, Japanese,  
781 Tujia, Yi, Hezhen, She and Naxi. We then used Han as an outgroup for the outgroup  $f_3$ -statistics  
782 in the form:  $f_3(\text{outgroup}; \text{population1}, \text{Saudi cluster})$ , whereby population1 was an HGDP  
783 population or another Saudi cluster except the target.

784 **Saudi and ancient genome analyses.** To assess the genetic affinity of Saudi clusters to present-  
785 day and ancient African groups, we computed the  $f_3$ -statistic of form  $f_3(\text{Han}; \text{African groups},$   
786 Saudi clusters) using the ADMIXTOOLS 2 package v2.0.4<sup>110,112</sup>. We removed individuals that  
787 were indicated on the AADR metadata as relatives, contaminated, duplicated, or have low  
788 coverage. We tested the below allele sharing pattern using  $f_4$ -statistics in ADMIXTOOLS 2 with  
789 parameters  $f_4\text{mode} = \text{TRUE}$ ,  $af\text{prod} = \text{TRUE}$ ,  $allsnps = \text{TRUE}$ . To estimate ‘Basal Eurasian’  
790 ancestry in the Saudi clusters using a selection of eight ancient African groups:  
791 MAR\_Taforalt\_EpiP<sup>49</sup>, Egypt\_ThirdIntermediatePeriod<sup>51</sup>, Ethiopia\_4500BP<sup>54,113</sup>,  
792 Kenya\_Nyarindi\_LSA\_Kansyore<sup>114</sup>, Tanzania\_Zanzibar\_1300BP<sup>115</sup>, Malawi\_Fingira\_LSA\_6000BP  
793<sup>54</sup>, South\_Africa\_400BP.SG<sup>55</sup>, and Cameroon\_SMA<sup>54</sup> in the  $f_4$ -statistic of form  $f_4(\text{Saudi cluster},$   
794 Han.DG<sup>34,116,117</sup> Ust-Ishim, African group). The North African Epipalaeolithic Moroccan  
795 Iberomaurusian Taforalt group (Taforalt EpiP) represents the best proxy of Basal Eurasian  
796 ancestry<sup>7</sup>, exhibiting genetic connections to both early Holocene Near Easterners, such as  
797 Levantine Epipaleolithic Natufians (Natufian EpiP), and sub-Saharan Africans. Thus, to further  
798 assess the genetic affinity of Saudi clusters to African ancestry relative to the shared ancestry of  
799 North African Upper Paleolithic Taforalt Moroccan and Epipaleolithic Natufian Levantine groups,  
800 we ran the  $f_4$ -statistic of form  $f_4(\text{Saudi cluster}, \text{Upper Paleolithic Taforalt}; \text{African group},$   
801 Epipaleolithic Natufian). Finally, to assess the relative shared drift between Epipaleolithic  
802 Levantine Natufian<sup>8</sup> and Neolithic Central Zagros<sup>8,118</sup> and CHG<sup>119</sup> ancestries, we used the  $f_4$ -  
803 statistic of the form  $f_4(\text{Saudi cluster}, \text{Yoruba.DG}^{\text{34,120}}; \text{Ganj Dareh N/CHG}, \text{Natufian EpiP})$ .

804 We employed replacement qpAdm <sup>112</sup> with parameters *allsnps=TRUE* and *fudge\_twice=TRUE* to  
805 model the ancestry for each of the Saudi clusters as it is partitioned in ancient groups across  
806 four broad and approximate periods. For each period we kept the following core fixed right  
807 group set of populations: Mbuti.DG <sup>121</sup>, Papuan.DG <sup>120</sup>, Russia\_Ust\_Ishim.DG <sup>122,123</sup>,  
808 Russia\_MA1\_HG.SG <sup>124</sup>, Russia\_Kostenki14 <sup>125</sup>, WHG <sup>10,119,125,126</sup>, CHG <sup>119</sup>, EHG <sup>127</sup>,  
809 Turkey\_Epipaleolithic <sup>128</sup>, Iran\_GanjDareh\_N <sup>8</sup>, and ISR\_Natufian\_EpiP <sup>8</sup>. For each of the four-  
810 time bins we modeled one to five source models, cycling through the source populations to  
811 form each qpAdm model. In evaluating the qpAdm models, we preferentially selected the  
812 model with the least number of sources with the largest p-value (with a plausibility threshold  
813 cut-off of 0.01 and admixture weights between 0 and 1), iteratively evaluating more complex  
814 models.

815 **Pre-Pottery Neolithic to Neolithic Sources:** Italy\_Sardinia\_N <sup>38,40</sup>, Levant\_PPN <sup>8,129</sup>,  
816 Mesopotamia\_PPN <sup>130</sup>, Anatolia\_Marmara\_Barcin\_N <sup>127,129</sup>, Turkey\_Catalhoyuk\_N\_Ceramic.SG  
817 <sup>131</sup>, MAR\_Taforalt\_EpiP <sup>49</sup>, Armenia\_Aknashen\_N <sup>129</sup>, Armenia\_MasisBlur\_N . **Chalcolithic to**  
818 **Bronze Age Sources:** Iran\_C\_SehGabi <sup>8</sup>, Turkey\_TellKurdu\_EC, Turkey\_C <sup>8,129</sup>, Israel\_C <sup>132</sup>,  
819 Armenia\_C <sup>8</sup>, Steppe\_Eneolithic <sup>133</sup>, and Ethiopia\_4500BP <sup>54,113</sup>. **Bronze Age Sources:**  
820 Ethiopia\_4500BP, Italy\_Sardinia\_EBA <sup>38,40</sup>, Mesopotamia\_LBA <sup>129</sup>, Israel\_MLBA <sup>134</sup>, Jordan\_LBA  
821 <sup>134</sup>, Lebanon\_MBA.SG <sup>135</sup>, Turkey\_EBA <sup>129</sup>, Armenia\_EBA\_KuraAraxes <sup>129,133</sup>, Armenia\_MBA <sup>8,129</sup>,  
822 and Germany\_BellBeaker <sup>127,136</sup>. **Bronze Age to Historical Sources:** Ethiopia\_4500BP,  
823 Italy\_Sardinia\_LBA <sup>40</sup>, Germany\_BellBeaker <sup>127,136</sup>, Iran\_Hasanlu\_IA <sup>129</sup>, Turkey\_IA <sup>129</sup>,  
824 Egypt\_ThirdIntermediatePeriod, AS\_EMT <sup>12</sup>, MH2\_MH1\_LT <sup>12</sup>, MH3\_LT <sup>12</sup>, Hungary\_IA\_LaTene  
825 <sup>137</sup>, Israel\_Ashkelon\_IA2 <sup>138</sup>, and Jordan\_LBA\_IA <sup>134</sup>.

826 We sought to date the formation of the plausible qpAdm models with DATES v4010 <sup>46</sup> using  
827 parameters *binsize: 0.001*, *maxdis: 1.0*, *qbin: 10*, *runfit: YES*, *qbin: 10*, *runfit: YES*, *afffit: YES*,  
828 *loalfit: 0.45*, *samecoeffs: NO*, and *jackknife: YES*. For each of the plausible qpAdm models, we  
829 ran the combination of sources through DATES, evaluating models with a normalized root mean  
830 standard deviation of < 0.7 & Z-score > 2 as well fitting. To obtain admixture dates in calendar  
831 years we used a generation time of 28 years <sup>139</sup>.

832 **Runs of homozygosity.** ROH are continuous segments of homozygous genotypes inherited from  
833 common ancestor <sup>56</sup>. Following Choudhury et al., <sup>140</sup> we used PLINK function *--option-homozyg*  
834 to identify runs of homozygosity (ROH) using the following parameters: we considered at least  
835 100 SNPs for ROH, with a total length  $\geq$  100 kilobases and at least one SNP per 50 kb on  
836 average; we set a scanning window to contain 100 SNPs, allowed 1 heterozygous call and 5  
837 missing calls per scanning window. We used three component Gaussian mixture model from the  
838 Mclust package (v.6.1) in R <sup>106</sup> following Pemberton et al., <sup>57</sup> to classify the ROHs into short,  
839 intermediate and long sizes.

840 **Demographic history.** Utilizing the phased WGS data, we estimated effective population sizes at  
841 different time points within the Saudi sub-clusters using RELATE v1.1.9 <sup>141</sup>. We used the  
842 *RelateFileFormats* in the *Relate* package to convert files from VCF format into haps/sample file  
843 format. For ancestral allele flipping, we provided RELATE with the human ancestor sequences  
844 release 107. We computed the genealogical trees using the parameters *-m 1.25e-8 -N 30,000*  
845 and subsequently used the *EstimatePopulationSize.sh* script provided with the *Relate* package  
846 to estimate the effective population sizes.

#### 847 **Evaluating imputation accuracy of Saudi genotypes**

848 We evaluated the impact of Saudi's population demographics on the imputation accuracy of its  
849 haplotypes, using the Trans-Omics for Precision Medicine (TOPMed) panel <sup>61</sup>. We compared the  
850 imputation accuracy between the Saudi and the Europeans in UK10K dataset <sup>62</sup>. For this  
851 comparison, we first selected the variants that are present in both dataset and then subsetted  
852 the Europeans to 3,252 individuals in order to match the same number of individuals as our  
853 Saudi data.

#### 854 **Enrichment of functionally deleterious alleles**

855 We compared the allelic architecture between Saudi Arabian and the gnomAD v4 <sup>64,65</sup>  
856 African/African American (gnomAD-AFR), non-Finnish European (gnomAD-EUR) and Middle  
857 Eastern (gnomAD-MID) populations. To check for potential enrichment or purging of deleterious

858 alleles in the Saudi, we computed the ratio of the proportional site frequency spectra for the  
859 deleterious alleles in Saudi to gnomAD-AFR or gnomAD-EUR, and contrasted it to the same ratio  
860 based on neutral or benign alleles. Utilizing the gnomAD exomes, which have a larger number of  
861 Middle Easterners compared to the genomes, we also made comparisons between the Middle  
862 Easterns and gnomAD-AFR and gnomAD-EUR. Significance differences in the ratios between  
863 variants functional classes were tested through bootstrapping.

864 For every comparison between populations or subpopulations, we used Hypergeometric (v  
865 3.6.2) distribution in R <sup>106</sup> to downsample both populations to equal sample sizes. All exome  
866 comparisons were downsampled to gnomAD-MID sample size. To account for technical  
867 differences in data generation of WGS call sets between gnomAD and Saudi data, we used the  
868 proportions of variants from the normalized allele frequency spectra rather than number of  
869 variants to compare the ratio between the Saudi and the gnomAD populations at a given allele  
870 count or frequency bin. However, when comparisons were made between two gnomAD  
871 populations or between two Saudi subpopulations, the actual number of variants were used.  
872

### 873 **DATA AVAILABILITY**

874 In compliance with Saudi privacy legislation and the protection of human subject confidentiality,  
875 the sharing of raw genotyping and clinical data is restricted. Access to this data requires prior  
876 approval from the Saudi National Bioethics Committee. The Saudi variants discovered through  
877 WGS and their estimated allele frequencies are deposited in the Figshare repository and can be  
878 accessed through this link: <https://doi.org/10.6084/m9.figshare.28059686.v1>.

879

### 880 **ACKNOWLEDGMENTS**

881 We are very grateful to the study participants who donated the samples used in the study. We  
882 thank Michael Campbell and Arun Durvasula for providing feedback on the drafted manuscript,

883 and Rui Leite Portela Martiniano for providing us with the Bahrain ancient DNA data. This work  
884 was supported by National Institute of General Medical Sciences (NIGMS) of the National  
885 Institute of Health under award number R35GM142783 (to C.W.K.C.). Computation for this work  
886 was supported by University of Southern California’s Center for Advanced Research Computing  
887 (<https://www.carc.usc.edu>).

888

## 889 **AUTHOR CONTRIBUTIONS**

890 D.K.M. and C. W. K. C. conceived and designed the study. C. W. K. C., S.M., and M.A. acquired  
891 funding for the data generation and analysis in this study. M.A. performed sample acquisition  
892 and data generation. D.K.M. and M.P.W. analyzed the data. C.D.H. provided analysis tools.  
893 D.K.M., M.P.W., C.D.H. and C.W.K.C. interpreted the results. D.K.M., M.P.W. and C.W.K.C. wrote  
894 the manuscript with input from all co-authors.

895

## 896 **DECLARATION OF INTERESTS**

897 The authors declare no competing interests.

898

## 899 **SUPPLEMENTAL INFORMATION**

900 Word document: Supplementary Figures S1 – 17

901 Word document: Supplementary Tables S1, 2, 3 and 18

902 Excel spreadsheet: Supplementary Tables S4 – 17

903

904 **REFERENCES**

- 905 1. Armitage, S.J., Jasim, S.A., Marks, A.E., Parker, A.G., Usik, V.I., and Uerpmann H-P (2011). The  
906 Southern Route“Out of Africa”:Evidence for an Early Expansionof Modern Humans into Arabia.  
907 *Science* (1979) *331*, 453–456. <https://doi.org/10.1594/PANGAEA.755114>.
- 908 2. Henn, B.M., Cavalli-Sforza, L.L., and Feldman, M.W. (2012). The great human expansion. Preprint,  
909 <https://doi.org/10.1073/pnas.1212380109> <https://doi.org/10.1073/pnas.1212380109>.
- 910 3. Rodriguez-Flores, J.L., Fakhro, K., Agosto-Perez, F., Ramstetter, M.D., Arbiza, L., Vincent, T.L.,  
911 Robay, A., Malek, J.A., Suhre, K., Chouchane, L., et al. (2016). Indigenous Arabs are descendants of  
912 the earliest split from ancient Eurasian populations. *Genome Res* *26*, 151–162.  
913 <https://doi.org/10.1101/gr.191478.115>.
- 914 4. Groucutt, H.S., Grün, R., Zalmout, I.A.S., Drake, N.A., Armitage, S.J., Candy, I., Clark-Wilson, R.,  
915 Louys, J., Breeze, P.S., Duval, M., et al. (2018). Homo sapiens in Arabia by 85,000 years ago. *Nat*  
916 *Ecol Evol* *2*, 800–809. <https://doi.org/10.1038/s41559-018-0518-2>.
- 917 5. Fernandes, V., Alshamali, F., Alves, M., Costa, M.D., Pereira, J.B., Silva, N.M., Cherni, L., Harich, N.,  
918 Cerny, V., Soares, P., et al. (2012). The Arabian cradle: Mitochondrial relicts of the first steps along  
919 the Southern route out of Africa. *Am J Hum Genet* *90*, 347–355.  
920 <https://doi.org/10.1016/j.ajhg.2011.12.010>.
- 921 6. Almarri, M.A., Haber, M., Lootah, R.A., Hallast, P., Al Turki, S., Martin, H.C., Xue, Y., and Tyler-  
922 Smith, C. (2021). The genomic history of the Middle East. *Cell* *184*, 4612-4625.e14.  
923 <https://doi.org/10.1016/j.cell.2021.07.013>.
- 924 7. Ferreira, J.C., Alshamali, F., Montinaro, F., Cavadas, B., Torroni, A., Pereira, L., Raveane, A., and  
925 Fernandes, V. (2021). Projecting Ancient Ancestry in Modern-Day Arabians and Iranians: A Key  
926 Role of the Past Exposed Arabo-Persian Gulf on Human Migrations. *Genome Biol Evol* *13*.  
927 <https://doi.org/10.1093/gbe/evab194>.
- 928 8. Lazaridis, I., Nadel, D., Rollefson, G., Merrett, D.C., Rohland, N., Mallick, S., Fernandes, D., Novak,  
929 M., Gamarra, B., Sirak, K., et al. (2016). Genomic insights into the origin of farming in the ancient  
930 Near East. *Nature* *536*, 419–424. <https://doi.org/10.1038/nature19310>.
- 931 9. Sankararaman, S., Patterson, N., Li, H., Pääbo, S., and Reich, D. (2012). The Date of Interbreeding  
932 between Neandertals and Modern Humans. *PLoS Genet* *8*.  
933 <https://doi.org/10.1371/journal.pgen.1002947>.
- 934 10. Lazaridis, I., Patterson, N., Mittnik, A., Renaud, G., Mallick, S., Kirsanow, K., Sudmant, P.H.,  
935 Schraiber, J.G., Castellano, S., Lipson, M., et al. (2014). Ancient human genomes suggest three  
936 ancestral populations for present-day Europeans. *Nature* *513*, 409–413.  
937 <https://doi.org/10.1038/nature13673>.
- 938 11. Vallini, L., Zampieri, C., Shoaee, M.J., Bortolini, E., Marciani, G., Aneli, S., Pievani, T., Benazzi, S.,  
939 Barausse, A., Mezzavilla, M., et al. (2024). The Persian plateau served as hub for Homo sapiens

- 940 after the main out of Africa dispersal. *Nat Commun* 15. [https://doi.org/10.1038/s41467-024-](https://doi.org/10.1038/s41467-024-46161-7)  
941 46161-7.
- 942 12. Martiniano, R., Haber, M., Almarri, M.A., Mattiangeli, V., Kuijpers, M.C.M., Chamel, B., Breslin,  
943 E.M., Littleton, J., Almahari, S., Aloraifi, F., et al. (2024). Ancient genomes illuminate Eastern  
944 Arabian population history and adaptation against malaria. *Cell Genomics* 4.  
945 <https://doi.org/10.1016/j.xgen.2024.100507>.
- 946 13. Fernandes, V., Brucato, N., Ferreira, J.C., Pedro, N., Cavadas, B., Ricaut, F.X., Alshamali, F., and  
947 Pereira, L. (2019). Genome-Wide Characterization of Arabian Peninsula Populations: Shedding  
948 Light on the History of a Fundamental Bridge between Continents. *Mol Biol Evol* 36, 575–586.  
949 <https://doi.org/10.1093/molbev/msz005>.
- 950 14. Khayat, A.M., Alshareef, B.G., Alharbi, S.F., AlZahrani, M.M., Alshangity, B.A., and Tashkandi, N.F.  
951 (2024). Consanguineous Marriage and Its Association With Genetic Disorders in Saudi Arabia: A  
952 Review. *Cureus*. <https://doi.org/10.7759/cureus.53888>.
- 953 15. Tadmouri, G.O., Nair, P., Obeid, T., Al Ali, M.T., Al Khaja, N., and Hamamy, H.A. (2009).  
954 Consanguinity and reproductive health among Arabs. *Reprod Health* 6.  
955 <https://doi.org/10.1186/1742-4755-6-17>.
- 956 16. El-Hazmi, M.A.F., Al-Swailem, A.R., Warsy, A.S., Al-Swailem, A.M., Sulaimani, R., Al-Meshari, A.A.,  
957 El-Hazmi, F., and Arabia, S. (1995). Consanguinity among the Saudi Arabian population.
- 958 17. Ben Halim, N., Ben Alaya Bouafif, N., Romdhane, L., Kefi Ben Atig, R., Chouchane, I., Bouyacoub,  
959 Y., Arfa, I., Cherif, W., Nouira, S., Talmoudi, F., et al. (2013). Consanguinity, endogamy, and genetic  
960 disorders in Tunisia. *J Community Genet* 4, 273–284. [https://doi.org/10.1007/s12687-012-0128-](https://doi.org/10.1007/s12687-012-0128-7)  
961 7.
- 962 18. Bittles, A.H. (2008). A community genetics perspective on consanguineous marriage. Preprint,  
963 <https://doi.org/10.1159/000133304> <https://doi.org/10.1159/000133304>.
- 964 19. Alkuraya, F.S. (2014). Genetics and genomic medicine in Saudi Arabia. *Mol Genet Genomic Med* 2,  
965 369–378. <https://doi.org/10.1002/mgg3.97>.
- 966 20. Mineta, K., Goto, K., Gojobori, T., and Alkuraya, F.S. (2021). Population structure of indigenous  
967 inhabitants of Arabia. *PLoS Genet* 17, e1009210.  
968 <https://doi.org/10.1371/JOURNAL.PGEN.1009210>.
- 969 21. Temaj, G., Nuhii, N., and Sayer, J.A. (2022). The impact of consanguinity on human health and  
970 disease with an emphasis on rare diseases. *Journal of Rare Diseases* 1.  
971 <https://doi.org/10.1007/s44162-022-00004-5>.
- 972 22. Aleissa, M., Aloraini, T., Alsubaie, L.F., Hassoun, M., Abdulrahman, G., Swaid, A., Al Eyaid, W., Al  
973 Mutairi, F., Ababneh, F., Alfadhel, M., et al. (2022). Common disease-associated gene variants in a  
974 Saudi Arabian population. *Ann Saudi Med* 42, 29–35. [https://doi.org/10.5144/0256-](https://doi.org/10.5144/0256-4947.2022.29)  
975 4947.2022.29.
- 976 23. Delatycki, M.B., Alkuraya, F., Archibald, A., Castellani, C., Cornel, M., Grody, W.W., Henneman, L.,  
977 Ioannides, A.S., Kirk, E., Laing, N., et al. (2020). International perspectives on the implementation



- 978 of reproductive carrier screening. Preprint at John Wiley and Sons Ltd,  
979 <https://doi.org/10.1002/pd.5611> <https://doi.org/10.1002/pd.5611>.
- 980 24. Hedrick, P.W., and Garcia-Dorado, A. (2016). Understanding Inbreeding Depression, Purging, and  
981 Genetic Rescue. Preprint at Elsevier Ltd, <https://doi.org/10.1016/j.tree.2016.09.005>  
982 <https://doi.org/10.1016/j.tree.2016.09.005>.
- 983 25. Sahoo, S.A., Zaidi, A.A., Anagol, S., and Mathieson, I. (2021). Long Runs of Homozygosity Are  
984 Correlated with Marriage Preferences across Global Population Samples. *Hum Biol* 93, 201–216.  
985 <https://doi.org/10.1353/hub.2021.0011>.
- 986 26. Scott, E.M., Halees, A., Itan, Y., Spencer, E.G., He, Y., Azab, M.A., Gabriel, S.B., Belkadi, A., Boisson,  
987 B., Abel, L., et al. (2016). Characterization of greater middle eastern genetic variation for  
988 enhanced disease gene discovery. Preprint at Nature Research, <https://doi.org/10.1038/ng.3592>  
989 <https://doi.org/10.1038/ng.3592>.
- 990 27. Alsalem, A.B., Halees, A.S., Anazi, S., Alshamekh, S., and Alkuraya, F.S. (2013). Autozygome  
991 Sequencing Expands the Horizon of Human Knockout Research and Provides Novel Insights into  
992 Human Phenotypic Variation. *PLoS Genet* 9. <https://doi.org/10.1371/journal.pgen.1004030>.
- 993 28. Overall, A., Ahmad, M., and Nichols, R.A. (2002). The effect of reproductive compensation on  
994 recessive disorders within consanguineous human populations. *Heredity (Edinb)* 88, 474–479.  
995 <https://doi.org/10.1038/sj/hdy/6800090>.
- 996 29. Ober, C., Hyslop, T., and Hauck, W.W. (1999). Inbreeding Effects on Fertility in Humans: Evidence  
997 for Reproductive Compensation.
- 998 30. Lohmueller, K.E. (2014). The distribution of deleterious genetic variation in human populations.  
999 Preprint at Elsevier Ltd, <https://doi.org/10.1016/j.gde.2014.09.005>  
1000 <https://doi.org/10.1016/j.gde.2014.09.005>.
- 1001 31. Castellano, S., Parra, G., Sánchez-Quinto, F.A., Racimo, F., Kuhlwilm, M., Kircher, M., Sawyer, S., Fu,  
1002 Q., Heinze, A., Nickel, B., et al. (2014). Patterns of coding variation in the complete exomes of  
1003 three Neandertals. *Proc Natl Acad Sci U S A* 111, 6666–6671.  
1004 <https://doi.org/10.1073/pnas.1405138111>.
- 1005 32. Simons, Y.B., and Sella, G. (2016). The impact of recent population history on the deleterious  
1006 mutation load in humans and close evolutionary relatives. Preprint at Elsevier Ltd,  
1007 <https://doi.org/10.1016/j.gde.2016.09.006> <https://doi.org/10.1016/j.gde.2016.09.006>.
- 1008 33. Laurent, R., Gineau, L., Utge, J., Lafosse, S., Phoeng, C.L., Hegay, T., Olaso, R., Boland, A., Deleuze,  
1009 J.F., Toupance, B., et al. (2024). Measuring the Efficiency of Purging by non-random Mating in  
1010 Human Populations. *Mol Biol Evol* 41. <https://doi.org/10.1093/molbev/msae094>.
- 1011 34. Bergström, A., McCarthy, S.A., Hui, R., Almarri, M.A., Ayub, Q., Danecek, P., Chen, Y., Felkel, S.,  
1012 Hallast, P., Kamm, J., et al. (2020). Insights into human genetic variation and population history  
1013 from 929 diverse genomes. *Science (1979)* 367. <https://doi.org/10.1126/science.aay5012>.
- 1014 35. Elliott, K.S., Haber, M., Daggag, H., Busby, G.B., Sarwar, R., Kennet, D., Petraglia, M., Petherbridge,  
1015 L.J., Yavari, P., Heard-Bey, F.U., et al. (2022). Fine-Scale Genetic Structure in the United Arab

- 1016 Emirates Reflects Endogamous and Consanguineous Culture, Population History, and Geography.  
1017 *Mol Biol Evol* 39. <https://doi.org/10.1093/molbev/msac039>.
- 1018 36. Hunwick, J.O. (2006). Arab views of black Africans and slavery. In *West Africa, Islam, and the Arab*  
1019 *world*, pp. 75–90.
- 1020 37. Charati, H., and Ori, R.J. (2021). Patterns of Genetic Structure and Evidence of Gene Flow  
1021 between Arabian Peninsula and European Populations. *Am J Biomed Sci Res* 12, 285–291.  
1022 <https://doi.org/10.34297/ajbsr.2021.12.001759>.
- 1023 38. Marcus, J.H., Posth, C., Ringbauer, H., Lai, L., Skeates, R., Sidore, C., Beckett, J., Furtwängler, A.,  
1024 Olivieri, A., Chiang, C.W.K., et al. (2020). Genetic history from the Middle Neolithic to present on  
1025 the Mediterranean island of Sardinia. *Nat Commun* 11. [https://doi.org/10.1038/s41467-020-](https://doi.org/10.1038/s41467-020-14523-6)  
1026 [14523-6](https://doi.org/10.1038/s41467-020-14523-6).
- 1027 39. Chiang, C.W.K., Marcus, J.H., Sidore, C., Biddanda, A., Al-Asadi, H., Zoledziewska, M., Pitzalis, M.,  
1028 Busonero, F., Maschio, A., Pistis, G., et al. (2018). Genomic history of the Sardinian population.  
1029 *Nat Genet* 50, 1426–1434. <https://doi.org/10.1038/s41588-018-0215-8>.
- 1030 40. Fernandes, D.M., Mittnik, A., Olalde, I., Lazaridis, I., Cheronet, O., Rohland, N., Mallick, S.,  
1031 Bernardos, R., Broomandkshobacht, N., Carlsson, J., et al. (2020). The spread of steppe and  
1032 Iranian-related ancestry in the islands of the western Mediterranean. *Nat Ecol Evol* 4, 334–345.  
1033 <https://doi.org/10.1038/s41559-020-1102-0>.
- 1034 41. Razali, R.M., Rodriguez-Flores, J., Ghorbani, M., Naeem, H., Aamer, W., Aliyev, E., Jubran, A.,  
1035 Ismail, S.I., Al-Muftah, W., Badji, R., et al. (2021). Thousands of Qatari genomes inform human  
1036 migration history and improve imputation of Arab haplotypes. *Nat Commun* 12.  
1037 <https://doi.org/10.1038/s41467-021-25287-y>.
- 1038 42. Lebling, R.W. (2009). “The Saracens of St. Tropez.” *Saudi Aramco World*.
- 1039 43. Mallick, S., Micco, A., Mah, M., Ringbauer, H., Lazaridis, I., Olalde, I., Patterson, N., and Reich, D.  
1040 (2024). The Allen Ancient DNA Resource (AADR) a curated compendium of ancient human  
1041 genomes. *Sci Data* 11. <https://doi.org/10.1038/s41597-024-03031-7>.
- 1042 44. Pagani, L., and Pagani, L. & C.I. (2019). What is Africa? A human perspective in Modern human  
1043 origins and dispersal. In *Morden human origins and dispersal*, H. Katerina and J. Gerhard, eds.  
1044 (Kerns Verlag), pp. 15–24.
- 1045 45. Yang, M.A., and Fu, Q. (2018). Insights into Modern Human Prehistory Using Ancient Genomes.  
1046 Preprint at Elsevier Ltd, <https://doi.org/10.1016/j.tig.2017.11.008>  
1047 <https://doi.org/10.1016/j.tig.2017.11.008>.
- 1048 46. Chintalapati, M., Patterson, N., and Moorjani, P. (2022). The spatiotemporal patterns of major  
1049 human admixture events during the European Holocene. *Elife* 11.  
1050 <https://doi.org/10.7554/ELIFE.77625>.
- 1051 47. Antonio, M.L., Weiß, C.L., Gao, Z., Sawyer, S., Oberreiter, V., Moots, H.M., Spence, J.P., Cheronet,  
1052 O., Zagorc, B., Praxmarer, E., et al. (2024). Stable population structure in Europe since the Iron  
1053 Age, despite high mobility. *Elife* 13. <https://doi.org/10.7554/ELIFE.79714>.

- 1054 48. Rodríguez-Varela, R., Günther, T., Krzewińska, M., Storå, J., Gillingwater, T.H., MacCallum, M.,  
1055 Arsuaga, J.L., Dobney, K., Valdiosera, C., Jakobsson, M., et al. (2017). Genomic Analyses of Pre-  
1056 European Conquest Human Remains from the Canary Islands Reveal Close Affinity to Modern  
1057 North Africans. *Curr Biol* 27, 3396-3402.e5. <https://doi.org/10.1016/J.CUB.2017.09.059>.
- 1058 49. Van De Loosdrecht, M., Bouzouggar, A., Humphrey, L., Posth, C., Barton, N., Aximu-Petri, A.,  
1059 Nickel, B., Nagel, S., Talbi, E.H., Abdeljalil, M., et al. (2018). Pleistocene North African genomes  
1060 link Near Eastern and sub-Saharan African human populations. *Science* (1979), 548–552.
- 1061 50. Moots, H.M., Antonio, M., Sawyer, S., Spence, J.P., Oberreiter, V., Weiß, C.L., Lucci, M., Cherifi,  
1062 Y.M.S., La Pastina, F., Genchi, F., et al. (2023). A genetic history of continuity and mobility in the  
1063 Iron Age central Mediterranean. *Nature Ecology & Evolution* 2023 7:9 7, 1515–1524.  
1064 <https://doi.org/10.1038/s41559-023-02143-4>.
- 1065 51. Schuenemann, V.J., Peltzer, A., Welte, B., Van Pelt, W.P., Molak, M., Wang, C.C., Furtwängler, A.,  
1066 Urban, C., Reiter, E., Nieselt, K., et al. (2017). Ancient Egyptian mummy genomes suggest an  
1067 increase of Sub-Saharan African ancestry in post-Roman periods. *Nat Commun* 8.  
1068 <https://doi.org/10.1038/ncomms15694>.
- 1069 52. Prendergast, M.E., Lipson, M., Sawchuk, E.A., Olalde, I., Ogola, C.A., Rohland, N., Sirak, K.A.,  
1070 Adamski, N., Bernardos, R., Broomandkshobacht, N., et al. (2019). Ancient DNA reveals a  
1071 multistep spread of the first herders into sub-Saharan Africa. *Science* (1979) 364.
- 1072 53. Sirak, K.A., Fernandes, D.M., Lipson, M., Mallick, S., Mah, M., Olalde, I., Ringbauer, H., Rohland,  
1073 N., Hadden, C.S., Harney, É., et al. (2021). Social stratification without genetic differentiation at  
1074 the site of Kulubnarti in Christian Period Nubia. *Nature Communications* 2021 12:1 12, 1–14.  
1075 <https://doi.org/10.1038/s41467-021-27356-8>.
- 1076 54. Lipson, M., Sawchuk, E.A., Thompson, J.C., Oppenheimer, J., Tryon, C.A., Ranhorn, K.L., de Luna,  
1077 K.M., Sirak, K.A., Olalde, I., Ambrose, S.H., et al. (2022). Ancient DNA and deep population  
1078 structure in sub-Saharan African foragers. *Nature* 2022 603:7900 603, 290–296.  
1079 <https://doi.org/10.1038/s41586-022-04430-9>.
- 1080 55. Schlebusch, C.M., Malmström, H., Günther, T., Sjödin, P., Coutinho, A., Edlund, H., Munters, A.R.,  
1081 Vicente, M., Steyn, M., Soodyall, H., et al. (2017). Southern African ancient genomes estimate  
1082 modern human divergence to 350,000 to 260,000 years ago. *Science* (1979) 358, 652–655.
- 1083 56. Ceballos, F.C., Joshi, P.K., Clark, D.W., Ramsay, M., and Wilson, J.F. (2018). Runs of homozygosity:  
1084 Windows into population history and trait architecture. Preprint at Nature Publishing Group,  
1085 <https://doi.org/10.1038/nrg.2017.109> <https://doi.org/10.1038/nrg.2017.109>.
- 1086 57. Pemberton, T.J., Absher, D., Feldman, M.W., Myers, R.M., Rosenberg, N.A., and Li, J.Z. (2012).  
1087 Genomic patterns of homozygosity in worldwide human populations. *Am J Hum Genet* 91, 275–  
1088 292. <https://doi.org/10.1016/j.ajhg.2012.06.014>.
- 1089 58. Thompson, E.A. (2013). Identity by descent: Variation in meiosis, across genomes, and in  
1090 populations. Preprint, <https://doi.org/10.1534/genetics.112.148825>  
1091 <https://doi.org/10.1534/genetics.112.148825>.

- 1092 59. Ceballos, F.C., Gürün, K., Altınışık, N.E., Gemici, H.C., Karamurat, C., Koptekin, D., Vural, K.B.,  
1093 Mapelli, I., Sağlıcan, E., Sürer, E., et al. (2021). Human inbreeding has decreased in time through  
1094 the Holocene. *Current Biology* 31, 3925-3934.e8. <https://doi.org/10.1016/j.cub.2021.06.027>.
- 1095 60. Petraglia, M.D., Groucutt, H.S., Guagnin, M., Breeze, P.S., and Boivin, N. (2020). Human responses  
1096 to climate and ecosystem change in ancient Arabia. *Proceedings of the National Academy of  
1097 Sciences* 117, 8263–8270. <https://doi.org/10.1073/pnas.1920211117/-/DCSupplemental>.
- 1098 61. Taliun, D., Harris, D.N., Kessler, M.D., Carlson, J., Szpiech, Z.A., Torres, R., Taliun, S.A.G., Corvelo,  
1099 A., Gogarten, S.M., Kang, H.M., et al. (2021). Sequencing of 53,831 diverse genomes from the  
1100 NHLBI TOPMed Program. *Nature* 590, 290–299. <https://doi.org/10.1038/s41586-021-03205-y>.
- 1101 62. Walter, K., Min, J.L., Huang, J., Crooks, L., Memari, Y., McCarthy, S., Perry, J.R.B., Xu, C., Futema,  
1102 M., Lawson, D., et al. (2015). The UK10K project identifies rare variants in health and disease.  
1103 *Nature* 526, 82–89. <https://doi.org/10.1038/nature14962>.
- 1104 63. Cahoon, J.L., Rui, X., Tang, E., Simons, C., Langie, J., Chen, M., Lo, Y.-C., and Chiang, C.W.K. (2024).  
1105 Imputation Accuracy Across Global Human Populations. *Am J Hum Genet* 111, P979-989.  
1106 <https://doi.org/10.1101/2023.05.22.541241>.
- 1107 64. Chen, S., Francioli, L.C., Goodrich, J.K., Collins, R.L., Kanai, M., Wang, Q., Alföldi, J., Watts, N.A.,  
1108 Vittal, C., Gauthier, L.D., et al. (2023). A genomic mutational constraint map using variation in  
1109 76,156 human genomes. *Nature* 625, 92–100. <https://doi.org/10.1038/s41586-023-06045-0>.
- 1110 65. Karczewski, K.J., Francioli, L.C., Tiao, G., Cummings, B.B., Alföldi, J., Wang, Q., Collins, R.L.,  
1111 Laricchia, K.M., Ganna, A., Birnbaum, D.P., et al. (2020). The mutational constraint spectrum  
1112 quantified from variation in 141,456 humans. *Nature* 581, 434–443.  
1113 <https://doi.org/10.1038/s41586-020-2308-7>.
- 1114 66. Ljungdahl, A., Kohani, S., Page, N.F., Wells, E.S., Wigdor, E.M., Dong, S., and Sanders, S.J. (2023).  
1115 AlphaMissense is better correlated with functional assays of missense impact than earlier  
1116 prediction algorithms. *bioRxiv* , 562294. <https://doi.org/10.1101/2023.10.24.562294>.
- 1117 67. Benegas, G.I., Singh Batra, S.I., Song, Y.S., and Edited by Kathryn Roeder, I. (2023). BIOPHYSICS  
1118 AND COMPUTATIONAL BIOLOGY OPEN ACCESS DNA language models are powerful predictors of  
1119 genome-wide variant effects. 120. <https://doi.org/10.1073/pnas>.
- 1120 68. Lohmueller, K.E., Indap, A.R., Schmidt, S., Boyko, A.R., Ryan, D., Hubisz, M.J., Sninsky, J.J., White,  
1121 T.J., Sunyaev, S.R., Nielsen, R., et al. (2008). Proportionally More Deleterious Genetic Variation In  
1122 European than in African Populations. *Nature* 451, 994–997.  
1123 <https://doi.org/10.1038/nature06611>.Proportionally.
- 1124 69. Henn, B.M., Botigué, L.R., Peischl, S., Dupanloup, I., Lipatov, M., Maples, B.K., Martin, A.R.,  
1125 Musharoff, S., Cann, H., Snyder, M.P., et al. (2016). Distance from sub-Saharan Africa predicts  
1126 mutational load in diverse human genomes. *Proceedings of the National Academy of Sciences*  
1127 113, E440–E449. <https://doi.org/10.1073/pnas.1510805112>.

- 1128 70. Simons, Y.B., Turchin, M.C., Pritchard, J.K., and Sella, G. (2014). The deleterious mutation load is  
1129 insensitive to recent population history. *Nat Genet* 46, 220–224.  
1130 <https://doi.org/10.1038/ng.2896>.
- 1131 71. Do, R., Balick, D., Li, H., Adzhubei, I., Sunyaev, S., and Reich, D. (2015). No evidence that selection  
1132 has been less effective at removing deleterious mutations in Europeans than in Africans. *Nat*  
1133 *Genet* 47, 126–131. <https://doi.org/10.1038/ng.3186>.
- 1134 72. Lim, E.T., Würtz, P., Havulinna, A.S., Palta, P., Tukiainen, T., Rehnström, K., Esko, T., Mägi, R.,  
1135 Inouye, M., Lappalainen, T., et al. (2014). Distribution and Medical Impact of Loss-of-Function  
1136 Variants in the Finnish Founder Population. *PLoS Genet* 10.  
1137 <https://doi.org/10.1371/journal.pgen.1004494>.
- 1138 73. Pedersen, C.E.T., Lohmueller, K.E., Grarup, N., Bjerregaard, P., Hansen, T., Siegismund, H.R.,  
1139 Moltke, I., and Albrechtsen, A. (2017). The effect of an extreme and prolonged population  
1140 bottleneck on patterns of deleterious variation: Insights from the Greenlandic Inuit. *Genetics* 205,  
1141 787–801. <https://doi.org/10.1534/genetics.116.193821>.
- 1142 74. Locke, A.E., Steinberg, K.M., Chiang, C.W.K., Service, S.K., Havulinna, A.S., Stell, L., Pirinen, M.,  
1143 Abel, H.J., Chiang, C.C., Fulton, R.S., et al. (2019). Exome sequencing of Finnish isolates enhances  
1144 rare-variant association power. *Nature* 572, 323–328. [https://doi.org/10.1038/s41586-019-1457-](https://doi.org/10.1038/s41586-019-1457-z)  
1145 [z](https://doi.org/10.1038/s41586-019-1457-z).
- 1146 75. Subramanian, S. (2016). Europeans have a higher proportion of high-frequency deleterious  
1147 variants than Africans. *Hum Genet* 135, 1–7. <https://doi.org/10.1007/s00439-015-1604-z>.
- 1148 76. Eaaswarkhanth, M., Pathak, A.K., Ongaro, L., Montinaro, F., Prashantha Hebbar, •, Osama  
1149 Alsmadi, •, Metspalu, M., Al-Mulla, F., Thangavel, •, and Thanaraj, A. (2021). Unraveling a fine-  
1150 scale high genetic heterogeneity and recent continental connections of an Arabian Peninsula  
1151 population. *European Journal of Human Genetics* 30, 307–319. [https://doi.org/10.1038/s41431-](https://doi.org/10.1038/s41431-021-00861-6)  
1152 [021-00861-6](https://doi.org/10.1038/s41431-021-00861-6).
- 1153 77. Khubrani, Y.M., Wetton, J.H., and Jobling, M.A. (2018). Extensive geographical and social structure  
1154 in the paternal lineages of Saudi Arabia revealed by analysis of 27 Y-STRs. *Forensic Sci Int Genet*  
1155 33, 98–105. <https://doi.org/10.1016/j.fsigen.2017.11.015>.
- 1156 78. Moreno-Estrada, A., Gravel, S., Zakharia, F., McCauley, J.L., Byrnes, J.K., Gignoux, C.R., Ortiz-Tello,  
1157 P.A., Martínez, R.J., Hedges, D.J., Morris, R.W., et al. (2013). Reconstructing the Population  
1158 Genetic History of the Caribbean. *PLoS Genet* 9. <https://doi.org/10.1371/journal.pgen.1003925>.
- 1159 79. Browning, S.R., Grinde, K., Plantinga, A., Gogarten, S.M., Stilp, A.M., Kaplan, R.C., Avilés-Santa,  
1160 M.L., Browning, B.L., and Laurie, C.C. (2016). Local ancestry inference in a large US-based  
1161 Hispanic/Latino study: Hispanic community health study/study of Latinos (HCHS/SOL). *G3: Genes,*  
1162 *Genomes, Genetics* 6, 1525–1534. <https://doi.org/10.1534/g3.116.028779>.
- 1163 80. El-Mouzan, M.I., Al-Salloum, A.A., Al-Herbish, A.S., Qurachi, M.M., and Al-Omar, A.A. (2007).  
1164 Regional variations in the prevalence of consanguinity in Saudi Arabia. *Saudi Med J* 28.

- 1165 81. Hellenthal, G., Busby, G.B.J., Band, G., Wilson, J.F., Capelli, C., Falush, D., and Myers, S. (2014). A  
1166 genetic atlas of human admixture history. *Science* (1979) *343*, 747–751.  
1167 <https://doi.org/10.1126/science.1243518>.
- 1168 82. Miran, J. (2022). *Red Sea Slave Trade* (Oxford University Press)  
1169 <https://doi.org/10.1093/acrefore/9780190277734.013.868>.
- 1170 83. Casals, F., Hodgkinson, A., Hussin, J., Idaghdour, Y., Bruat, V., de Maillard, T., Grenier, J.C., Gbeha,  
1171 E., Hamdan, F.F., Girard, S., et al. (2013). Whole-Exome Sequencing Reveals a Rapid Change in the  
1172 Frequency of Rare Functional Variants in a Founding Population of Humans. *PLoS Genet* *9*.  
1173 <https://doi.org/10.1371/journal.pgen.1003815>.
- 1174 84. Warsy, A.S., Al-Jaser, M.H., Albass, A., Al-Daihan, S., and Alanazi, M. (2014). Is consanguinity  
1175 prevalence decreasing in Saudis?: A study in two generations. *Afr Health Sci* *14*, 314–321.  
1176 <https://doi.org/10.4314/ahs.v14i2.5>.
- 1177 85. Albanghali, M.A. (2023). Prevalence of Consanguineous Marriage among Saudi Citizens of Albaha,  
1178 a Cross-Sectional Study. *Int J Environ Res Public Health* *20*.  
1179 <https://doi.org/10.3390/ijerph20043767>.
- 1180 86. Al-Gazali, L., Hamamy, H., and Al-Arrayad, S. (2006). Genetic disorders in the Arab world.
- 1181 87. Saffi, M., and Howard, N. (2015). Exploring the Effectiveness of Mandatory Premarital Screening  
1182 and Genetic Counselling Programmes for  $\beta$ -Thalassaemia in the Middle East: A Scoping Review.  
1183 Preprint at S. Karger AG, <https://doi.org/10.1159/000430837>  
1184 <https://doi.org/10.1159/000430837>.
- 1185 88. Tadmouri, G.O., Nair, P., Obeid, T., Al Ali, M.T., Al Khaja, N., and Hamamy, H.A. (2009).  
1186 Consanguinity and reproductive health among Arabs. *Reprod Health* *6*.  
1187 <https://doi.org/10.1186/1742-4755-6-17>.
- 1188 89. Elfatih, A., Saad, C., Ismail, S., Al-Muftah, W., Badji, R., Darwish, D., Fadl, T., Yasin, H., Ennaifar, M.,  
1189 Abdel-atif, R., et al. (2024). Analysis of 14,392 whole genomes reveals 3.5% of Qataris carry  
1190 medically actionable variants. *European Journal of Human Genetics*.  
1191 <https://doi.org/10.1038/s41431-024-01656-1>.
- 1192 90. Mbarek, H., Gandhi, G.D., Selvaraj, S., Al-Muftah, W., Badji, R., Al-Sarraj, Y., Saad, C., Darwish, D.,  
1193 Alvi, M., Fadl, T., et al. (2022). Qatar genome: Insights on genomics from the Middle East. *Hum*  
1194 *Mutat.* <https://doi.org/10.1002/HUMU.24336>.
- 1195 91. Thareja, G., Al-Sarraj, Y., Belkadi, A., Almotawa, M., Ismail, S., Al-Muftah, W., Badji, R., Mbarek, H.,  
1196 Darwish, D., Fadl, T., et al. (2021). Whole genome sequencing in the Middle Eastern Qatari  
1197 population identifies genetic associations with 45 clinically relevant traits. *Nature*  
1198 *Communications* *2021* *12*:1 *12*, 1–10. <https://doi.org/10.1038/s41467-021-21381-3>.
- 1199 92. Kousathanas, A., Pairo-Castineira, E., Rawlik, K., Stuckey, A., Odhams, C.A., Walker, S., Russell,  
1200 C.D., Malinauskas, T., Wu, Y., Millar, J., et al. (2022). Whole-genome sequencing reveals host  
1201 factors underlying critical COVID-19. *Nature* *607*, 97–103. [https://doi.org/10.1038/s41586-022-](https://doi.org/10.1038/s41586-022-04576-6)  
1202 [04576-6](https://doi.org/10.1038/s41586-022-04576-6).

- 1203 93. Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M.A.R., Bender, D., Maller, J., Sklar, P.,  
1204 de Bakker, P.I.W., Daly, M.J., et al. (2007). PLINK: A tool set for whole-genome association and  
1205 population-based linkage analyses. *Am J Hum Genet* 81, 559–575.  
1206 <https://doi.org/10.1086/519795>.
- 1207 94. Chang, C.C., Chow, C.C., Tellier, L.C., Vattikuti, S., Purcell, S.M., and Lee, J.J. (2015). Second-  
1208 generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience* 4, 7.  
1209 <https://doi.org/10.1186/s13742-015-0047-8>.
- 1210 95. Browning, B.L., Tian, X., Zhou, Y., and Browning, S.R. (2021). Fast two-stage phasing of large-scale  
1211 sequence data. <https://doi.org/10.1016/j.ajhg.2021.08.005>.
- 1212 96. Manichaikul, A., Mychaleckyj, J.C., Rich, S.S., Daly, K., Sale, M., and Chen, W.M. (2010). Robust  
1213 relationship inference in genome-wide association studies. *Bioinformatics* 26, 2867–2873.  
1214 <https://doi.org/10.1093/bioinformatics/btq559>.
- 1215 97. Fang, H., Hui, Q., Lynch, J., Honerlaw, J., Assimes, T.L., Huang, J., Vujkovic, M., Damrauer, S.M.,  
1216 Pyarajan, S., Gaziano, J.M., et al. (2019). Harmonizing Genetic Ancestry and Self-identified  
1217 Race/Ethnicity in Genome-wide Association Studies. *Am J Hum Genet* 105, 763–772.  
1218 <https://doi.org/10.1016/j.ajhg.2019.08.012>.
- 1219 98. Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows–Wheeler  
1220 transform. *Bioinformatics* 25, 1754–1760. <https://doi.org/10.1093/bioinformatics/btp324>.
- 1221 99. Mckenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., Garimella, K.,  
1222 Altshuler, D., Gabriel, S., Daly, M., et al. (2010). The Genome Analysis Toolkit: A MapReduce  
1223 framework for analyzing next-generation DNA sequencing data. *Genome Res* 20, 1297–1303.  
1224 <https://doi.org/10.1101/gr.107524.110>.
- 1225 100. Jónsson, H., Sulem, P., Kehr, B., Kristmundsdottir, S., Zink, F., Hjartarson, E., Hardarson, M.T.,  
1226 Hjorleifsson, K.E., Eggertsson, H.P., Gudjonsson, S.A., et al. (2017). Data Descriptor: Whole  
1227 genome characterization of sequence diversity of 15,220 Icelanders. *Sci Data* 4.  
1228 <https://doi.org/10.1038/sdata.2017.115>.
- 1229 101. Van der Auwera, G.A., Carneiro, M.O., Hartl, C., Poplin, R., del Angel, G., Levy-Moonshine, A.,  
1230 Jordan, T., Shakir, K., Roazen, D., Thibault, J., et al. (2013). From fastQ data to high-confidence  
1231 variant calls: The genome analysis toolkit best practices pipeline. *Curr Protoc Bioinformatics*.  
1232 <https://doi.org/10.1002/0471250953.bi1110s43>.
- 1233 102. Loh, P.R., Danecek, P., Palamara, P.F., Fuchsberger, C., Reshef, Y.A., Finucane, H.K., Schoenherr, S.,  
1234 Forer, L., McCarthy, S., Abecasis, G.R., et al. (2016). Reference-based phasing using the Haplotype  
1235 Reference Consortium panel. *Nat Genet* 48, 1443–1448. <https://doi.org/10.1038/ng.3679>.
- 1236 103. McLaren, W., Pritchard, B., Rios, D., Chen, Y., Flicek, P., and Cunningham, F. (2010). Deriving the  
1237 consequences of genomic variants with the Ensembl API and SNP Effect Predictor. *Bioinformatics*  
1238 26, 2069–2070. <https://doi.org/10.1093/bioinformatics/btq330>.

- 1239 104. Cheng, J., Novati, G., Pan, J., Bycroft, C., Žemgulyte, A., Applebaum, T., Pritzel, A., Wong, L.H.,  
1240 Zielinski, M., Sargeant, T., et al. (2023). Accurate proteome-wide missense variant effect  
1241 prediction with AlphaMissense. *Science* (1979) *381*. <https://doi.org/10.1126/science.adg7492>.
- 1242 105. McInnes, L., Healy, J., Saul, N., and Großberger, L. (2018). UMAP: Uniform Manifold  
1243 Approximation and Projection Software • Review • Repository • Archive.  
1244 <https://doi.org/10.21105/joss.00861>.
- 1245 106. R Core Team (2015). R: A language and environment for statistical computing. R Foundation for  
1246 Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.
- 1247 107. Rousseeuw, P.J. (1987). Silhouettes: a graphical aid to the interpretation and validation of cluster  
1248 analysis.
- 1249 108. Batool, F., and Hennig, C. (2019). Clustering with the Average Silhouette Width.
- 1250 109. Alexander, D.H., and Novembre, J. (2009). Fast Model-Based Estimation of Ancestry in Unrelated  
1251 Individuals. 1655–1664. <https://doi.org/10.1101/gr.094052.109.vidual>.
- 1252 110. Maier, R., and Patterson, N. (2024). admixtools: Inferring demographic history from genetic data.  
1253 R package version 2.0.4. Preprint.
- 1254 111. Patterson, N., Moorjani, P., Luo, Y., Mallick, S., Rohland, N., Zhan, Y., Genschoreck, T., Webster, T.,  
1255 and Reich, D. (2012). Ancient admixture in human history. *Genetics* *192*, 1065–1093.  
1256 <https://doi.org/10.1534/genetics.112.145037>.
- 1257 112. Maier, R., Flegontov, P., Flegontova, O., Işıldak, U., Changmai, P., and Reich, D. (2023). On the  
1258 limits of fitting complex models of population history to f-statistics. *Elife* *12*.  
1259 <https://doi.org/10.7554/eLife.85492>.
- 1260 113. Gallego Llorente, M., Jones, E.R., Eriksson, A., Siska, V., Arthur, K.W., Arthur, J.W., Curtis, M.C.,  
1261 Stock, J.T., Coltorti, M., Pieruccini, P., et al. (2015). Ancient Ethiopian genome reveals extensive  
1262 Eurasian admixture throughout the African continent. *Science* (1979) *350*, 820–822.
- 1263 114. Wang, K., Goldstein, S., Bleasdale, M., Clist, B., Clist, B., Bostoen, K., Bakwa-Lufu, P., Buck, L.T.,  
1264 Buck, L.T., Crowther, A., et al. (2020). Ancient genomes reveal complex patterns of population  
1265 movement, interaction, and replacement in sub-Saharan Africa. *Sci Adv* *6*, 183–195.
- 1266 115. Skoglund, P., Thompson, J.C., Prendergast, M.E., Mitnik, A., Sirak, K., Hajdinjak, M., Salie, T.,  
1267 Rohland, N., Mallick, S., Peltzer, A., et al. (2017). Reconstructing Prehistoric African Population  
1268 Structure. *Cell* *171*, 59-71.e21. <https://doi.org/10.1016/J.CELL.2017.08.049>.
- 1269 116. Jakobsson, M., Scholz, S.W., Scheet, P., Gibbs, J.R., VanLiere, J.M., Fung, H.C., Szpiech, Z.A.,  
1270 Degnan, J.H., Wang, K., Guerreiro, R., et al. (2008). Genotype, haplotype and copy-number  
1271 variation in worldwide human populations. *Nature* *2008* *451*:7181 *451*, 998–1003.  
1272 <https://doi.org/10.1038/nature06742>.
- 1273 117. Li, J.Z., Absher, D.M., Tang, H., Southwick, A.M., Casto, A.M., Ramachandran, S., Cann, H.M.,  
1274 Barsh, G.S., Feldman, M., Cavalli-Sforza, L.L., et al. (2008). Worldwide human relationships



- 1275 inferred from genome-wide patterns of variation. *Science* 319, 1100–1104.  
1276 <https://doi.org/10.1126/SCIENCE.1153717>.
- 1277 118. Narasimhan, V.M., Patterson, N., Moorjani, P., Rohland, N., Bernardos, R., Mallick, S., Lazaridis, I.,  
1278 Nakatsuka, N., Olalde, I., Lipson, M., et al. (2019). The formation of human populations in South  
1279 and Central Asia. *Science* (1979) 365.
- 1280 119. Jones, E.R., Gonzalez-Fortes, G., Connell, S., Siska, V., Eriksson, A., Martiniano, R., McLaughlin,  
1281 R.L., Gallego Llorente, M., Cassidy, L.M., Gamba, C., et al. (2015). Upper Palaeolithic genomes  
1282 reveal deep roots of modern Eurasians. *Nature Communications* 2015 6:1 6, 1–8.  
1283 <https://doi.org/10.1038/ncomms9912>.
- 1284 120. Skoglund, P., Mallick, S., Bortolini, M.C., Chennagiri, N., Hünemeier, T., Petzl-Erler, M.L., Salzano,  
1285 F.M., Patterson, N., and Reich, D. (2015). Genetic evidence for two founding populations of the  
1286 Americas. *Nature* 2015 525:7567 525, 104–108. <https://doi.org/10.1038/nature14895>.
- 1287 121. Mallick, S., Li, H., Lipson, M., Mathieson, I., Gymrek, M., Racimo, F., Zhao, M., Chennagiri, N.,  
1288 Nordenfelt, S., Tandon, A., et al. (2016). The Simons Genome Diversity Project: 300 genomes from  
1289 142 diverse populations. *Nature* 538, 201–206. <https://doi.org/10.1038/nature18964>.
- 1290 122. Fu, Q., Li, H., Moorjani, P., Jay, F., Slepchenko, S.M., Bondarev, A.A., Johnson, P.L.F., Aximu-Petri,  
1291 A., Prüfer, K., De Filippo, C., et al. (2014). Genome sequence of a 45,000-year-old modern human  
1292 from western Siberia. *Nature* 2014 514:7523 514, 445–449.  
1293 <https://doi.org/10.1038/nature13810>.
- 1294 123. Prüfer, K., De Filippo, C., Grote, S., Mafessoni, F., Korlević, P., Hajdinjak, M., Vernot, B., Skov, L.,  
1295 Hsieh, P., Peyrégne, S., et al. (2017). A high-coverage Neandertal genome from Vindija Cave in  
1296 Croatia. *Science* (1979) 358, 655–658.
- 1297 124. Raghavan, M., Skoglund, P., Graf, K.E., Metspalu, M., Albrechtsen, A., Moltke, I., Rasmussen, S.,  
1298 Stafford, T.W., Orlando, L., Metspalu, E., et al. (2013). Upper Palaeolithic Siberian genome reveals  
1299 dual ancestry of Native Americans. *Nature* 2013 505:7481 505, 87–91.  
1300 <https://doi.org/10.1038/nature12736>.
- 1301 125. Fu, Q., Posth, C., Hajdinjak, M., Petr, M., Mallick, S., Fernandes, D., Furtwängler, A., Haak, W.,  
1302 Meyer, M., Mittnik, A., et al. (2016). The genetic history of Ice Age Europe. *Nature* 2016 534:7606  
1303 534, 200–205. <https://doi.org/10.1038/nature17993>.
- 1304 126. Olalde, I., Allentoft, M.E., Sánchez-Quinto, F., Santpere, G., Chiang, C.W.K., DeGiorgio, M., Prado-  
1305 Martinez, J., Rodríguez, J.A., Rasmussen, S., Quilez, J., et al. (2014). Derived immune and ancestral  
1306 pigmentation alleles in a 7,000-year-old Mesolithic European. *Nature* 2014 507:7491 507, 225–  
1307 228. <https://doi.org/10.1038/nature12960>.
- 1308 127. Mathieson, I., Lazaridis, I., Rohland, N., Mallick, S., Patterson, N., Roodenberg, S.A., Harney, E.,  
1309 Stewardson, K., Fernandes, D., Novak, M., et al. (2015). Genome-wide patterns of selection in 230  
1310 ancient Eurasians. *Nature* 528, 499–503. <https://doi.org/10.1038/NATURE16152>.
- 1311 128. Feldman, M., Fernández-Domínguez, E., Reynolds, L., Baird, D., Pearson, J., Hershkovitz, I., May,  
1312 H., Goring-Morris, N., Benz, M., Gresky, J., et al. (2019). Late Pleistocene human genome suggests

- 1313 a local origin for the first farmers of central Anatolia. *Nature Communications* 2019 10:1 10, 1–10.  
1314 <https://doi.org/10.1038/s41467-019-09209-7>.
- 1315 129. Lazaridis, I., Alpaslan-Roodenberg, S., Acar, A., Açikkol, A., Agelarakis, A., Aghikyan, L., Akyüz, U.,  
1316 Andreeva, D., Andrijašević, G., Antonović, D., et al. (2022). The genetic history of the Southern  
1317 Arc: A bridge between West Asia and Europe. *Science* (1979) 377.
- 1318 130. Lazaridis, I., Alpaslan-Roodenberg, S., Acar, A., Açikkol, A., Agelarakis, A., Aghikyan, L., Akyüz, U.,  
1319 Andreeva, D., Andrijašević, G., Antonović, D., et al. (2022). Ancient DNA from Mesopotamia  
1320 suggests distinct Pre-Pottery and Pottery Neolithic migrations into Anatolia. *Science* (1979) 377,  
1321 982–987.
- 1322 131. Yaka, R., Mapelli, I., Kaptan, D., Doğu, A., Chyleński, M., Erdal, Ö.D., Koptekin, D., Vural, K.B.,  
1323 Bayliss, A., Mazzucato, C., et al. (2021). Variable kinship patterns in Neolithic Anatolia revealed by  
1324 ancient genomes. *Curr Biol* 31, 2455–2468.e18. <https://doi.org/10.1016/J.CUB.2021.03.050>.
- 1325 132. Harney, É., May, H., Shalem, D., Rohland, N., Mallick, S., Lazaridis, I., Sarig, R., Stewardson, K.,  
1326 Nordenfelt, S., Patterson, N., et al. (2018). Publisher Correction: Ancient DNA from Chalcolithic  
1327 Israel reveals the role of population mixture in cultural transformation. *Nature Communications*  
1328 2018 9:1 9, 1–1. <https://doi.org/10.1038/s41467-018-06484-8>.
- 1329 133. Wang, C.C., Reinhold, S., Kalmykov, A., Wissgott, A., Brandt, G., Jeong, C., Cheronet, O., Ferry, M.,  
1330 Harney, E., Keating, D., et al. (2019). Ancient human genome-wide data from a 3000-year interval  
1331 in the Caucasus corresponds with eco-geographic regions. *Nature Communications* 2019 10:1 10,  
1332 1–13. <https://doi.org/10.1038/s41467-018-08220-8>.
- 1333 134. Agranat-Tamir, L., Waldman, S., Martin, M.A.S., Gokhman, D., Mishol, N., Eshel, T., Cheronet, O.,  
1334 Rohland, N., Mallick, S., Adamski, N., et al. (2020). The Genomic History of the Bronze Age  
1335 Southern Levant. *Cell* 181, 1146–1157.e11. <https://doi.org/10.1016/j.cell.2020.04.024>.
- 1336 135. Haber, M., Doumet-Serhal, C., Scheib, C., Xue, Y., Danecek, P., Mezzavilla, M., Youhanna, S.,  
1337 Martiniano, R., Prado-Martinez, J., Szpak, M., et al. (2017). Continuity and Admixture in the Last  
1338 Five Millennia of Levantine History from Ancient Canaanite and Present-Day Lebanese Genome  
1339 Sequences. *Am J Hum Genet* 101, 274–282. <https://doi.org/10.1016/J.AJHG.2017.06.013>.
- 1340 136. Olalde, I., Brace, S., Allentoft, M.E., Armit, I., Kristiansen, K., Booth, T., Rohland, N., Mallick, S.,  
1341 Szécsényi-Nagy, A., Mittnik, A., et al. (2018). The Beaker phenomenon and the genomic  
1342 transformation of northwest Europe. *Nature* 555, 190–196.  
1343 <https://doi.org/10.1038/NATURE25738>.
- 1344 137. Patterson, N., Isakov, M., Booth, T., Büster, L., Fischer, C.E., Olalde, I., Ringbauer, H., Akbari, A.,  
1345 Cheronet, O., Bleasdale, M., et al. (2021). Large-scale migration into Britain during the Middle to  
1346 Late Bronze Age. *Nature* 2021 601:7894 601, 588–594. <https://doi.org/10.1038/s41586-021-04287-4>.
- 1348 138. Feldman, M., Master, D.M., Bianco, R.A., Burri, M., Stockhammer, P.W., Mittnik, A., Aja, A.J.,  
1349 Jeong, C., and Krause, J. (2019). Ancient DNA sheds light on the genetic origins of early Iron Age  
1350 Philistines. *Sci Adv* 5, 61–64.

- 1351 139. Moorjani, P., Sankararaman, S., Fu, Q., Przeworski, M., Patterson, N., and Reich, D. (2016). A  
1352 genetic method for dating ancient genomes provides a direct estimate of human generation  
1353 interval in the last 45,000 years. *Proc Natl Acad Sci U S A* *113*, 5652–5657.  
1354 <https://doi.org/10.1073/pnas.1514696113>.
- 1355 140. Choudhury, A., Aron, S., Botigué, L.R., Sengupta, D., Botha, G., Bensellak, T., Wells, G., Kumuthini,  
1356 J., Shriner, D., Fakim, Y.J., et al. (2020). High-depth African genomes inform human migration and  
1357 health. *Nature* *586*, 741–748. <https://doi.org/10.1038/s41586-020-2859-7>.
- 1358 141. Speidel, L., Forest, M., Shi, S., and Myers, S.R. (2019). A method for genome-wide genealogy  
1359 estimation for thousands of samples. *Nat Genet* *51*, 1321–1329. [https://doi.org/10.1038/s41588-](https://doi.org/10.1038/s41588-019-0484-x)  
1360 [019-0484-x](https://doi.org/10.1038/s41588-019-0484-x).
- 1361