1

# The Genomic History of the Middle East

Mohamed A. Almarri[1,2*], Marc Haber[3,4*], Reem A. Lootah[2], Pille Hallast[1,5],
Saeed Al Turki[6,7], Hilary C. Martin[1], Yali Xue[1], Chris Tyler-Smith[1]

[1]Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton CB10 1SA, UK
[2]Department of Forensic Science and Criminology, Dubai Police GHQ, Dubai, United Arab Emirates.
[3]Institute of Cancer and Genomic Sciences, University of Birmingham, Birmingham B15 2TT, UK
[4]Centre for Computational Biology, University of Birmingham, Birmingham B15 2TT, UK
[5]Institute of Biomedicine and Translational Medicine, University of Tartu, Tartu, 50411, Estonia
[6]Translational Pathology, Department of Pathology and Laboratory Medicine, King Abdulaziz Medical City, Ministry of National Guard- Health Affairs, Riyadh, Saudi Arabia
[7]Department of Genetics & Genomics, College of Medicine and Health Sciences, United Arab Emirates University. Al Ain, United Arab Emirates

*Correspondence: ma17@sanger.ac.uk (M.A.A), m.haber@bham.ac.uk (M.H.)

## Abstract

The Middle East is an important region to understand human evolution and migrations, but is underrepresented in genetic studies. We generated and analysed 137 high-coverage physically-phased genome sequences from eight Middle Eastern populations using linked-read sequencing. We found no genetic traces of early expansions out-of-Africa in present-day populations, but find Arabians have elevated Basal Eurasian ancestry that dilutes their Neanderthal ancestry. A divergence in population size within the region starts before the Neolithic, when Levantines expanded while Arabians maintained small populations that could have derived ancestry from local epipaleolithic hunter-gatherers. All populations suffered a bottleneck overlapping documented aridification events, while regional migrations increased genetic structure, and may have contributed to the spread of the Semitic languages. We identify new variants that show evidence of selection, some dating from the onset of the desert climate in the region. Our results thus provide detailed insights into the genomic and selective histories of the Middle East.

**Introduction**

Global whole-genome sequencing projects have provided insights into human diversity, dispersals, and past admixture events (Bergström *et al.*, 2020; Mallick *et al.*, 2016; GenomeAsia100K Consortium, 2019; 1000 Genomes Project Consortium *et al.*, 2015). However, many populations remain understudied, which restricts our understanding of genetic variation and population history, and may exacerbate health inequalities (Sirugo *et al.*, 2019). A region particularly understudied by large-scale sequencing projects is the Middle East. Situated between Africa, Europe and South Asia, it forms an important region to understand human evolution, history and migrations. The demographic history and prehistoric population movements of Middle Easterners are poorly understood, as are their relationships among themselves and to other global populations. The region contains some of the earliest evidence of modern humans outside Africa, with fossils dated to ~180 thousand years ago (kya) and ~85 kya identified in the Levant and North West Arabia, respectively (Hershkovitz *et al.*, 2018; Groucutt *et al.*, 2018). In addition, tool kits suggesting their presence have been identified in South East Arabia dating to ~125 kya (Armitage *et al.*, 2011). Although most of Arabia is a hyper-arid desert today, there were several humid periods resulting in a 'green Arabia' in the past which facilitated human dispersals, with the onset of the current desert climate thought to have started around 6 kya (Petraglia *et al.*, 2020). The toggling from humid to arid periods has been proposed to result in population movements adapting to the climate. The Neolithic transition within Arabia may have developed independently within the region, or resulted from an expansion of Levantine Neolithic farmers southwards (Drechsler, 2009; Uerpmann, *et al.,* 2010; Crassard *et al.*, 2013a; Crassard *et al.,* 2013b; Hilbert *et al.,* 2015). To address such questions, we generated and analysed a high-coverage physically-phased open-access dataset of populations from the Arabian Peninsula, the Levant and Iraq. In addition to creating a catalogue of genetic variation in an understudied region that will assist future medical studies, we have investigated the population structure, demographic and selective histories, and admixture events with modern and archaic humans.

68  **Results**

69  **Dataset and Sample Sequencing**

70  We sequenced 137 whole genomes from eight Middle Eastern populations (Figure 1A) to an

71  average coverage of 32x using a library preparation method that preserves long-range

72  information from short reads, and aligned them to the GRCh38 reference (Methods). An

73  advantage of using this 'linked-read' technology is the reconstruction of physically-phased

74  haplotypes and improved alignments at repetitive regions which confound short-read

75  aligners (Figure S1). All populations investigated speak Arabic, a Semitic language of the

76  Afro-Asiatic language family, with the exception of the Iraqi Kurdish group who speak

77  Kurdish, an Iranian language belonging to the Indo-European family. After quality control

78  (Methods) we identify 23.1 million single nucleotide variants (SNVs). We compared our

79  dataset to variants identified in the recently released Human Genome Diversity Project

80  (HGDP-CEPH) study (Bergström *et al.*, 2020). We find 4.9 million autosomal SNVs in our

81  dataset that are not found in the HGDP. As expected, most of the new variants are rare

82  (93%, < 1% minor allele frequency); however, ~370,000 are common (> 1%). Interestingly,

83  most of these common variants are outside the accessibility mask defined by Bergstrom *et*

84  *al.*, 2020 (~246,000). This illustrates the importance of sequencing genetically under-

85  represented populations such as Middle Easterners and the inclusion of regional-private

86  variants in future medical studies. It also demonstrates that a significant amount of unknown

87  variation resides in regions that are not accessible to standard short-reads.

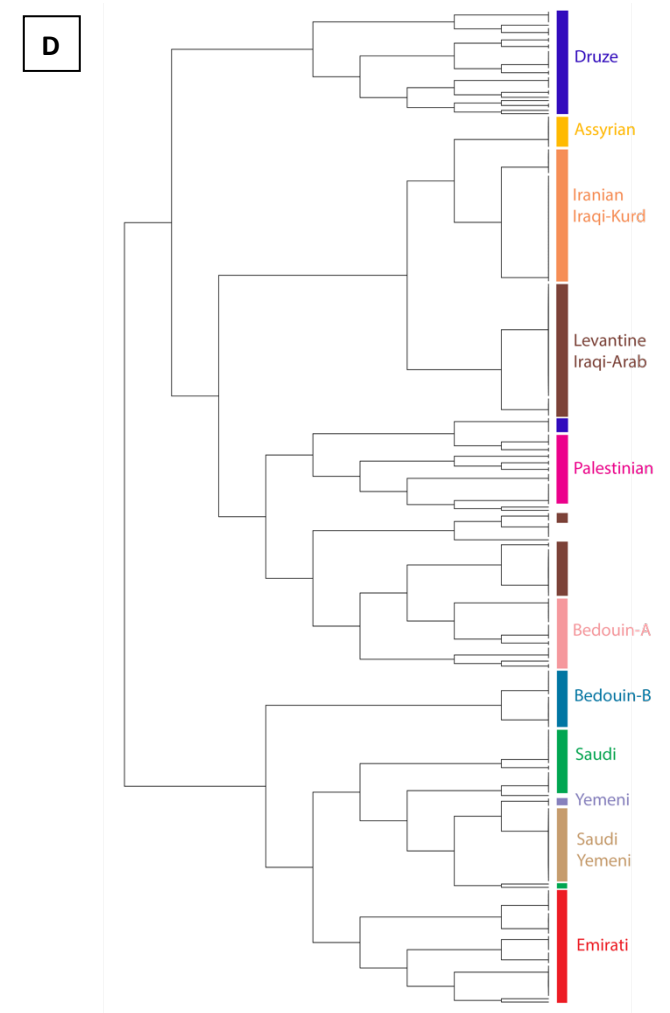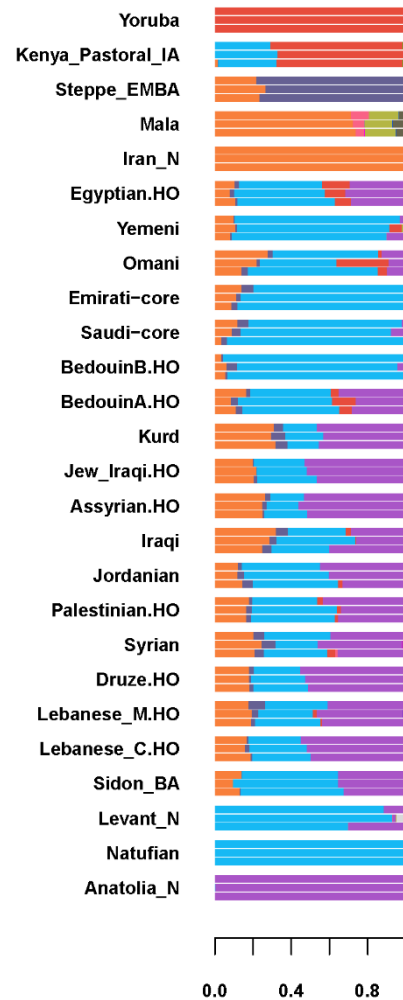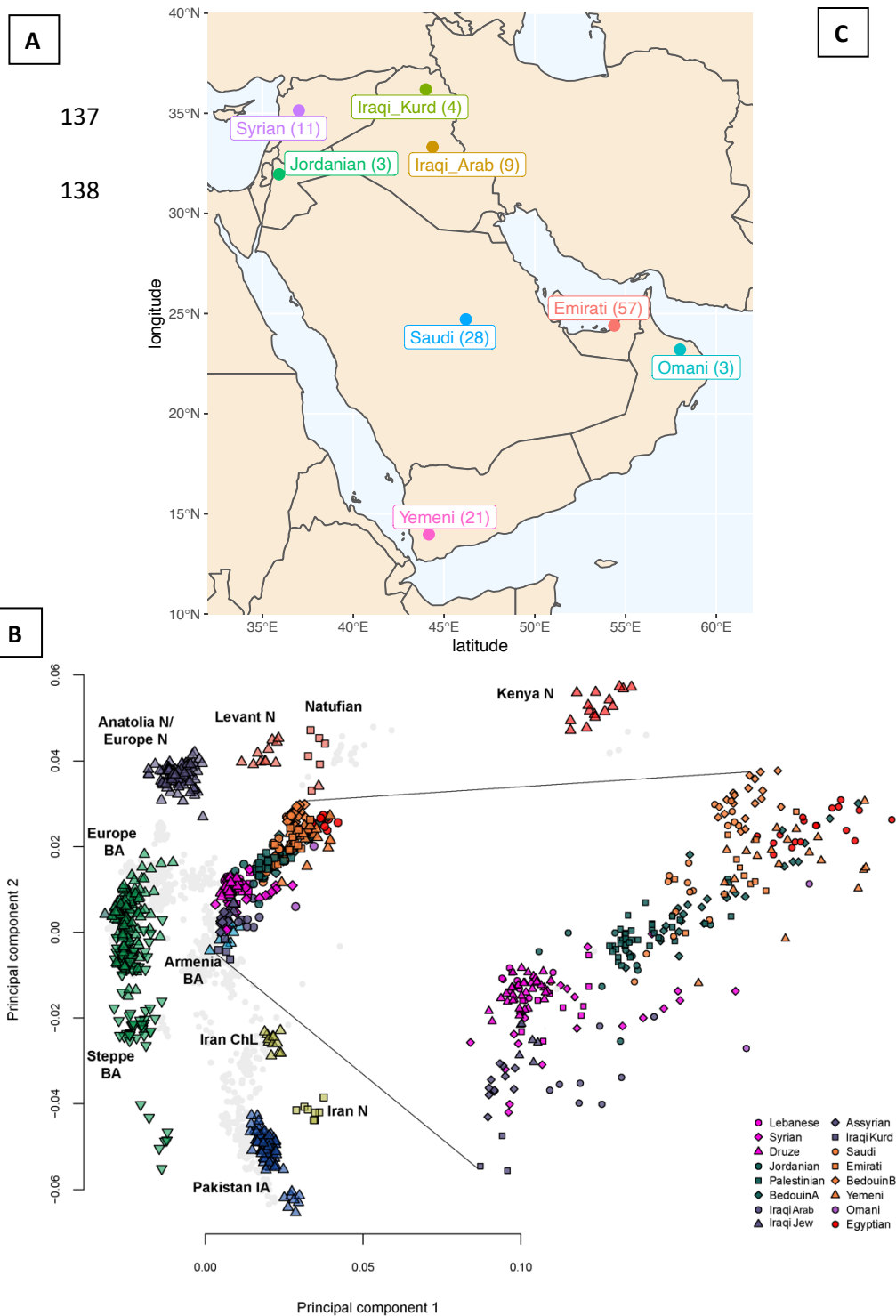88  **Population Structure and Admixture**

89  Uncovering population structure and past admixture events is important for understanding

90  population history and for designing and interpreting medical studies. We explored the

91  structure and diversity of our dataset using both single-variant and haplotype-based

92  methods. After merging our dataset with global populations, fineSTRUCTURE (Lawson *et*

93  *al.*, 2012) identified genetic clusters that are concordant with geography, and showed that

94  self-labelled populations generally formed distinct clusters (Figure 1D and S2). Populations

95  from the Levant and Iraq (Lebanese, Syrians, Jordanians, Druze and Iraqi-Arabs) clustered

96  together, while Iraqi-Kurds clustered with Central Iranian populations. Arabian populations

97  (Emiratis, Saudis, Yemenis and Omanis) clustered with Bedouins (BedouinB) from the

98  HGDP. The fineSTRUCTURE analysis thus allowed us to identify subpopulations who show

99  minimal admixture, which we herein label 'core'.

100  We next analysed our samples in the context of ancient regional and global populations.

101  Principal component analysis (Figures 1B and S3) shows that present-day Middle

3

102    Easterners are positioned between ancient Levantine hunter-gatherers (Natufians), Neolithic

103    Levantines (Levant_N), Bronze Age Europeans and ancient Iranians. Arabians and

104    Bedouins are positioned close to ancient Levantines, while present-day Levantines are

105    drawn towards Bronze Age Europeans. Iraqi Arabs, Iraqi Kurds and Assyrians appear

106    relatively closer to ancient Iranians and are positioned near Bronze Age Armenians. We find

107    that most present-day Middle Easterners can be modelled as deriving their ancestry from

108    four ancient populations (Table S1): Levant_N, Neolithic Iranians (Iran_N), Eastern Hunter

109    Gatherers (EHG), and a ~4,500 year old East African (Mota). We observe a contrast

110    between the Levant and Arabia: Levantines have excess EHG ancestry (Figure S4), which

111    we showed previously had arrived in the Levant after the Bronze Age along with people

112    carrying ancient south-east European and Anatolian ancestry (Haber *et al.,* 2017, Haber *et*

113    *al.* 2020). Our results here show this ancestry remained mostly confined to the Levant

114    region. Another contrast between the Levant and Arabia is the excess of African ancestry in

115    Arabian populations. We find that the closest source of African ancestry for most populations

116    in our dataset is Bantu Speakers from Kenya, in addition to contributions from Nilo-Saharan

117    speakers from Ethiopia specifically in the Saudi population. We estimate that African

118    admixture in the Middle East occurred within the last 2,000 years, with most populations

119    showing signals of admixture around 500-1,000 years ago (Figure S5 and Table S2).

120    In addition to differences in EHG and African ancestry, we observe an excess of Natufian

121    ancestry in the South compared with the North (Figure S4). Model-based clustering also

122    shows that Arabian populations have little Anatolia Neolithic (Anatolia_N) ancestry compared

123    with the modern-day Levantines (purple component in Figure 1C). This result is intriguing

124    since Levant_N shares significant ancestry with Anatolia_N compared with the preceding

125    local Natufian population (Lazaridis *et al.*, 2016), and a hypothesized Neolithic expansion

126    from the Levant to Arabia should have also carried Anatolia_N ancestry. The difference in

127    ancient Anatolian ancestry could also be from post-Bronze Age events, which resulted in

128    differences in EHG ancestry in the region (Haber *et al.*, 2020). When we substitute Levant_N

129    with Natufians, we found that Arabians could be successfully modelled (Table S1 and Figure

130    S7), suggesting that they could derive all of their local ancestry from Natufians without

131    requiring additional ancestry from Levant_N. On the other hand, none of the present-day

132    Levantines could be modelled as such.

133    In addition to the local ancestry from Epipaleolithic/Neolithic people, we find an ancestry

134    related to ancient Iranians that is ubiquitous today in all Middle Easterners (orange

135    component in Figure 1C; Table S1). Previous studies showed that this ancestry was not

136    present in the Levant during the Neolithic period, but appears in the Bronze Age where

**Figure 1. Overview of the dataset and population structure of the Middle East**. **A)** Map illustrating the populations sampled in this study, with numbers in brackets illustrating number of individuals. **B)** Principal component analysis of ancient and modern populations. Eigenvectors were inferred with present-day populations from the Middle East, North and East Africa, Europe, Central and South Asia. The ancient samples were then projected onto the plot (all modern non-Middle Easterners shown as grey points). Plot also shows a magnification of the modern Middle Eastern cluster. See Figure S3 for more details. **C)** Temporally-aware model-based clustering using ~80,000 transversions and 9 time points. Showing K=13 when the Anatolia_N and Natufian components split. See Figure S5 for more details. ".HO" suffix refers to samples from the Human Origins Dataset. **D)** Finestructure tree of modern-day Middle Easterners with population clusters highlighted. See Figure S2 for more details.

139    ~50% of the local ancestry was replaced by a population carrying ancient Iran-related

140    ancestry (Lazaridis *et al.*, 2016). We explored whether this ancestry penetrated both the

141    Levant and Arabia at the same time, and found that admixture dates mostly followed a North

142    to South cline, with the oldest admixture occurring in the Levant region between 3,900 and

143    5,600 ya (Table S3), followed by admixture in Egypt (2,900-4,700 ya), East Africa (2,200-

144    3,300) and Arabia (2,000-3,800). These times overlap with the dates for the Bronze Age

145    origin and spread of Semitic languages in the Middle East and East Africa estimated from

146    lexical data (Kitchen *et al.*, 2009; Figure S8). This population potentially introduced the Y-

147    chromosome haplogroup J1 into the region (Chiaroni *et al.*, 2010; Lazaridis *et al.*, 2016). The

148    majority of the J1 haplogroup chromosomes in our dataset coalesce around ~5.6 [95% CI,

149    4.8-6.5] kya, agreeing with a potential Bronze Age expansion; however, we do find rarer

150    earlier diverged lineages coalescing ~17 kya (Figure S9). The haplogroup common in

151    Natufians, E1b1b, is also frequent in our dataset, with most lineages coalescing ~8.3 [7-9.7]

152    kya, though we also find a rare deeply divergent Y-chromosome which coalesces 39 kya

153    (Figure S9).

154    **Effective Population size and Separation History**

155    Historical effective population sizes can be inferred through the distribution of coalescence

156    times between chromosomes sampled from a population (Li and Durbin, 2011). However,

157    there is limited resolution in recent periods using single human genomes, while errors in

158    haplotype phasing create artefacts when using multiple genomes (Schiffels and Durbin,

159    2014; Terhorst *et al.*, 2017). Although methods have been developed that extend these

160    approaches by incorporating the allele frequency spectrum from unphased genomes, they

161    do not have resolution at recent times, for e.g. through the metal ages (Terhorst *et al.*, 2017;

162    Bergström *et al.*, 2020). By leveraging recent advances in generating genome-wide

163    genealogies (Speidel *et al.*, 2019), and the large number of physically-phased samples in

164    our study, we could estimate the effective population size of each population in our dataset

165    up to very recent times - 1 kya (Figure 2A and S16A). We found all Middle Easterners had a

166    significant decrease in population size, around the out-of-Africa event ~50-70 kya. The

167    recovery from this bottleneck follows a similar pattern until 15-20kya, when a contrast

168    between the Levant and Arabia started to emerge. All Levantine and Iraqi populations

169    continued to show a substantial population expansion, while Arabians maintained similar

170    sizes. This contrast is noteworthy since it starts after the end of the Last Glacial Maximum

171    and becomes prominent during the Neolithic, when agriculture developed in the Fertile

172    Crescent and led to settled societies supporting larger populations. Following the Neolithic,

173    and with the start of the aridification of Arabia around 6kya, Arabian populations experienced
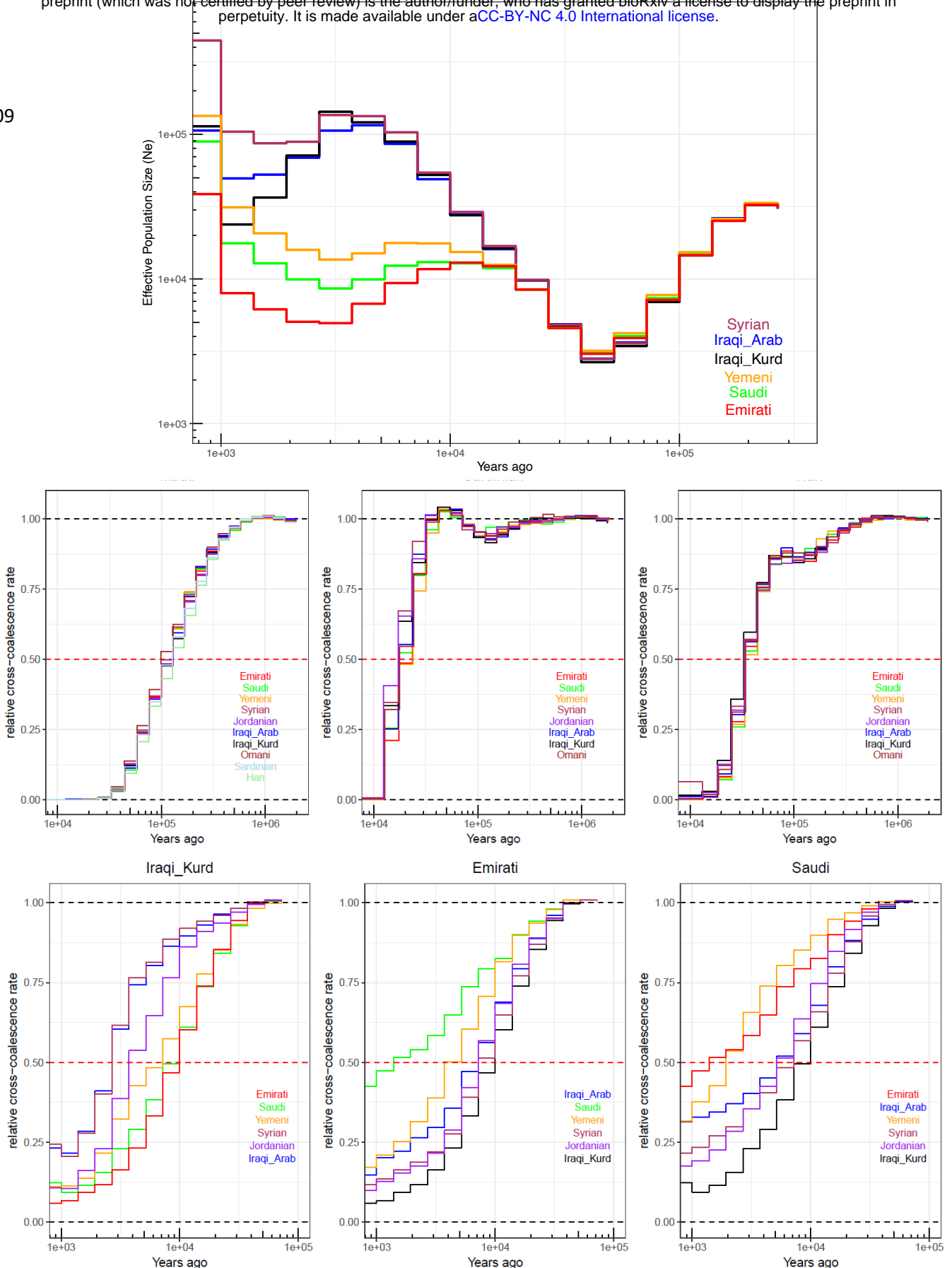
174  a bottleneck while Levantines continued to increase in size. The expansion in Levantines

175  then plateaus and their population size decreases around the 4.2 kiloyear aridification event

176  (Weiss et al., 1993). The decline in Emiratis is especially prominent, reaching an effective

177  population size of ~5,000, more than 20 times smaller than Levantines and Iraqis at the

178  same time period. A recovery can be observed in the past 2 ky.

179  We next studied the population separation history of Middle Eastern populations among

180  themselves and from global populations. The importance of accurate phasing in this analysis

181  is illustrated by an earlier finding that suggested, based on statistically phased data, that

182  modern-day Papuans harbour ancestry of an early expansion of modern humans out of

183  Africa (Pagani *et al.*, 2016). However, this was not replicated using physically-phased

184  genomes, suggesting it was caused by a statistical phasing artefact (Bergström *et al.*, 2020).

185  Conversely, when exploring population separation history at recent times, rare variants

186  become more informative but are less accurately phased by statistical methods, and are

187  unlikely to be present in reference panels. We first tested whether present-day Middle

188  Easterners harbour ancestry from an early human expansion out of Africa by comparing the

189  split times of our populations with physically-phased samples from the HGDP (Figure 2B and

190  S10). Using a relative cross-coalescent rate (rCCR) of 0.5 as a heuristic estimate of split

191  time, we found that Levantines, Arabians, Sardinians and Han Chinese share the same split

192  time, and additionally the same gradual pattern of separation, from Mbuti ~120kya. We then

193  compared the populations in our dataset with Sardinians and found they split ~20 kya, with

194  Levantines showing a slightly more recent divergence than Arabians. In contrast to the

195  gradual separation patterns to Mbuti, Sardinians show more of a clean split to all Middle

196  Eastern populations. Notably, all lineages within the Levant and Arabia, and in addition to

197  lineages within all Middle Eastern populations and Sardinians, coalesce within 40 kya. These

198  results collectively suggest that present-day Middle Eastern populations do not harbour any

199  significant traces from an earlier expansion out of Africa, and all descend from the same

200  population that expanded out of the continent ~50-60 kya.

201  We then compared the separation times of populations within the Middle East, and found the

202  oldest divergence times were between Arabia and the Levant/Iraq (Figure 2C and S16B).

203  The Emiratis split from Iraqi Kurds around 10 kya, and more recently around 7 kya from

204  Jordanians, Syrians and Iraqi Arabs. Saudi split times from the same populations appear

205  more recent, around 5-7 kya, while the Yemeni separation curves are intermediate between

206  the Emirati and Saudi curves. The split times between Arabia and the Levant predate the

207  Bronze Age, agreeing with our phylogenetic modelling that if a Bronze Age expansion into

208

209



**Figure 2. Population size and separation history. Top)** Effective population size histories for Middle Eastern populations. More details in Figure S16A. **Center)** Separation history between Mbuti, Sardinians and Han (indicated at the top of each panel) with each of the Middle Eastern populations (identified within each panel). All Middle Eastern populations show similar split time with each of these global populations. **Bottom)** Separation history within the Middle East (population indicated at the top of each panel, and within each panel). More comparisons show in Figure S16B. Note the different X-axis scales.

8

210    Arabia occurred, it did not result in a complete replacement of ancestry.
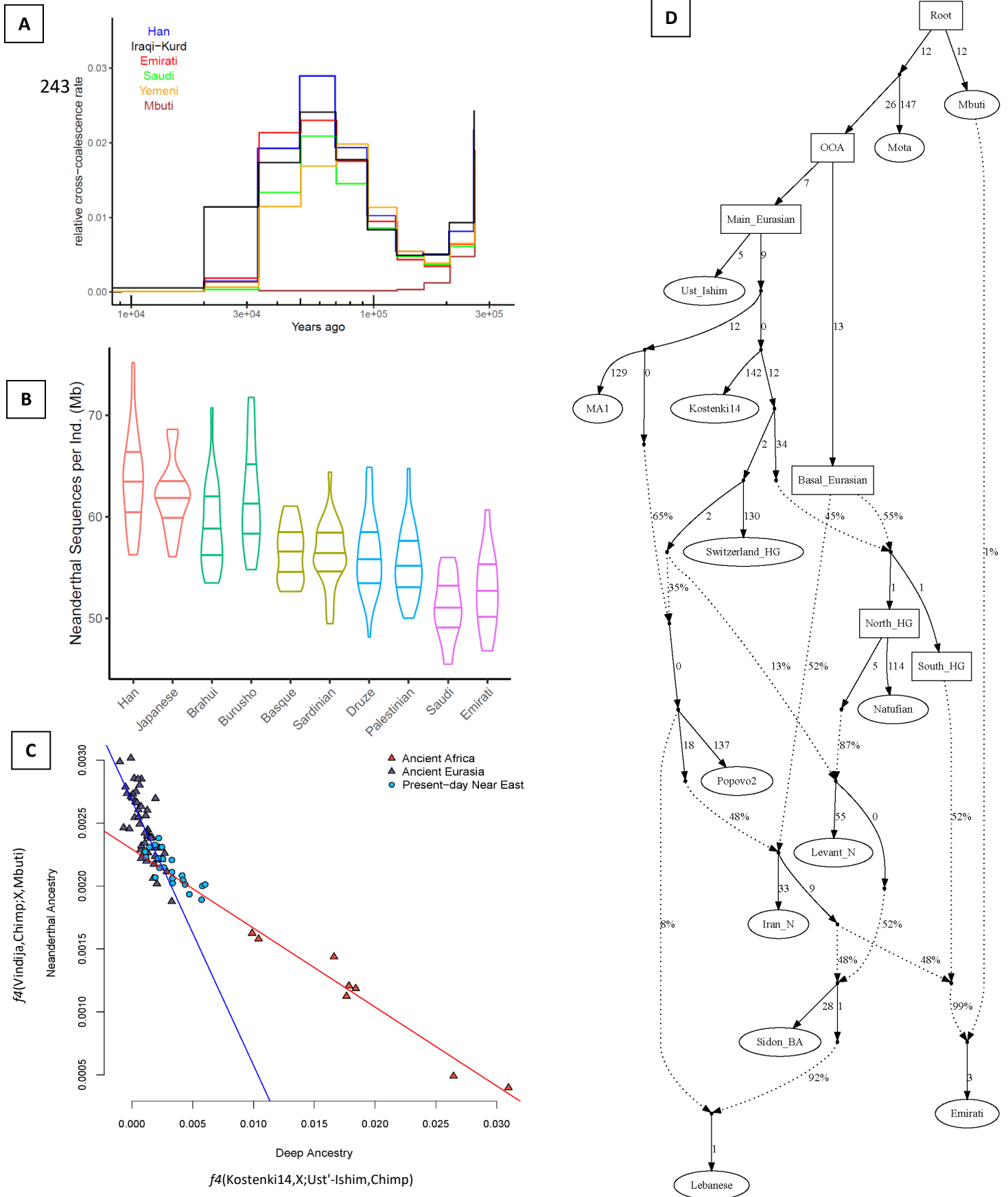
211    Within the Levant and Iraq, all splits occurred in the past 3-4 ky. Within Arabia, Yemenis split

212    from Emiratis ~4 kya and Saudis appear as the least divergent population to both the

213    Emiratis and Yemenis, with recent splits within the last 2ky. The separation history of the

214    region suggests continuous historical gene flow occurring between the Levant/Iraq and

215    Central Arabia, and in addition between Central Arabia to the Southeast, and separately to

216    the Southwest in Yemen.

217    **Archaic introgression and deep ancestry in the Middle East**

218    The similar amount of Neanderthal ancestry in most non-African populations and the low

219    diversity of introgressed haplotypes suggest that modern humans likely experienced a single

220    pulse of Neanderthal admixture as they expanded out of Africa (Bergström *et al.*, 2020).

221    Middle Eastern populations have previously been shown to have lower Neanderthal ancestry

222    than European and East Asian populations (Rodriguez-Flores *et al.*, 2016; Bergström *et al.*,

223    2020); however, the interpretation of this finding is complicated by recent African admixture

224    'diluting' Neanderthal ancestry (Haber et al., 2016). In addition, some analyses require the

225    use of an outgroup, which, if it itself contains Neanderthal ancestry, can bias estimates

226    (Chen *et al.*, 2020). To investigate Neanderthal introgression in our dataset, we exploited the

227    accurate phasing of our samples and compared cross-coalescent rates with the high

228    coverage Vindija Neanderthal genome (Prüfer *et al.*, 2017). All Middle Easterners showed an

229    archaic admixture signal at a time point similar to other Eurasians (Figure 3A).

230    We then used an identity-by-descent-based method, IBDmix, which directly compares a

231    target population to the Neanderthal genome to detect haplotypes of Neanderthal origin

232    (Chen *et al.*, 2020). We ran IBDmix on our samples and the HGDP dataset, recovering

233    segments totalling ~1.27 Gb that are of likely Neanderthal origin. When comparing the

234    amount of Neanderthal haplotypes that are private to our dataset but not present in other

235    non-Middle Eastern Eurasians, we found only ~25 Mb in total, illustrating that the vast

236    majority of Neanderthal haplotypes in the region are shared with other populations.

237    However, we do find relatively large introgressed haplotypes (~500kb) that are very rare

238    globally, but reach high frequencies in Arabia (Figure S12).

239    We then compared the average number of total Neanderthal bases per population, and

240    found lower values in Arabia in comparison to other Eurasian populations, including

241    Levantines. The Druze and Sardinians, for example, have similar amounts (average ~56.4

242    Mb per individual) of Neanderthal ancestry (Figure 3B). In contrast, in Arabia, Emirati.core

243



**Figure 3. Archaic introgression and deep structure in the Middle East**. **A**) Relative cross coalescent rate against Vindija Neanderthal. Note the y-axis range. **B**) Distribution of total length of Neanderthal sequences (Mb) per sample in each population. Horizontal lines depict 25%, 50%, and 75% quantiles. Colors reflect regional grouping. **C**). Neanderthal ancestry is negatively correlated with a deep ancestry in the Middle East. Two clines explain the depletion of Neanderthal Ancestry in Middle Easterners; one formed by basal Eurasian ancestry and the other is African ancestry. We plot regression lines using the ancient Africans (red) and the ancient Eurasians (blue). **D**) A possible model for the population formation in the Middle East. Populations in ellipses are sampled populations while populations in boxes are hypothetical. Worst f-statistics: (Lebanese, Emirati; Lebanese, Emirati) Z score = -2.83. See Figure S11 for alternative graph models. BA: Bronze Age; HG: Hunter-gatherer.

244 and Saudi.core have an average of 52.7 Mb and 52.1 Mb Neanderthal ancestry respectively,

245 which is ~8% lower than the Druze and Sardinians, and ~20% less than Han Chinese. Since

246 Emirati.core and Saudi.core have less than 3% of African ancestry, the depletion of

247 Neanderthal ancestry in Arabia cannot be explained by the African ancestry alone. Lazaridis

248 *et al.*, (2014) proposed that a basal Eurasian population, with low-to-no Neanderthal

249 ancestry, had contributed different proportions to ancient and modern Eurasians, reaching

250 ~50% in Neolithic Iranians and Natufians. Since Arabians have an excess of Natufian-like

251 ancestry compared to elsewhere in the Middle East, we find they also carry an excess of

252 basal Eurasian ancestry which will reduce their Neanderthal ancestry. In addition, most

253 modern Middle Easterners carry African ancestry from recent admixture which also

254 contributes to their deep ancestry (relative to the time of a main Eurasian ancestry). We find

255 a negative correlation (Pearson's r = -0.81, *P* = 2.76e-06) between the increase in deep

256 ancestry and the amount of Neanderthal ancestry in the modern Middle Easterners. When

257 testing all ancient populations we find two clines (Figure 3C) explaining the depletion of

258 Neanderthal ancestry: The first is formed by African ancestry while the second is formed by

259 a Basal Eurasian ancestry in ancient Eurasians. Middle Easterners appear to be affected by

260 both clines since they harbour both ancestries.

261 **Selection**

262 There is currently a limited understanding of the effects of selection in Arabian populations,

263 with the current hyper-arid climate and a long-term nomad-like subsistence potentially

264 exerting selective pressure for adaptations. To explore this, we searched genome-wide

265 genealogies for lineages carrying mutations that have spread unusually quickly (Speidel *et*

266 *al.*, 2019) at a conservative genome-wide threshold ($P < 5\times10^{-8}$). Previous studies identified

267 two correlated variants (rs41380347 and rs55660827), distinct from the known European

268 variant (rs4988235), that are associated with lactase persistence in Arabia (Imtiaz et al.

269 2007; Enattah et al. 2008). For the Arabian variant rs41380347, we found evidence for

270 strong selection (s = 0.011, logLR = 13.27), similar to, but slightly weaker than, the reported

271 strength of selection at rs4988235 in Europeans (s = 0.016-0.018; Mathieson and Mathieson

272 2018; Stern et al. 2019). The variant is present at highest frequency in the core Arabian

273 populations: ~50% in Saudis and Emiratis, and at a much lower frequency in the Levant and

274 Iraq (4%). Remarkably, the variant is not present in any Eurasian or African population in the

275 1000 Genome Project (1KG). We also did not find the variant in published ancient Eurasian

276 whole genomes, including ancient Levantines and Iranians, consistent with a recent origin of

277 the haplotype within the Middle East and subsequent increase in frequency due to selection.

278 We find the variant had a rapid increase in frequency between 9 kya and the present day

11

279    (Figure 4B). Notably, this period overlaps with the transition from a hunter-gatherer to a

280    herder-gatherer lifestyle in Arabia (Petraglia *et al*. 2020).

281    We also identified additional variants that show an increase in frequency recently (Figure

282    4C-D. A variant within *LMTK2*, rs11762534, which is also an eQTL for many genes*,* displays

283    evidence of selection (s=0.005; logLR = 16.49) and is associated with blood cell percentages

284    and malignant neoplasm of prostate. *LMTK2* encodes a serine/threonine kinase that is

285    implicated in diverse cellular processes including apoptosis, growth factor signalling and

286    appears essential for spermatogenesis in mice (Kawa *et al.,* 2006; Cruz *et al.,* 2019).

287    Outside the Middle East, the variant is highly stratified and is present at the highest

288    frequency in Europeans (1KG, 45%), but we find it at 66% frequency in the Arabian

289    populations. Intriguingly, the variant also shows differentiation in BedouinB (81%), while

290    appearing less frequent in Druze and Palestinians (both ~55%). We additionally looked for

291    strongly differentiated variants between Arabia and the Levant/Iraq (Figure S13). The variant

292    showing the most extreme population branch statistic in Yemenis is rs2814778, where the

293    derived allele results in the Duffy-null phenotype and is almost exclusively found in African

294    populations in the 1000 Genome Project. However, the variant is very common in Yemenis

295    (74%), and decreases in frequency moving northwards in the peninsula (59% in Saudis

296    while reaching 6% in Iraqi-Arabs). We find that across the genome this locus shows the

297    highest enrichment of African ancestry in the Middle East (Methods). As the average amount

298    of African ancestry in Yemenis and Saudis is ~9% and ~3% respectively, the high frequency

299    of this variant appears consistent with positive selection after African admixture. It has been

300    thought that the derived allele protects against *Plasmodium vivax* infection (Miller et al.,

301    1976), which has been historically present in Arabia.

302    An advantage of using genome-wide genealogies is its power to detect relatively weak

303    selection. We subsequently searched for evidence of polygenic adaptation in Arabian

304    populations across 20 polygenic traits specifically over the past 2,000 years (Methods). For

305    most traits, we find no, or inconclusive, evidence for recent directional selection, including

306    height, skin colour, and BMI (Figure 4A). However a few traits do show evidence, with

307    selection for higher years of education (EduYears) showing the strongest signal consistent

308    across all Arabian populations ($P$ = 0.0002 in Saudis). This has also been reported in the

309    British population (Stern *et al.*, 2020); however, the signal was shown to become attenuated

310    after conditioning on other traits, suggesting indirect selection via a correlated trait. In

311    contrast to findings in the British population (Stern *et al.*, 2020), we do not find selection

312    acting on traits such as sunburn, hair color and tanning ability. Within Arabia, the direction of

313    selection on most traits appears to be similar across populations, likely as a result of shared

12

314    ancestry; however, we note that the current varied environments across the region can

315    potentially cause different recent selective pressures. In Emiratis, we find evidence of

316    selection on variants increasing type 2 diabetes (T2D, $P$ = 0.004). This result is intriguing, as

317    the prevalence of T2D in Emiratis is among the highest globally and is partly thought to

318    result from strong recent shift to a sedentary lifestyle (Malik $et$ $al.$, 2005). We also find

319    nominal evidence of selection acting to increase levels of low-density lipoproteins (LDL; $P$ =

320    0.01) and decrease levels of Apoliprotein B (APOB; $P$ = 0.01) in the same population; but

321    they appear suggestive after adjusting for multiple testing ($P_{adj}$ = 0.06 at 5% FDR).
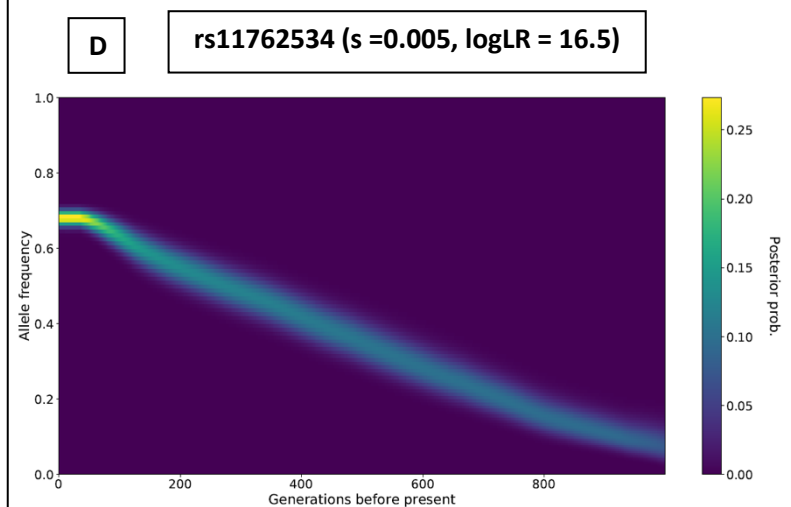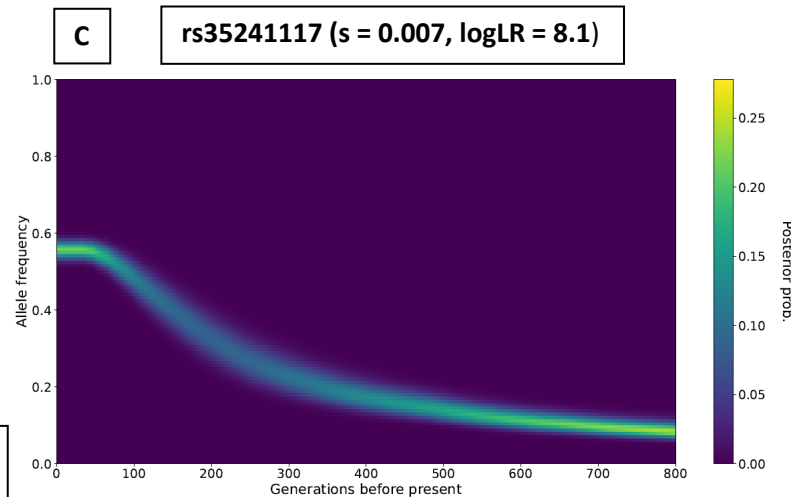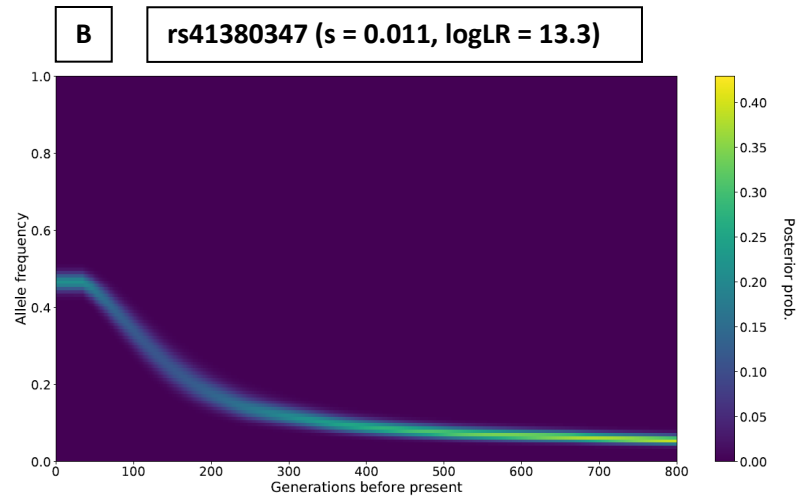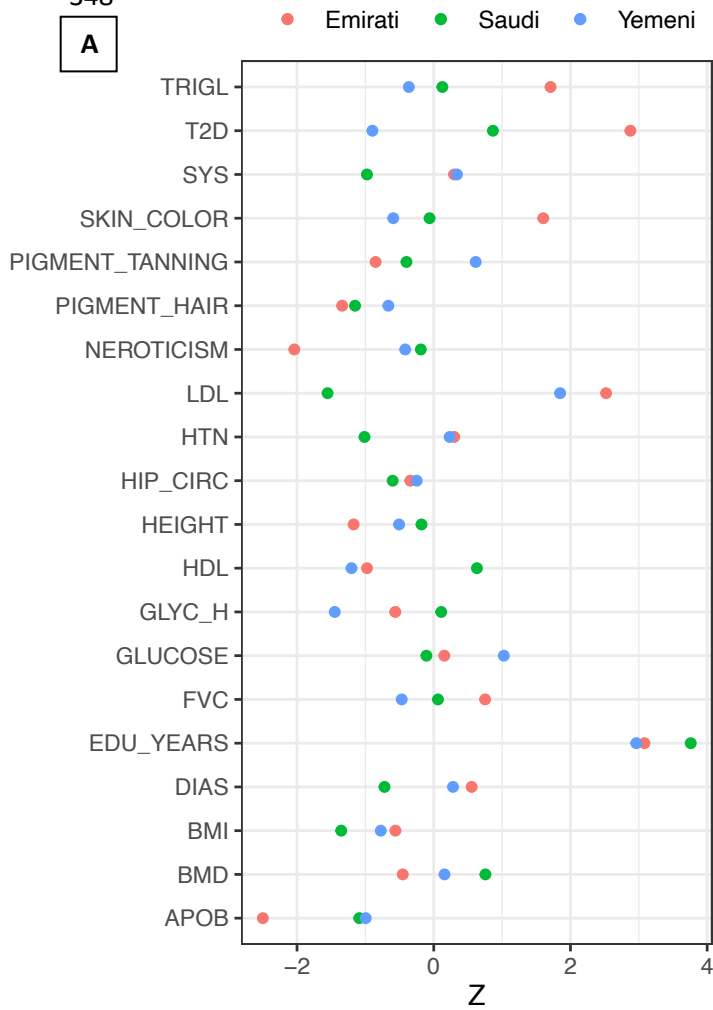
## Discussion

323    In this study we have generated a high-coverage open-access resource from the genetically

324    understudied Middle East region. To our knowledge, this is the first study where the whole

325    population investigated is experimentally-phased, allowing the reconstruction of large and

326    accurate haplotypes. We find millions of variants that are not catalogued in previous global

327    sequencing projects, with a significant proportion being common in the Middle East. A

328    majority of these common variants reside outside of short-read accessibility masks,

329    highlighting the limitation of standard short-read sequencing based studies.

330    The large number of physically-phased haplotypes allowed us to study population history

331    from relatively old periods (>100 kya) to very recent times (1 kya). We find no evidence that

332    an early expansion of humans out of Africa has contributed genetically to present-day

333    populations in the region. This finding adds to the growing consensus that all contemporary

334    non-African modern humans descend from a single expansion out-of-Africa, quickly followed

335    by admixture with Neanderthals, before populating the rest of the world (Mallick $et$ $al.,$ 2016;

336    Bergstrom $et$ $al.,$ 2020). We find that Middle Eastern populations have very little Neanderthal

337    DNA that is private to the region, with the vast majority shared with other Eurasians. We

338    demonstrate that Arabian populations have lower Neanderthal ancestry than Levantine,

339    European and East Asian populations and attribute this difference to elevated ancestry from

340    a basal Eurasian population, which did not admix with Neanderthals, in addition to recent

341    African admixture.

342    By modelling contemporary populations using ancient genomes, we identify differences

343    between the Levant and Arabia. The Levant today have higher European/Anatolian-related

344    ancestry and Arabia having higher African and Natufian-like ancestry. The contrast between

345    the regions is also illustrated by their population-size histories which diverge before the

346    Neolithic and suggest that the transition to a sedentary agricultural lifestyle allowed the

347

348



**Figure 4. Selection in Arabia. A)** Testing for recent polygenic selection, over the past 2000 years, on 20 traits within Arabian populations. Asterisks indicate the test is significant after correcting for multiple testing (FDR = 5%). TRIGL: Triglycerides; T2D: Type2 Diabetes; SYS: Systemic Blood Pressure; LDL: Low-density lipoproteins; HTN: Hypertension; HIP_CIRC: Hip circumference; HDL: High-density lipoproteins; GLYC_H: Glycosylated haemoglobin; FVC: Forced Vital Capacity; EDU_YEARS: Years of Education; DIAS: Diastolic blood pressure; BMI: Body Mass Index; BMD: Bone Mass Density; APOB: Apoliprotein B **B)** Historical allele trajectory of rs41380347 which is associated with lactase persistence and almost private to the Middle East. s = selection coefficient.  **C)** Frequency trajectory of rs35241117, located near *TNKS,* which is present at the highest frequency in Arabia globally and is associated with multiple traits including glomerular filtration rate, bone mineral density, BMI, standing height and hypertension. **D)** Frequency trajectory of rs11762534 which is associated with lymphocyte and neutrophil percentages and prostate neoplasm malignancy and is also present at the highest frequency in Arabia. s = selection coefficient.

14

349  growth of populations in the Levant, but was not paralleled in Arabia. It has been suggested

350  that population discontinuity occurred between the late Pleistocene and Early Holocene in

351  Arabia, and that the peninsula was repopulated by Neolithic farmers from the Fertile

352  Crescent (Uerpmann *et al.,* 2010). Our results do not support a complete replacement of the

353  Arabian populations by Levantine farmers. In addition our models suggest that Arabians

354  could have derived their ancestry from Natufian-like local hunter-gatherer populations

355  instead of Levantine farmers.

356  An additional source of ancestry needed to model modern Middle Easterners is related to

357  ancient Iranians. Our admixture tests show that this ancestry first reached the Levant, and

358  subsequently reached Egypt, East Africa and Arabia. The timings of these events

359  interestingly overlap with the origin and spread of the Semitic languages (Kitchen *et al.*,

360  2009), suggesting a potential population carrying this ancestry may have spread the

361  language. We find climate change associated aridification events to coincide with population

362  bottlenecks, with Arabians decreasing in size 6kya with the onset of the desert climate while

363  Levantines around the 4.2 kiloyear aridification event. This severe drought has been

364  suggested to have caused the collapse of kingdoms and empires in the Middle East and

365  South Asia, potentially reflected genetically in the signal we identify (Weiss, 2017). Future

366  ancient DNA studies from Arabia are needed to refine the formation of the Arabian

367  populations.

368  The application of ancestral recombination graphs to reconstruct the evolutionary history of

369  variants offers a powerful method to study natural selection. We refine and identity new

370  signals of selection in Arabian populations. The example of the lactase persistence

371  associated variant, which during the past few thousand years increased to a frequency

372  reaching 50% and is almost absent outside the region, demonstrates the importance of

373  studying underrepresented populations to understand human history and adaptations. Our

374  results indicate that polygenic selection might have played a role in increasing the frequency

375  of variants that were potentially beneficial in the past, but today are associated with diseases

376  such as T2D. We find few signals of polygenic selection in Arabian populations, which may

377  be a consequence of their long-term small effective population sizes which will theoretically

378  reduce the strength of selection. We also note that Middle Eastern populations are among

379  the most understudied populations included in GWAS (Sirugo *et al.*, 2019), which limits the

380  analysis of polygenic traits. Our study and the recent establishment of national biobanks in

381  the region are a step forward to reduce these disparities and offer an exciting opportunity to

382  explore, in the future, complex and disease traits in the Middle East.

383

15

## Author contributions

M.A.A., Y.X. and C.T-S. conceived this study. M.A.A. and M.H. designed and performed the analyses with contributions from P.H. M.A.A., M.H. Y.X. and C.T-S interpreted the results with input from H.C.M. R.A.L coordinated sample collection and extraction. S.A.T assisted in study design. M.A.A. and M.H. wrote the manuscript. Y.X. and C.T-S. supervised the work. All authors approved the final version of the paper. All authors declare no conflict of interest.

## Data availability

Raw read alignments are available from the European Nucleotide Archive (ENA) under study accession number xxxx. Phased VCFs are available on xxxx.

## References

1000 Genomes Project Consortium *et al.* (2015) 'A global reference for human genetic variation', *Nature*, 526(7571), pp. 68–74.

Armitage, S. J. *et al.* (2011) 'The Southern Route "Out of Africa": Evidence for an Early Expansion of Modern Humans into Arabia', *Science*, pp. o3–456. doi: 10.1126/science.1199113.

Bergström, A. *et al.* (2020) 'Insights into human genetic variation and population history from 929 diverse genomes', *Science*, 367(6484). doi: 10.1126/science.aay5012.

Chen, L. *et al.* (2020) 'Identifying and Interpreting Apparent Neanderthal Ancestry in African Individuals', *Cell*, pp. 677–687.e16. doi: 10.1016/j.cell.2020.01.012.

Chiaroni, J. *et al.* (2010) 'The emergence of Y-chromosome haplogroup J1e among Arabic-speaking populations', *European journal of human genetics: EJHG*, 18(3), pp. 348–353.

Crassard, R. *et al.* (2013) 'Beyond the Levant: first evidence of a pre-pottery Neolithic incursion into the Nefud Desert, Saudi Arabia', *PloS one*, 8(7), p. e68061.

Crassard, R. and Drechsler, P. (2013) 'Towards new paradigms: multiple pathways for the Arabian Neolithic', *Arabian Archaeology and Epigraphy*, pp. 3–8. doi: 10.1111/aae.12021.

Cruz, D. F., Farinha, C. M. and Swiatecka-Urban, A. (2019) 'Unraveling the Function of Lemur Tyrosine Kinase 2 Network', *Frontiers in pharmacology*, 10, p. 24.

Drechsler, P. (2009) 'The Dispersal of the Neolithic over the Arabian Peninsula'. doi: 10.30861/9781407305028.

424   Enattah, N. S. *et al.* (2008) 'Independent introduction of two lactase-persistence alleles into
425   human populations reflects different history of adaptation to milk culture', *American journal of*
426   *human genetics*, 82(1), pp. 57–72.

427   GenomeAsia100K Consortium (2019) 'The GenomeAsia 100K Project enables genetic
428   discoveries across Asia', *Nature*, 576(7785), pp. 106–111.

429   Groucutt, H. S. *et al.* (2018) 'Homo sapiens in Arabia by 85,000 years ago', *Nature Ecology*
430   *& Evolution*, pp. 800–809. doi: 10.1038/s41559-018-0518-2.

431   Haber, M. *et al.* (2016) 'Chad Genetic Diversity Reveals an African History Marked by
432   Multiple Holocene Eurasian Migrations', *American journal of human genetics*, 99(6), pp.
433   1316–1324.

434   Haber, M. *et al.* (2017) 'Continuity and Admixture in the Last Five Millennia of Levantine
435   History from Ancient Canaanite and Present-Day Lebanese Genome Sequences', *The*
436   *American Journal of Human Genetics*, pp. 274–282. doi: 10.1016/j.ajhg.2017.06.013.

437   Haber, M. *et al.* (2020) 'A Genetic History of the Middle East from an aDNA Time Course
438   Sampling Eight Points in the Past 4,000 Years', *The American Journal of Human Genetics*,
439   pp. 149–157. doi: 10.1016/j.ajhg.2020.05.008.

440   Hershkovitz, I. *et al.* (2018) 'The earliest modern humans outside Africa', *Science*,
441   359(6374), pp. 456–459.

442   Hilbert, Y. H. *et al.* (2015) 'Archaeological evidence for indigenous human occupation of
443   Southern Arabia at the Pleistocene/ Holocene transition: The case of al-Hatab in Dhofar,
444   Southern Oman', *Paléorient*, pp. 31–49. doi: 10.3406/paleo.2015.5674.

445   Imtiaz, F. *et al.* (2007) 'The T/G 13915 variant upstream of the lactase gene (LCT) is the
446   founder allele of lactase persistence in an urban Saudi population', *Journal of medical*
447   *genetics*, 44(10), p. e89.

448   Kitchen, A. *et al.* (2009) 'Bayesian phylogenetic analysis of Semitic languages identifies an
449   Early Bronze Age origin of Semitic in the Middle East', *Proceedings. Biological sciences /*
450   *The Royal Society*, 276(1668), pp. 2703–2710.

451   Lawson, D. J. *et al.* (2012) 'Inference of population structure using dense haplotype data',
452   *PLoS genetics*, 8(1), p. e1002453.

453   Lazaridis, I. *et al.* (2014) 'Ancient human genomes suggest three ancestral populations for
454   present-day Europeans', *Nature*, 513(7518), pp. 409–413.

455   Lazaridis, I. *et al.* (2016) 'Genomic insights into the origin of farming in the ancient Middle
456   East', *Nature*, 536(7617), pp. 419–424.

457   Li, H. and Durbin, R. (2011) 'Inference of human population history from individual whole-
458   genome sequences', *Nature*, 475(7357), pp. 493–496.

459   Malik, M. *et al.* (2005) 'Glucose intolerance and associated factors in the multi-ethnic
460   population of the United Arab Emirates: results of a national survey', *Diabetes research and*
461   *clinical practice*, 69(2), pp. 188–195.

462   Mallick, S. *et al.* (2016) 'The Simons Genome Diversity Project: 300 genomes from 142
463   diverse populations', *Nature*, 538(7624), pp. 201–206.

464  Mathieson, S. and Mathieson, I. (2018) 'FADS1 and the Timing of Human Adaptation to
465  Agriculture', *Molecular Biology and Evolution*, pp. 2957–2970. doi: 10.1093/molbev/msy180.

466  Miller, L. H. *et al.* (1976) 'The resistance factor to Plasmodium vivax in blacks. The Duffy-
467  blood-group genotype, FyFy', *The New England journal of medicine*, 295(6), pp. 302–304.

468  Pagani, L. *et al.* (2016) 'Genomic analyses inform on migration events during the peopling of
469  Eurasia', *Nature*, 538(7624), pp. 238–242.

470  Petraglia, M. D. *et al.* (2020) 'Human responses to climate and ecosystem change in ancient
471  Arabia', *Proceedings of the National Academy of Sciences*, pp. 8263–8270. doi:
472  10.1073/pnas.1920211117.

473  Prüfer, K. *et al.* (2017) 'A high-coverage Neandertal genome from Vindija Cave in Croatia',
474  *Science*, 358(6363), pp. 655–658.

475  Rodriguez-Flores, J. L. *et al.* (2016) 'Indigenous Arabs are descendants of the earliest split
476  from ancient Eurasian populations', *Genome research*, 26(2), pp. 151–162.

477  Schiffels, S. and Durbin, R. (2014) 'Inferring human population size and separation history
478  from multiple genome sequences', *Nature genetics*, 46(8), pp. 919–925.

479  Sirugo, G., Williams, S. M. and Tishkoff, S. A. (2019) 'The Missing Diversity in Human
480  Genetic Studies', *Cell*, p. 1080. doi: 10.1016/j.cell.2019.04.032.

481  Speidel, L. *et al.* (2019) 'A method for genome-wide genealogy estimation for thousands of
482  samples', *Nature genetics*, 51(9), pp. 1321–1329.

483  Stern, A. J. *et al.* (2020) 'Disentangling selection on genetically correlated polygenic traits
484  using whole-genome genealogies'. doi: 10.1101/2020.05.07.083402. *biorxiv*

485  Stern, A. J., Wilton, P. R. and Nielsen, R. (2019) 'An approximate full-likelihood method for
486  inferring selection and allele frequency trajectories from DNA sequence data'. *PLOS*
487  *Genetics* doi: 10.1101/592675.

488  Terhorst, J., Kamm, J. A. and Song, Y. S. (2017) 'Robust and scalable inference of
489  population history from hundreds of unphased whole genomes', *Nature Genetics*, pp. 303–
490  309. doi: 10.1038/ng.3748.

491  Uerpmann, H.-P., Potts, D. T. and Uerpmann, M. (2010) 'Holocene (Re-)Occupation of
492  Eastern Arabia', *The Evolution of Human Populations in Arabia*, pp. 205–214. doi:
493  10.1007/978-90-481-2719-1_15.

494  Weiss, H. *et al.* (1993) 'The genesis and collapse of third millennium north mesopotamian
495  civilization', *Science*, 261(5124), pp. 995–1004.

496  Harvey Weiss (ed.). (2017) Megadrought and collapse: from early agriculture to Angkor.
497  Oxford: Oxford University Press; 978-0-19-932919-9

498