1    On the post-glacial spread of human commensal *Arabidopsis thaliana*: journey to the east

2

3

4    Che-Wei Hsu[1], Cheng-Yu Lo[1], Cheng-Ruei Lee[1,2,3]

5

6    1. Institute of Ecology and Evolutionary Biology, National Taiwan University, No 1, Sec 4,

7    Roosevelt Rd, Taipei 10617, Taiwan ROC

8    2. Institute of Plant Biology, National Taiwan University, No 1, Sec 4, Roosevelt Rd, Taipei

9    10617, Taiwan ROC

10    3. Genome and Systems Biology Degree Program, National Taiwan University, No 1, Sec 4,

11    Roosevelt Rd, Taipei 10617, Taiwan ROC

12

13    Author of correspondence:

14    Cheng-Ruei Lee

15    Room 1129, Life Science Building, No 1, Sec 4, Roosevelt Rd, Taipei 10617, Taiwan ROC

16    886-2-33662535

17    chengrueilee@ntu.edu.tw

18

19

**Abstract**

With the availability of more sequenced genomes, our understanding of the evolution and demographic history of the model plant *Arabidopsis thaliana* is rapidly expanding. Here we compile previously published data to investigate global patterns of genetic variation. While the Southeast African accessions were reported to be the most divergent among worldwide populations, we found accessions from Yunnan, China to be genetically close to the sub-Saharan accessions. Our further investigation of worldwide chloroplast genomes identified several deeply diverged haplogroups existing only in Eurasia, and the African populations have lower variation in many haplogroups they shared with the Eurasian populations. Bayesian inferences of chloroplast demography showed that representative haplogroups of Africa exhibited long-term stable population size, suggesting recent selective sweep or bottleneck is not able to explain the lower chloroplast variation in Africa. Taken together, these patterns cannot be easily explained by a single out-of-Africa event. Several Eurasian chloroplast haplogroups had rapid population growth since 10 kya, presumably reflecting the recent expansion of the weedy non-relicts across Eurasia. Our demographic analysis on a chromosomal region un-affected by relict introgression also suggested the European, Central Asian, and Chinese Yangtze populations diverged no earlier than 15 kya, in contrast to previous estimates of 45 kya inferred from whole genome that likely contains relict admixture. The most recent expansion is observed in the Yangtze population of China less than 2000 years ago. Similar to Iberia, the western end of non-relict expansion reported in our previous study, in this eastern end of Eurasia we find clear traces of gene flow between the Yangtze non-relicts and the Yunnan relicts. Genes under strong selection and previously suggested to contribute to adaptation in the Yangtze valley are enriched for traces of relict introgression, especially those related with biotic and immune responses. The results suggest the ability of non-relicts to obtain locally adaptive alleles through admixture with relicts is an important factor contributing to the rapid expansion across the environmental gradients spanning the eastern to the western coast of Eurasia.

## Introduction

*Arabidopsis thaliana* is not only a model species in plant molecular biology, but also increasingly used to address major questions in ecology and evolution. The evolutionary history of this model plant is frequently revisited as more geographically diverse samples are collected and sequenced. The first continental-scale study of species-wide demography is published in 2016, where one globally distributed human commensal "non-relict" group as well as several "relict" groups located in relatively un-disturbed habitats were identified[1]. The follow-up study showed that the non-relicts originated recently near the Balkans and spread along the east-west axis of Eurasia, wiping out continental-wide relict populations while incorporating locally adaptive alleles from them[2]. As to the new world, the North American population arrived at around 1600 AD[3], and most of them likely came from the region near southeastern England and northwestern Germany, carrying a charismatic inversion in chromosome 4[4].

Durvasula *et al.*[5] investigated African *A. thaliana* and suggested an "out of Africa" demographic model, given the highest genomic variation and numbers of private alleles observed in African accessions. In this model, *A. thaliana* originated in Africa and diverged into three populations at ca. 90 kya: The Moroccan, Levantine and Southeast African groups. The migration of Moroccan population northwards to Iberia was illustrated as well as the Levantine migration wave westward into Europe and eastward to Central Asia at ca. 45 kya[6]. On the other hand, it remains unclear how *A. thaliana* first arrived Africa given all other *Arabidopsis* species were found in temperate Northern Hemisphere[7], and the pattern that Africa contains most variation can be equally likely explained by a non-African origin of *A. thaliana* followed by non-relict expansion wiping out most Eurasian variation[2].

Zou *et al.*[8] studied *A. thaliana* accessions from China and showed that the population in Yangtze River Basin arrived relatively recently. They also showed that genes associated with immune response as well as flowering time were significantly enriched in the list of selected genes in the Yangtze population. For flowering time, genetic mapping identified a candidate gene in chromosome 2, containing the *SVP* gene (AT2G22540) with a loss-of-function mutation accelerating flowering. It remains unclear what constitutes the source of adaptive allele in the Yangtze population – the sources of adaptation may be novel mutation, standing variation, or as we have shown for the Iberian non-relicts, introgression from locally adaptive relicts[2].

Here we compile global data and re-investigate the evolutionary history of *A. thaliana* with two specific aims. (1) We used the maternally inherited chloroplast genomes to study the species history from a different perspective and investigate the out-of-Africa hypothesis

86   in the context of global samples. (2) We wish to clear up the evolutionary history and timing

87   of non-relict expansion across Eurasia, especially whether the Chinese Yangtze population

88   represents the eastern end of non-relict expansion and whether adaptive introgression also

89   happened there.

90

91

92

93   **Results**

94

95   **Genetic variation in nuclear genomes**

96        To investigate the global patterns of *Arabidopsis thaliana* genomic variation, we

97   compiled data from the 1,001 genomes project[1], the African accessions[5], and the Chinese

98   accessions[8]. Phylogenetic tree using *Arabidopsis lyrata* as the outgroup[7] confirmed the

99   previous observation that the Tanzanian and South African accessions (hereafter the "TZSA"

100  clade) are most divergent to all others (Fig. 1a, Supplementary Fig. 1). Interestingly, two

101  accessions from Yunnan, China are also genetically close to the TZSA clade, inconsistent with

102  the single out-of-Africa event suggested previously[5,6].

103       Of the two Yunnan accessions, one (SRR2204703) has heterozygosity typical of

104  self-fertilizing *A. thaliana*, and the other (SRR2204316) has very high heterozygosity

105  (Supplementary Fig. 2). While the higher heterozygosity may result from recent outcrossing

106  events or DNA contamination, both samples have very low chloroplast heterozygosity as all

107  other accessions (< 0.001), making DNA contamination less likely. Since the number of

108  heterozygous sites in an individual reflects the number of SNPs between its two parents, we

109  suspected the high heterozygosity of SRR2204316 might result from the cross between two

110  genetically divergent groups, similar to a recent study in ancient humans[9].

111       The ADMIXTURE[10] K = 2 result supports this idea (Fig. 1b). While the Chinese Yangtze

112  population is highly similar to typical Eurasian non-relicts and the Yunnan accessions are

113  close to the TZSA group (Fig. 1a), admixture exists (Fig. 1b). We further investigated this with

114  ABBA-BABA tests (Table 1). We first used non-relicts from Western Europe, which in theory

115  had no gene flow with any of the Tanzanian/South African/Yunnan relict population, as a

116  reference group to test whether the Yunnan accessions had gene flow from non-relicts. The

117  sign of gene flow is highly significant for the highly heterozygous Yunnan accession (*P* =

118  5.17E-25, Table 1 Test A) but not the other (*P* = 0.539, Table 1 Test A). Using the relatively

119  un-admixed Yunnan accession as a reference, the Chinese Yangtze non-relicts showed strong

120  signs of gene flow from the Yunnan relicts (Table 1 Test B). Finally, since both the highly

4

121    heterozygous Yunnan accession and the Chinese Yangtze population showed signs of

122    admixture, the tests involving both groups are highly significant (Table 1 Test C). Therefore,

123    similar to Iberia, the far-eastern end of Eurasia is also affected by the rapid expansion of

124    non-relict population, with introgression from local relicts along the way.

125

126    **Genetic variation in chloroplast genomes**

127    　　Our results from the nulcear genome suggest a more complex demography than a

128    single out-of-Africa event[5,6]. To better understand this, we investigated the chloroplast

129    genomes, dated with 11 outgroup species[7,11]. Principal component analysis (PCA) of

130    chloroplast variation within *A. thaliana* identified several genetic groups (Fig. 2a,

131    Supplementary Fig. 3), which is consistent with the phylogenetic tree (Supplementary Fig. 4).

132    After collapsing branches with low approximate likelihood ratio support (Supplementary Fig.

133    5), we found that the *A. thaliana* chloroplast tree exhibits a basal polytomy (Fig. 2b), with

134    seven major haplogroups branched off at roughly the same time: groups 1, 2, 8, 9, 10, and

135    two monophyletic clades: 3 + 4 and 5 + 6 + 7. While studies based on the nuclear genomes

136    showed that the African accessions contain higher polymorphism and are highly diverged

137    from Eurasian populations[5], we did not observe such pattern in chloroplast. Instead, the

138    African accessions represent only a small subset of chloroplast variation (Groups 2, 3, and 7,

139    Fig. 2b), suggesting that the chloroplast phylogeny captures more ancient demographic

140    history than the divergence between African and European populations. Indeed, molecular

141    dating with BEAST[12] confirmed that the major chloroplast haplogroups diverged at ca. 227

142    kya (95% highest posterior density 121-340 kya, Supplementary Fig. 6,7). This considerably

143    predates the inferred divergence time between African and Eurasian nuclear genomes

144    (90-120 kya)[5,6].

145    　　The spatial distribution of chloroplast haplogroups is uneven, with several highly

146    diverged haplogroups existing only in Eurasia, and the African population containing only

147    group 2, 3, and 7 (Fig. 3a). While Morocco has highest nuclear genomic variation[5], we

148    observed this only for group 3, where the variation decreases from Morocco northwards (Fig.

149    3c), suggesting its Moroccan origin and later northward migration. Group 2 is only confined

150    in Iberia and Morocco, with the former having higher variation (Fig. 3b). The monophyletic

151    clade containing group 5, 6 and 7, on the other hand, has a global distribution with highest

152    variation in Europe, especially the Balkan Peninsula (Fig. 3e). Therefore, even among the

153    three haplogroups Africa shares with Eurasia, only one of them has African population

154    containing higher variation. In summary, our observation is consistent with both hypotheses

155    of (1) African origin of *A. thaliana* and complete lineage sorting between the African and the

5

156    Eurasian accessions or (2) a Eurasian origin and dispersal into different regions

157    (southwestern Europe and northwestern Africa, Balkan and Levant, and south Asia), after

158    which most Eurasian nuclear genomic variation was wiped out by the rapidly expanding

159    non-relicpts.

160

161    **Demography of chloroplast genomes**

162         We further used Extended Bayesian Skyline Plots[13] to investigate chloroplast

163    demographic histories. Haplogroups typical of the Iberian and Moroccan region (group 2, 3,

164    and 4) showed long-term stable population size (Fig. 4). The Eurasian haplogroups 1 and 5

165    had population size increase since 20 kya, consistent with the post-glacial expansion with

166    the retreat of ice sheet. Interestingly, for the globally distributed group 7, the central Asian

167    group 8, the European group 6W, and the Chinese Yangtze group 6E, all had rapid population

168    size increase since 10 kya, a time point close to the previously inferred rapid expansion of

169    weedy non-relicts[2,14]. In addition, haplogroups 6 and 7 had highest genetic variation near

170    the Balkan Peninsula (Supplementary Fig. 8c,d), corresponding to the inferred origin of

171    non-relict expansion[2].

172         For haplogroups shared by Morocco and Europe (groups 2, 3, and 7), we further

173    investigated their demographic histories separately. While the Morocco population of all

174    three groups still exhibited long-term stable population size (Supplementary Fig. 9), the

175    European group 2 had population size increase since 20 kya (the first post-glacial expansion),

176    and the European group 7 had size increase since 10 kya (the non-relict expansion). Taken

177    together, compared to Europe, the Moroccan population was less influenced by either

178    episode of demographic change, especially the second expansion wave wiping out most

179    nuclear genomic variation across Eurasia[2]. While the fact that Morocco possesses most

180    nuclear genomic variation could be interpreted as an African origin of *Arabidopsis thaliana*[5,6],

181    the complex demographic history we showed is an equally likely explanation.

182

183    **The eastern end of non-relict expansion**

184         While all current and previous[2,14] estimates suggested the non-relicts expanded around

185    10 kya and almost all Eurasian *A. thaliana* are descendants of this population, some studies

186    estimated the population divergence time between the European, Central Asian, and

187    Chinese non-relict populations to be around 45 kya[5,6,8]. While these studies performed the

188    multiple sequentially Markovian coalescent (MSMC) estimates[15] on the whole genome, we

189    wish to note there are clear evidences that non-relict populations across Eurasia had

190    introgression from distinct and highly diverged local relicts (Table 1 and ref. 2). Using the

191    whole genome, in some genomic regions one would be comparing between relicts and

192    non-relicts or two highly diverged relict groups (e.g. between Tanzania and Morocco),

193    thereby overestimating the true divergence time between the European, Central Asian, and

194    Chinese non-relict populations. We therefore focus on a unique chromosomal translocation,

195    where the non-relicts have a charismatic derived haplotype[2,16] (Supplementary Fig. 10).

196    Since the two structural variants of the translocation cannot recombine effectively, genetic

197    variation within the derived haplotype reflects demographic history of non-relicts[2]. Using

198    only this 750 kb region, we first performed the same set of MSMC comparisons as Durvasula

199    *et al.*[5] and successfully re-created similar patterns (Supplementary Fig. 11), demonstrating

200    this region alone contains enough information to trace the demographic history. Based on

201    the derived haplotypes in this region, the European, Central Asian, and Chinese Yangtze

202    non-relict populations diverged between 5 to 15 kya (Fig. 5), consistent with the time of

203    rapid population expansion in several chloroplast haplogroups (10 kya, Fig. 4) as well as the

204    inferred timing of non-relict expansion[2,14].

205        For non-relicts, the most notable recent expansion happened in the Chinese Yangtze

206    population (Fig. 4). Assuming the North American population had a common ancestor

207    around 1600 AD[3], using simple genetic distance and assuming the same mutation rate, we

208    estimated the Chinese Yangtze population having a common ancestor at 568 AD (estimated

209    from chloroplast) or 823 AD (estimated from the chromosomal translocation region in

210    chromosome 1[2]). The time point is consistent with *A. thaliana* entering China through

211    central Asia with human activities, with the Silk Road being one possibility.

212        Given that the Yangtze population clearly had introgression from the Yunnan relicts (Fig.

213    1b and Table 1), we further investigated whether introgression contributed to local

214    adaptation of this population. We calculated the $\hat{f}_d$ statistic[17] for 50-kb windows across the

215    genome, from which gene flow between Yangtze population and Yunnan relicts was inferred

216    (Supplementary Table 1). Then we compared the results with genes under strong selection

217    in Yangtze population identified by Zou *et al.*[8] (Supplementary Table 1,2). Windows with

218    these selected genes were found enriched in the top 5% tail of $\hat{f}_d$ distribution (Fisher's

219    exact test, odds ratio = 2.049, *P* = 0.007). Genes both under strong selection and with

220    evidences of relict origin were overrepresented for gene ontology (GO) terms associated

221    with biotic interaction, immune response, and programmed cell death (Supplementary Table

222    2), while strongly selected genes without strong traces of introgression (presumably

223    representing novel mutations or standing variations within the invading ancestral Yangtze

224    non-relicts) have no enrichment of any GO term. Taken together, much similar to the

225    western end of Eurasia[2], our results suggested the ability of the expanding non-relict

226  population to colonize the eastern end of Eurasia (the Yangtze River Basin) was also greatly

227  facilitated by introgressions from local relicts.

228  Interestingly, in whole-genome phylogenetic tree the Yangtze population has long

229  branches relative to other non-relict populations (Fig. 1a). To test whether this is caused by

230  introgression from a highly diverged group (the Yunnan relicts) or natural selection

231  accelerating the fixation of novel mutations in some genomic regions, we excluded the top

232  20% windows with highest introgression (Supplementary Fig. 12a), any window containing

233  positively selected genes (Supplementary Fig. 12b), or both (Fig. 12c). These trees remain

234  similar to the whole-genome tree where the Yangtze population still has long branch length.

235  It is likely that the Yangtze population exhibits higher mutation rate or more rapid life cycle

236  resulting in more than one generation per year, and both hypotheses need to be formally

237  tested. If so, time to the most recent common ancestor of Yangtze accessions would be

238  more recent than our estimation.

239

240  **Discussion**

241

242  **On ancient population structure**

243  Combining currently available data from genome resequencing projects of *Arabidopsis*

244  *thaliana*, here we revisit demographic history of *A. thaliana* from the global perspective. The

245  "out of Africa" hypothesis states that the African population first separated into the

246  Moroccan, Levantine, and Southeast African groups at ca. 90 kya followed by a migration

247  event from Levant into Eurasia[5,6]. However, we observed that the Chinese Yunnan

248  accessions are genetically closer to the Tanzanian and South African group than to any other

249  group, which suggests more than one "out of Africa" events if the ancestral population is

250  originated from Africa. For chloroplast, we observed several highly diverged haplogroups

251  existing only in Eurasia, and Africa contains merely a subset of overall chloroplast variation,

252  which hints that the ancestral population may not originate from Africa. Together, these

253  results suggest another demographic scenario that is as possible as the "out of Africa"

254  model (Fig. 6): Like all other species in the *Arabidopsis* genus, ancient *Arabidopsis thaliana*

255  originated in temperate Eurasia and separated into the Moroccan/Iberian, Levantine, and

256  South/Southwest Asian groups at ca. 90 kya. Later the Moroccan/Iberian and Levantine

257  group migrated northwards into Eurasia while the Asian group dispersed into Tanzania and

258  South Africa. At around 10 kya, the weedy non-relict group expanded across Eurasia. On the

259  other hand, we acknowledge while the existence of more ancient chloroplast variation in

260  Eurasia might indicate a Eurasian origin of *A. thaliana*, it is also likely that the ancient

261 variation once existed in Africa but was later lost due to strong bottleneck events or

262 selective sweeps favoring a few chloroplast haplogroups. Chloroplast demography, however,

263 shows long-term stable population size in Morocco for all haplogroups (Supplementary Fig.

264 9), thus not lending strong support to the Moroccan chloroplast bottleneck or sweep

265 scenario.

266 　　　For both the "out of Africa" and the "into Africa" models, the most notable

267 demographic turnover event is the recent replacement of many Eurasian relict populations

268 by the weedy "non-relicts"[1,2], which expanded along the east-west axis of Eurasia and left

269 more relict genomic fragments in southern and northern Europe. Interestingly, this is

270 supported by an independent study: Exposito-Alonso et al.[18] found that the same alleles

271 increasing survival under extreme drought are enriched in the relict accessions and are

272 concentrated in northern and southern Europe, a pattern predicted by the east-west

273 non-relit expansion. The drastic demographic turnover is also responsible for only few

274 private variants being observed in each regional Eurasian population[5], which is a logical

275 outcome since most Eurasian genomes descended only recently from a single population[1,2].

276 Meanwhile, the African population remained relatively isolated from Eurasia and retained

277 much of the ancestral nuclear-genome variation. Therefore, our results suggest that the

278 current patterns of global nuclear-genome variation (Africa containing most variation) is a

279 consequence of non-relicts wiping out most ancient genetic variation in Eurasia[2], which is

280 compatible with both scenarios about ancient A. thaliana population structure (Fig. 6).

281 　　　While no relict accession was found in northern central Asia, this region is enriched for

282 the ancestral haplotype of the chromosome 1 translocation (Supplementary Fig 10), and

283 these ancestral haplotypes form a unique genetic group when comparing to worldwide

284 ancestral haplotypes[2]. In the present study, several Eurasia-only chloroplast haplogroups

285 (group 8, 9, 10) also exist in this region (Fig 3). We therefore suspect another unique relict

286 group might have existed near this region, which later became very rare or extinct.

287

288 **On the recently established Chinese population**

289 　　　In the process of rapid expansion, populations in the expansion front constantly

290 encounter novel environments, which may impede the speed and extent of expansion.

291 Hence, how can a population rapidly spread across a wide geographical and environmental

292 range, and what is the source of adaptation to these drastically different environments? We

293 therefore focus on the origin and adaptation of Chinese Yangtze population for the second

294 part of our investigation. We show that the Yangtze population originated no more than

295 2000 years ago and spread rapidly across the basin. Using the properly rooted phylogenetic

296 tree, we showed that Yangtze population belongs to the non-relict group and are genetically

297 the closest to Central Asian non-relicts (Fig. 1).

298  Zou *et al.*[8] performed genome-wide scans for signal of selection in the Yangtze

299 population. Here we also investigated this result in the context of introgression from Yunnan

300 relicts. We found that selected genes are enriched with signs of relict introgression, and

301 genes with both signs of selection and introgression are overrepresented for

302 immune-related functions. On the other hand, selected genes without signs of introgression

303 do not have any significant gene ontology enrichment. Our results therefore suggest, among

304 the various aspects of adaptation to the novel Yangtze River Basin environment, the

305 adaptation to immune-related biotic stress is associated with gene flow with local relicts,

306 which might have co-existed with local pathogens for a long time. Interestingly, for the

307 western end of non-relict expansion in Iberia, relict introgression likely contributed to the

308 adaptation to abiotic factor, as highly introgressed genes in Iberian non-relicts are enriched

309 for GO terms including root development and ion metal transmembrane activity[2]. In the end,

310 how can the non-relicts, a population near the Balkans, occupy such broad environmental

311 gradient spanning more than 10,000 km across Eurasia within 10,000 years? While the

312 mal-adaptation to novel environments in the expansion front may impede non-relict spread,

313 our results suggest non-relicts frequently assimilated the biological distinctiveness of locally

314 adaptive relicts. Together with human's long-term disturbance of native Eurasian vegetation

315 and non-relicts' association with anthropogenically disturbed habitats[1,2], the environmental

316 resistance to non-relict expansion appears futile in most of Eurasia.

317

318

319 **Materials and Methods**

320

321 **Data source and SNP identification of nuclear genome**

322  In this study, we obtained *Arabidopsis thaliana* data from the 1,135 worldwide

323 genomes[1], 73 African genomes[5], 116 Chinese genomes[8] and one *Arabidopsis lyrata* sample

324 (SRR2040792)[7]. Reads were trimmed based on quality using SolexaQA[19], and possible

325 remaining adaptor sequences were removed with cutadapt[20]. Reads were mapped to the

326 TAIR 10 reference genome using BWA 0.7.15[21]. Picard Tools

327 (http://broadinstitute.github.io/picard) were used to mark duplicated read pairs, and the

328 genotypes of each site in each accession (including non-variant sites and SNPs) were called

329 following GATK 3.7 best practice[22].

330    We further filtered the SNPs with QUAL < 100, QD < 20, call rate < 0.99, DP < 3 or > 2

331    standard deviations from genome-wide average depth and removed 2 Chinese accessions

332    (SRR2204178, SRR2204343) with high missing rate, resulting in 5,915,870 SNPs and 1323

333    accessions.

334

335    **Alignment of *A. thaliana* population data with outgroups**

336    In addition to the *Arabidopsis thaliana* reference chloroplast genome (NC000932), we

337    obtained the outgroup chloroplast genomes from the genera *Arabidopsis*, *Capsella*, and

338    *Camelina*[7,11]: *Arabidopsis lyrata* subsp. *petraea* (LT161948), *Arabidopsis lyrata* subsp. *lyrata*

339    (LN877383), *Arabidopsis halleri* subsp. *halleri* (LN877382), *Arabidopsis carpatica* (LT161918),

340    *Arabidopsis arenosa* subsp. *arenosa* (LT161904), *Arabidopsis nitida* (LT161970), *Arabidopsis*

341    *pedemontana* (LN877384), *Arabidopsis cebennensis* (LN877381), *Capsella rubella*

342    (LN877385), *Capsella bursa-pastoris* (NC_009270), and *Camelina sativa* (LN877386).

343    All twelve chloroplast genomes were annotated with Verdant[23], and about 90 protein

344    coding genes were identified in each sequence. We retained genes existing in all species and

345    excluded those within the two inverted repeat regions, resulting in 67 orthologous genes

346    with one-to-one relationship in all species. The 67 genes were separately aligned with

347    MUSCLE 3.8.31[24]. Based on this alignment of *A. thaliana* reference chloroplast genome with

348    outgroups species, we used custom R scripts to "paste" the Illumina-based *A. thaliana*

349    accession data[1,5,8] onto the among-species alignment. All following analyses were based on

350    this concatenated dataset of 67 protein coding genes.

351    To remove possible confounding effect from heteroplasmy, we excluded accessions

352    with heterozygous SNPs > 0.5% of all SNPs. For all remaining accessions, any heterozygous

353    genotype call was transformed to missing data, and accessions with > 1% missing data

354    among all SNPs were excluded.

355

356    **Patterns of chloroplast and chromosome 1 translocation polymorphism**

357    For chloroplast, bi-allelic SNPs of 67 protein coding genes were called using vcftools[25]

358    after conversion of fasta alignment format to vcf format in TASSEL[26]. Focusing on bi-allelic

359    sites with no missing data and zero heterozygosity, we identified 760 sites polymorphic

360    within *A. thaliana* and 2124 fixed between *A. thaliana* and any one outgroup species. SNPs

361    were identified in chromosome 1 translocation following procedures of previous studies[2,16].

362    PCA was done with R package adegenet[27] and visualized with R package ggplot2[28].

363    Accessions of *Arabidopsis thaliana* were assorted into groups (geo-clusters) according

364    to geographical location where they were sampled (Supplementary Table 3). Diversity of

365    each chloroplast within each geo-cluster was then estimated by pair-wise genetic distance[29].

366    To account for uneven sampling, the average genetic variation of 100 resampling trials was

367    obtained for each geo-cluster. For each re-sampling trial, 100 samples were randomly drawn

368    with replacement within each geo-cluster. Groups with less than 3 samples in each geo

369    cluster were ignored. All the calculation and plots were completed in R with customize

370    scripts and package ggplot2[28]. Pie charts were plotted with R package ggplot2[28] as well.

371

**Phylogenetic reconstruction and divergence time estimation**

373    5,915,870 nuclear SNPs of 1323 accessions (including *Arabidopsis lyrata*) were used to

374    construct nuclear neighbor-joining tree. The pair-wise distance is calculated by dividing

375    number of SNP difference between pairs with total number of non-missing sites that are

376    polymorphic within 1323 accessions.

377    1312 chloroplast haplotypes of 2124 bi-allelic SNPs were converted into phylip format

378    and submitted to maximum-likelihood-based phyML 3.0[30] for reconstruction of phylogenetic

379    tree. Substitution model selection was done by SMS[31] with Bayesian Information Criterion

380    (BIC). Subtree pruning regrafting (SPR) was used as tree searching algorithm and the branch

381    support was estimated by approximate likelihood ratio test (aLRT SH-like). Branches with low

382    support (aLRT < 0.5) were collapsed with TreeGraph2[32]. The collapsed maximum likelihood

383    tree was visualized and colored in FigTree version 1.4.3

384    (http://tree.bio.ed.ac.uk/software/figtree/).

385    Divergence time among haplogroups was estimated with BEAST version 2.5.0[12]. The

386    collapsed ML tree of 426 unique chloroplast haplotypes was constructed as described

387    previously and served as starting tree for BEAST after it was converted ultrametric and had

388    the node ages fit within constraints of calibration points using R package ape[33]. Three

389    calibration points estimated in previous studies[7,11] were adopted, the root height

390    (divergence time between genus *Arabidopsis* and *Camelina*, *Capsella*) was set to 8.16 Mya;

391    divergence time between genus *Capsella* and *Camelina* was set to 7.36 Mya; and the

392    divergence time between *Arabidopsis thaliana* and other species in *Arabidopsis* genus was

393    set to 5.97 Mya. Normal distribution with 1 mya standard deviation was used for all the

394    three calibration points.

395    Two independent MCMC runs with $5 \times 10^8$ chain length were generated with Calibrated

396    Yule Model for priors and Relaxed Clock Log Normal for clock model. GTR was chosen as site

397    model according to SMS. Parameters of MCMC trees were sampled every $5 \times 10^4$ generations

398    and submitted to Tracer version 1.6 (http://tree.bio.ed.ac.uk/software/tracer/) for quality

399    control of MCMC chains. LogCombiner[12] was implemented to combine two independent

400    runs. A maximum clade credibility (MCC) tree was constructed using 18000 output trees of

401    LogCombiner with 10% burn-in in TreeAnnotator[12]. The MCC tree was then visualized and

402    colored in FigTree version 1.4.3 (http://tree.bio.ed.ac.uk/software/figtree/).

403

404    **ABBA-BABA and $\hat{f}_d$ estimation**

405         Python and R scripts were downloaded from

406    (https://github.com/simonhmartin/genomics_general) for ABBA-BABA and $\hat{f}_d$ estimation[17].

407    Genome-wide D statistics were estimated in R followed the instruction of

408    (http://evomics.org/learning/population-and-speciation-genomics/2018-population-and-sp

409    eciation-genomics/abba-baba-statistics/), West European population (EU) was treated as

410    Pop1, Yangtze population (YA) was treated as Pop2, Yunnan (YU) (the less admixed accession)

411    was treated as Pop3 and *Arabidopsis lyrata* was treated as outgroup. Sliding window analysis

412    of $\hat{f}_d$ estimation between Yangtze population and Yunnan (the less admixed accession) was

413    done using the Python scripts. Window size was set to 50 kb with 20 kb step, each window

414    containing at least 100 SNPs. Windows with top 5 % highest $\hat{f}_d$ values were viewed as

415    regions with strong introgression between Yangtze and Yunnan. The selected genes in

416    Yangtze population[8] within/outside the window of introgression were then submitted to

417    agriGO v2.0[34] for gene ontology analysis.

418

419    **Multiple sequentially Markovian coalescent analysis (MSMC)**

420         Relative cross coalescence rate was estimated using MSMC v2[15]. Sequences of

421    chromosome 1 translocation (Chr1:20271447-21032307) were used as input. Generation

422    time was set at 1 year and mutation rate of $7.1 \times 10^{-9}$ was assumed according to previous

423    studies[5,8]. Results were then plotted in R.

424

425    **Extended Bayesian skyline plot**

426         Extended Bayesian Skyline analysis was done using BEAST v2.5.0[12] with fixed number of

427    2124 chloroplast genome polymorphisms. Parameters were set in BEAUti 2[12], HYK

428    substitution model with empirical frequency was chosen and the clock model was set to

429    strict clock with clock rate estimated by previous Calibrated Yule Model (0.0223). Priors were

430    set to Coalescent Extended Bayesian Skyline with 0.5 population model factor and default

431    value for the rest of parameters. Sufficient length of MCMC chains were run to achieve

432    acceptable ESS values, which indicates the model is well-mixed. The ESS values were

433    estimated in Tracer v1.6.0[12], and the results were plotted in R.

434

13

435

436

437

**Acknowledgements**

We thank Thomas Mitchell-Olds and the Lee lab member for comments to this manuscript. We thank all researchers who have made the *Arabidopsis* genetic resources and data publicly available. We are grateful to Computer and Information Networking Center, National Taiwan University for the support of high-performance computing facilities. This work is supported by the Ministry of Science and Technology of Taiwan (105-2311-B-002-040-MY2 and 107-2636-B-002-004 to CRL).

**Author contributions**

Designed the study: CRL. Analyzed data: CWH, CRL, CYL. Wrote the paper: CRL, CWH.

**Conflict of interest statement**

The authors declare no conflict of interest

**Data Accessibility**

All data were downloaded from public database.

454

455

References

1.  1,135 Genomes Reveal the Global Pattern of Polymorphism in Arabidopsis thaliana. *Cell* 166, 481–491 (2016).

2.  Lee, C.-R. *et al.* On the post-glacial spread of human commensal Arabidopsis thaliana. *Nat. Commun.* 8, 14458 (2017).

3.  Exposito-Alonso, M. *et al.* The rate and potential relevance of new mutations in a colonizing plant lineage. *PLoS Genet.* 14, e1007155 (2018).

4.  Fransz, P. *et al.* Molecular, genetic and evolutionary analysis of a paracentric inversion in Arabidopsis thaliana. *Plant J.* 88, 159–178 (2016).

5.  Durvasula, A. *et al.* African genomes illuminate the early history and transition to selfing in Arabidopsis thaliana. *Proc. Natl. Acad. Sci. U.S.A.* 114, 5213–5218 (2017).

6.  Fulgione, A. & Hancock, A. M. Archaic lineages broaden our view on the history of Arabidopsis thaliana. *New Phytol.* 21, 1877 (2018).

7.  Novikova, P. Y. *et al.* Sequencing of the genus Arabidopsis identifies a complex history of nonbifurcating speciation and abundant trans-specific polymorphism. *Nat. Genet.* 48, 1077–1082 (2016).

8.  Zou, Y.-P. *et al.* Adaptation of Arabidopsis thaliana to the Yangtze River basin. *Genome Biol.* 18, 239 (2017).

9.  Slon, V. *et al.* The genome of the offspring of a Neanderthal mother and a Denisovan father. *Nature* 531, 504 (2018).

10. Alexander, D. H., Novembre, J. & Lange, K. Fast model-based estimation of ancestry in unrelated individuals. *Genome Res.* 19, 1655–1664 (2009).

11. Hohmann, N., Wolf, E. M., Lysak, M. A. & Koch, M. A. A Time-Calibrated Road Map of Brassicaceae Species Radiation and Evolutionary History. *Plant Cell* 27, 2770–2784 (2015).

12. Bouckaert, R. *et al.* BEAST 2: a software platform for Bayesian evolutionary analysis. *PLoS Comput Biol* 10, e1003537 (2014).

13. Heled, J. & Drummond, A. J. Bayesian inference of population size history from multiple loci. *BMC Evol. Biol.* 8, 289 (2008).

14. François, O., Blum, M. G. B., Jakobsson, M. & Rosenberg, N. A. Demographic history of european populations of Arabidopsis thaliana. *PLoS Genet.* 4, e1000075 (2008).

15. Schiffels, S. & Durbin, R. Inferring human population size and separation history from multiple genome sequences. *Nat. Genet.* 46, 919–925 (2014).

16. Long, Q. *et al.* Massive genomic variation and strong selection in Arabidopsis thaliana

491        lines from Sweden. *Nat. Genet.* 45, 884–890 (2013).

492   17.   Martin, S. H., Davey, J. W. & Jiggins, C. D. Evaluating the use of ABBA-BABA statistics

493        to locate introgressed loci. *Mol Biol Evol* 32, 244–257 (2015).

494   18.   Exposito-Alonso, M. *et al.* Genomic basis and evolutionary potential for extreme

495        drought adaptation in Arabidopsis thaliana. *Nat Ecol Evol* 2, 352–358 (2018).

496   19.   Cox, M. P., Peterson, D. A. & Biggs, P. J. SolexaQA: At-a-glance quality assessment of

497        Illumina second-generation sequencing data. *BMC Bioinformatics* 11, 485 (2010).

498   20.   Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing

499        reads. *EMBnet.journal* 17, 10 (2011).

500   21.   Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler

501        transform. *Bioinformatics* 25, 1754–1760 (2009).

502   22.   McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for

503        analyzing next-generation DNA sequencing data. *Genome Res.* 20, 1297–1303 (2010).

504   23.   McKain, M. R., Hartsock, R. H., Wohl, M. M. & Kellogg, E. A. Verdant: automated

505        annotation, alignment and phylogenetic analysis of whole chloroplast genomes.

506        *Bioinformatics* 33, 130–132 (2017).

507   24.   Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high

508        throughput. *Nucleic Acids Res.* 32, 1792–1797 (2004).

509   25.   Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* 27, 2156–2158

510        (2011).

511   26.   Bradbury, P. J. *et al.* TASSEL: software for association mapping of complex traits in

512        diverse samples. *Bioinformatics* 23, 2633–2635 (2007).

513   27.   Jombart, T. adegenet: a R package for the multivariate analysis of genetic markers.

514        *Bioinformatics* 24, 1403–1405 (2008).

515   28.   Wickham, H. ggplot2. *Wiley Interdisciplinary Reviews: Computational Statistics* 3,

516        180–185 (2011).

517   29.   Nei, M. Genetic Distance between Populations. *The American Naturalist* 106, 283–

518        292 (1972).

519   30.   Guindon, S. *et al.* New Algorithms and Methods to Estimate Maximum-Likelihood

520        Phylogenies: Assessing the Performance of PhyML 3.0. *Systematic Biology* 59, 307–

521        321 (2010).

522   31.   Lefort, V., Longueville, J.-E. & Gascuel, O. SMS: Smart Model Selection in PhyML. *Mol

523        Biol Evol* 34, 2422–2424 (2017).

524   32.   Stöver, B. C. & Müller, K. F. TreeGraph 2: combining and visualizing evidence from

525        different phylogenetic analyses. *BMC Bioinformatics* 11, 7 (2010).

526  33.  Paradis, E., Claude, J. & Strimmer, K. APE: Analyses of Phylogenetics and Evolution in

527      R language. *Bioinformatics* 20, 289–290 (2004).

528  34.  Tian, T. *et al.* agriGO v2.0: a GO analysis toolkit for the agricultural community, 2017

529      update. *Nucleic Acids Res.* 45, W122–W129 (2017).

530

531

532 **Figure legend**

533

534 Fig. 1. Differentiation of *Arabidopsis thaliana* nuclear genomes. (a) Neighbor-Joining tree. (b)

535 K=2 ADMIXTURE result.

536

537 Fig. 2. Differentiation of *Arabidopsis thaliana* chloroplast genomes. (a) Principal component

538 analysis. (b) Maximum likelihood cladogram, where branches with low aLRT support were

539 collapsed. Group 6 (East) consists of samples from Yangtze River Basin, China and Kashmir,

540 India. Group 6 (West) consists of samples in Eurasia.

541

542 Fig. 3. Geographical distribution and spatial genetic variation of chloroplast haplogroups. (a)

543 Diversity map of chloroplast haplogroups. Pie charts show the proportion of each group, and

544 chart size is proportional to sample size. (b), (c), (d), (e) are polymorphism maps correspond

545 to group 2, 3, 4, 5+6+7 respectively. The diameter of each circle is proportional to mean

546 pair-wise genetic distance of each geographical region.

547

548 Fig. 4. Variation of population size over time inferred from chloroplast polymorphism.

549 Extended Bayesian Skyline is plotted for each chloroplast haplogroup except group 10, which

550 has small sample size.

551

552 Fig. 5. Timing of population splits inferred from chromosome 1 translocation. Relative cross

553 coalescence rate (CCR) between populations is shown. EU: Western Europe, CA: Central Asia,

554 YA: Yangtze, TZ: Tanzania, TZSA: Tanzania and South Africa. Decrease of CCR from 1.0

555 indicates population split, steep slope of CCR from 1.0 to 0.0 indicates drastic and complete

556 isolation while mild one indicates slow and progressive isolation. (a) 4 haplotypes MSMC

557 that has better estimation of older splits. (b) 8 haplotypes MSMC that has better estimation

558 of more recent splits.

559

560 Fig. 6. Two scenarios of demographic history consistent with the present-day pattern of

561 spatial genetic variation in *Arabidopsis thaliana*. Scenario 1 represents the "Out of Africa"

562 model: Ancestral population of *Arabidopsis thaliana* in Africa split into 3 populations at ca.

563 90 kya, the Moroccan, Levantine and Sub-saharan African. Later, Moroccan expanded into

564 Europe through Iberia, Levantine dispersed into Central Asia and Europe while Sub-saharan

565 African migrated into Yunnan possibly through Southwest and South Asia. Scenario 2

566 represents the "Into Africa" model: Ancestral population of *Arabidopsis thaliana* in Europe

18

567  split into 3 populations, the Moroccan/Iberian, Levantine and a South/Southwest Asian

568  population. Later, Moroccan/Iberian expanded northwards and eastwards into Europe,

569  Levantine dispersed into Central Asia and Europe while the South/Southwest Asian migrated

570  into Yunnan and Sub-Saharan Africa. Since 10 kya, the weedy non-relicts from Balkan and

571  Eastern Europe spread rapidly westwards into Iberia and eastwards into Yangtze River Basin

572  of China, wiping out genetic variation along the way while obtaining adaptive genes through

573  gene flow between local relicts.

574

575  Supplementary Fig. 1. Geographical distribution of nuclear genetic variation in *Arabidopsis*

576  *thaliana*. The color of dots corresponds to the Neighbor-Joining tree in Figure 1a.

577

578  Supplementary Fig. 2. Distribution of nuclear genome heterozygosity of 1322 *Arabidopsis*

579  *thaliana* accessions.

580

581  Supplementary Fig. 3. Differentiation of *Arabidopsis thaliana* chloroplast genomes. (a) PC3

582  and PC4. (b) PC5 and PC6. Group 6 (East) consists of samples from Yangtze River Basin, China

583  and Kashmir, India. Group 6 (West) consists of samples in Eurasia.

584

585  Supplementary Fig. 4. Chloroplast uncollapsed maximum likelihood phylogram. Group 6

586  (East) consists of samples from Yangtze River Basin, China and Kashmir, India. Group 6 (West)

587  consists of samples in Eurasia. Note that this is an uncollapsed bifurcating tree. Some

588  internal branches are too short to be clearly visible. These branches also tend to have

589  extremely low branch support.

590

591  Supplementary Fig. 5. Chloroplast uncollapsed maximum likelihood cladogram with aLRT

592  branch support. Branches were colored according to chloroplast haplogroups defined in

593  Figure 2.

594

595  Supplementary Fig. 6. Chloroplast BEAST dated tree. Node values represent mean height of

596  divergence time in mya. Branches were colored according to chloroplast haplogroups

597  defined previously.

598

599  Supplementary Fig. 7. Chloroplast BEAST dated tree. Node values represent 95% highest

600  posterior density range of divergence time in mya. Branches were colored according to

601  chloroplast haplogroups defined previously.

602

603 Supplementary Fig. 8. Spatial genetic variation of chloroplast haplogroups. (a), (b), (c), (d),

604 (e), (f), (g) are polymorphism maps correspond to group 1, 5, 6, 7, 8, 9, 10 respectively. The

605 diameter of each circles is proportional to mean pair-wise genetic distance of each

606 geographical region.

607

608 Supplementary Fig. 9. Variation of population size over time inferred from chloroplast

609 polymorphism. Extended Bayesian Skyline is plotted for European and Moroccan population

610 of chloroplast (a) haplogroup 2, (b) haplogroup 3 and (c) haplogroup 7.

611

612 Supplementary Fig. 10. Genetic differentiation of chromosome 1 translocation. (a) Principal

613 component analysis. (b) Geographical distribution. Red dots are accessions with the

614 ancestral haplotype, and blue are accessions with the rearranged derived haplotype.

615

616 Supplementary Fig. 11. Reproducing 4-haplotype MSMC results of Durvasula *et al*. (2017)

617 using chromosome 1 translocation instead of whole genome. Relative cross coalescence rate

618 (CCR) between populations is shown: EU: West Europe, CA (ancestral): Central Asia

619 accession with ancestral allele of chromosome 1 translocation, MO: Morocco, TZ: Tanzania,

620 SA: South Africa. Decrease of CCR from 1.0 indicates population split, steep slope of CCR

621 from 1.0 to 0.0 indicates drastic and complete isolation while mild one indicates slow and

622 progressive isolation.

623

624 Supplementary Fig. 12. Nuclear Neighbor-Joining tree built (a) without SNPs located in

625 regions of top 20% $\hat{f}_d$, (b) without SNPs located in regions containing selected genes of

626 Yangtze population and (c) both.

627

628

629 Table 1. Results from the ABBA-BABA test in the form of (((P1,P2),P3),O) where O is the

630 outgroup *Arabidopsis lyrata* [a]

631

| Test | P1 | P2 | P3 | *D* statistic | *Z* score | *P* value |
|------|------|----------|----------|------|--------|-----------|
| A | TZSA | YU-admix | EU | 0.172 | 10.330 | 5.17E-25 |
| A | TZSA | YU-pure | EU | 0.011 | 0.614 | 0.539 |
| B | EU | YA | YU-pure | 0.081 | 3.967 | 7.27E-05 |
| B | TZSA | YU-pure | YA | 0.063 | 3.023 | 0.003 |
| C | EU | YA | YU-admix | 0.256 | 14.562 | 4.92E-48 |
| C | TZSA | YU-admix | YA | 0.365 | 21.756 | 6.05E-105 |

632

633 a. EU: Western European non-relicts. TZSA: Tanzanian and South African relicts. YA: Chinese

634 Yangtze River Basin non-relicts. YU-admix: The more-admixed Yunnan relict with high

635 heterozygosity. YU-pure: The less-admixed Yunnan relict.

636

21

Fig. 1

a



Iberian non-relict

Canary island
& Iberian relict

Europe,
America

Moroccan, Cape Verde island,
Italian & Levantine relict

*A. lyrata*

Yunnan, Tanzanian
& South African relict

Central Asia

Yangtze

b



Yangtze

Yunnan

Tanzania

South Africa

Fig. 2

a



1.0

Group 1
Group 2
Group 3
Group 4
Group 5
Group 6
Group 7
Group 8
Group 9
Group 10

b

Group 6 (East)    Group 5

Group 7

Group 6 (West)

Cap,Cam
outgroup

226.5
kya

Ara
outgroup

170.5
kya

Group 9

Group 10

Group 8

Group 1

Group 2

Moroccan

Group 3

Sub-saharan
African

Group 4

Fig. 3



Cp Group
1 2 3 4 5 6 7 8 9 10

Fig. 4

Fig. 5

Fig. 6



Scenario 1: Out of Africa

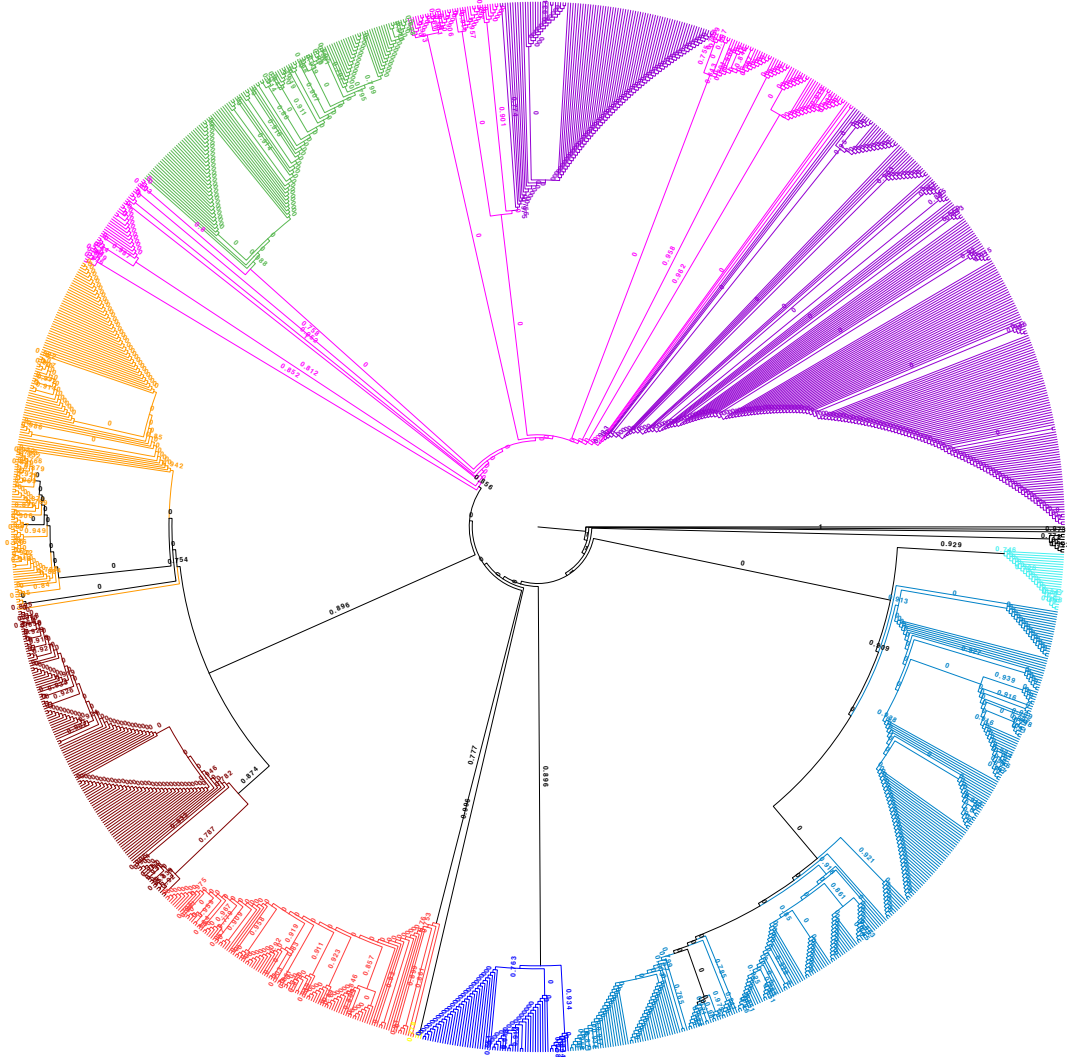Scenario 2: Into Africa

Supplementary Fig. 1
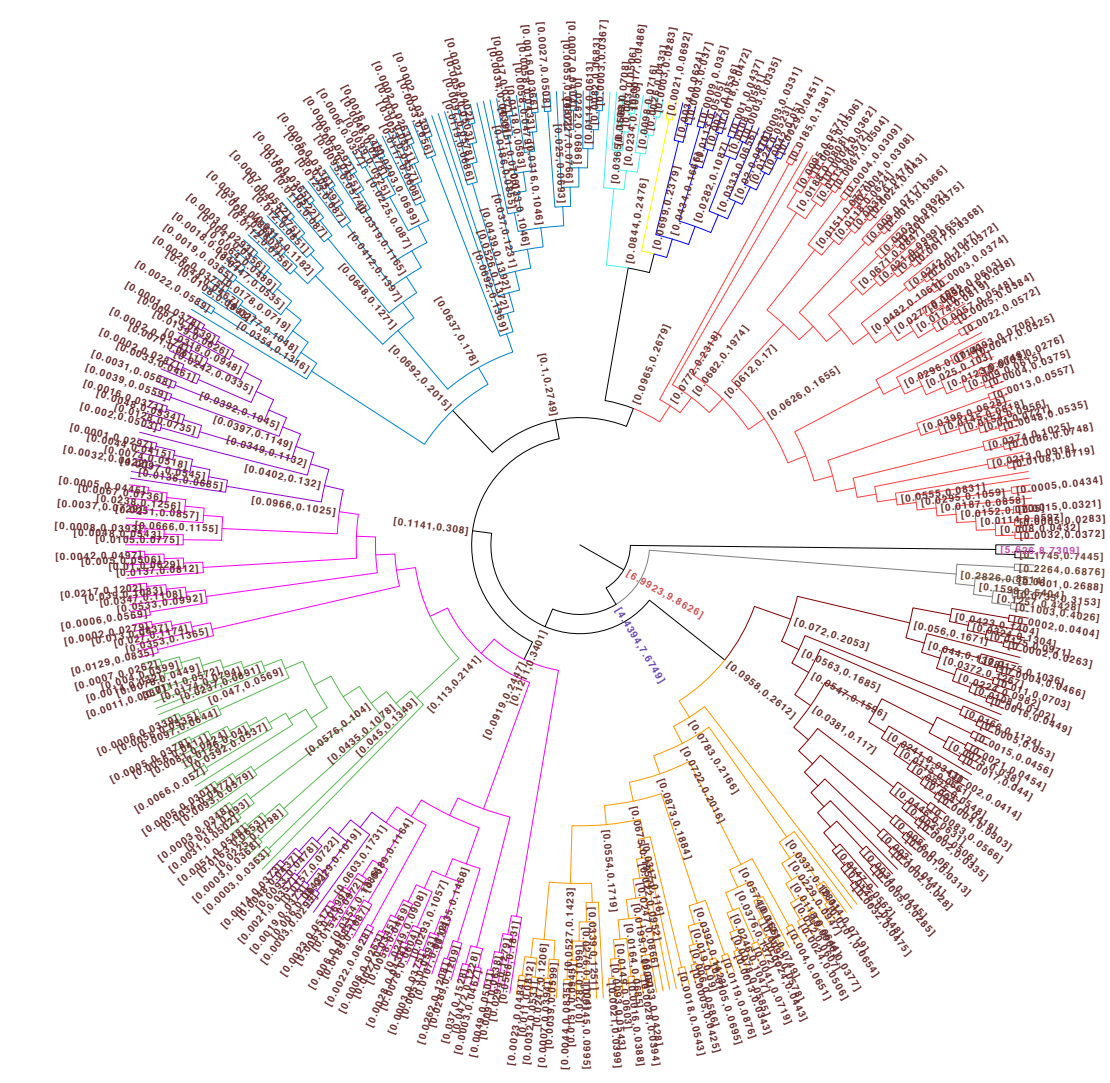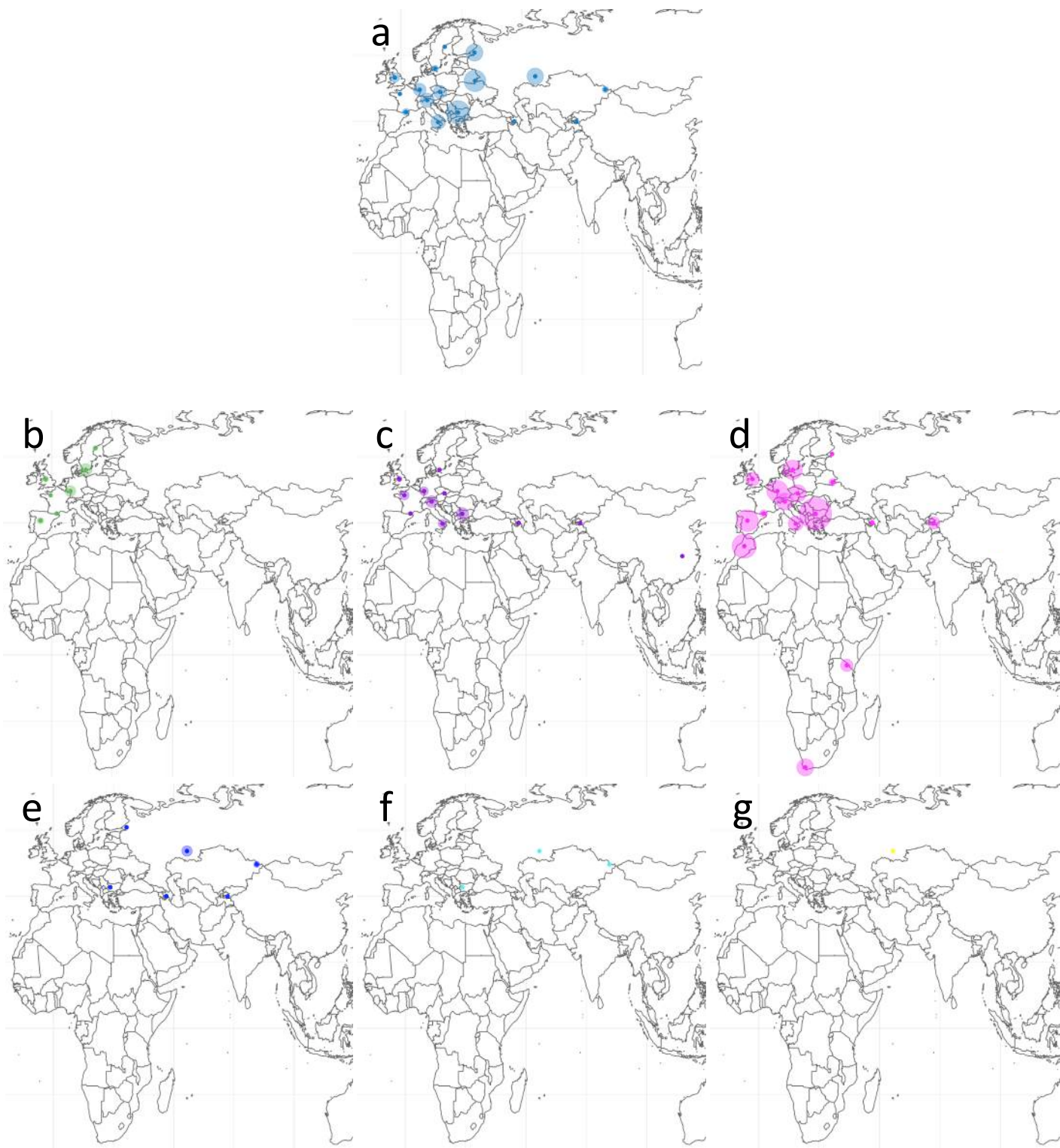
Supplementary Fig. 2

Supplementary Fig. 3

# Supplementary Fig. 4



Group 7

Group 5

Group 6 (East)

Group 7

Group 7

Group 6 (West)

Group 3

Outgroups

Group 9

Group 4

Group 1

Group 8

Group 2

Group 10

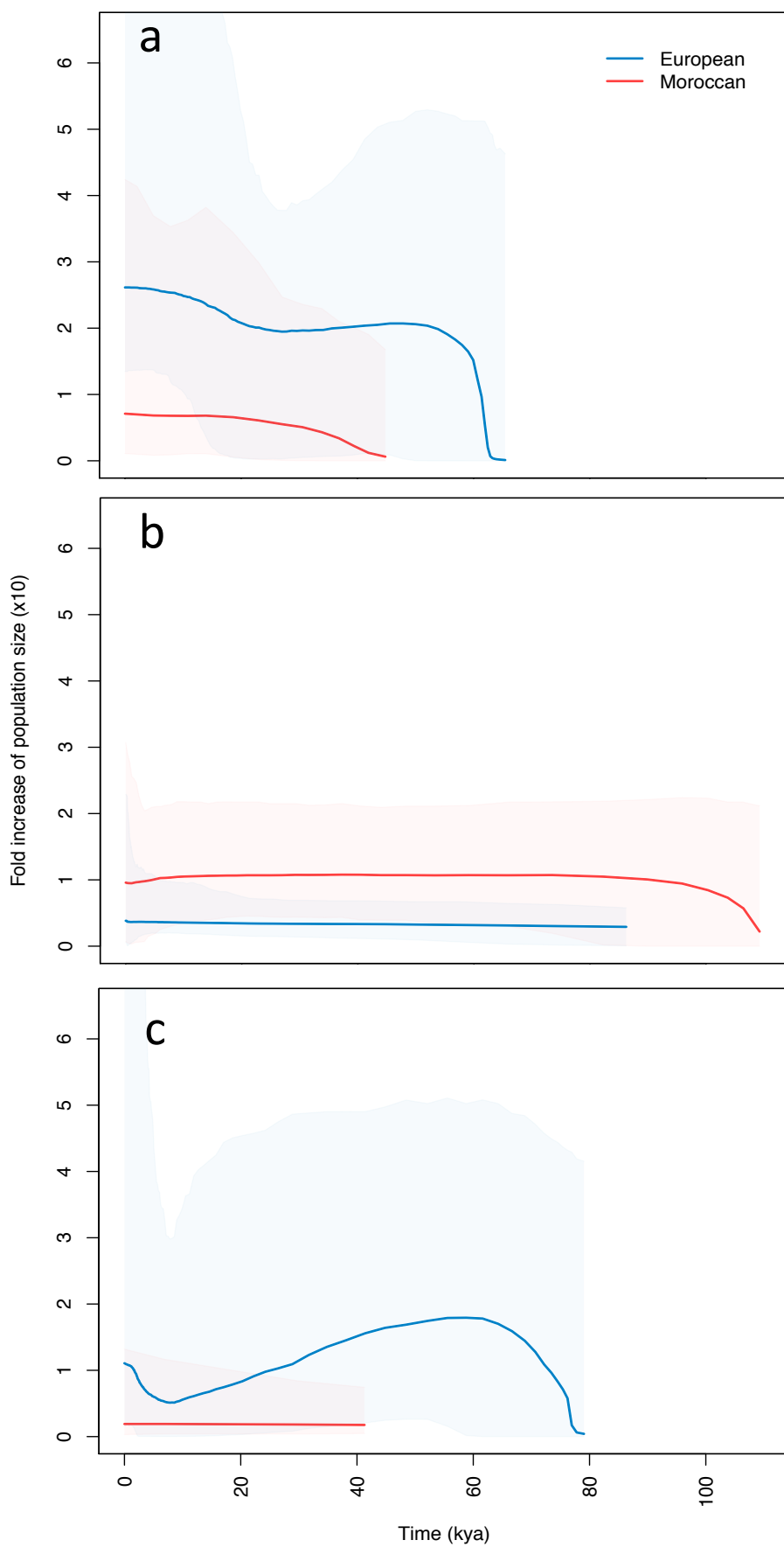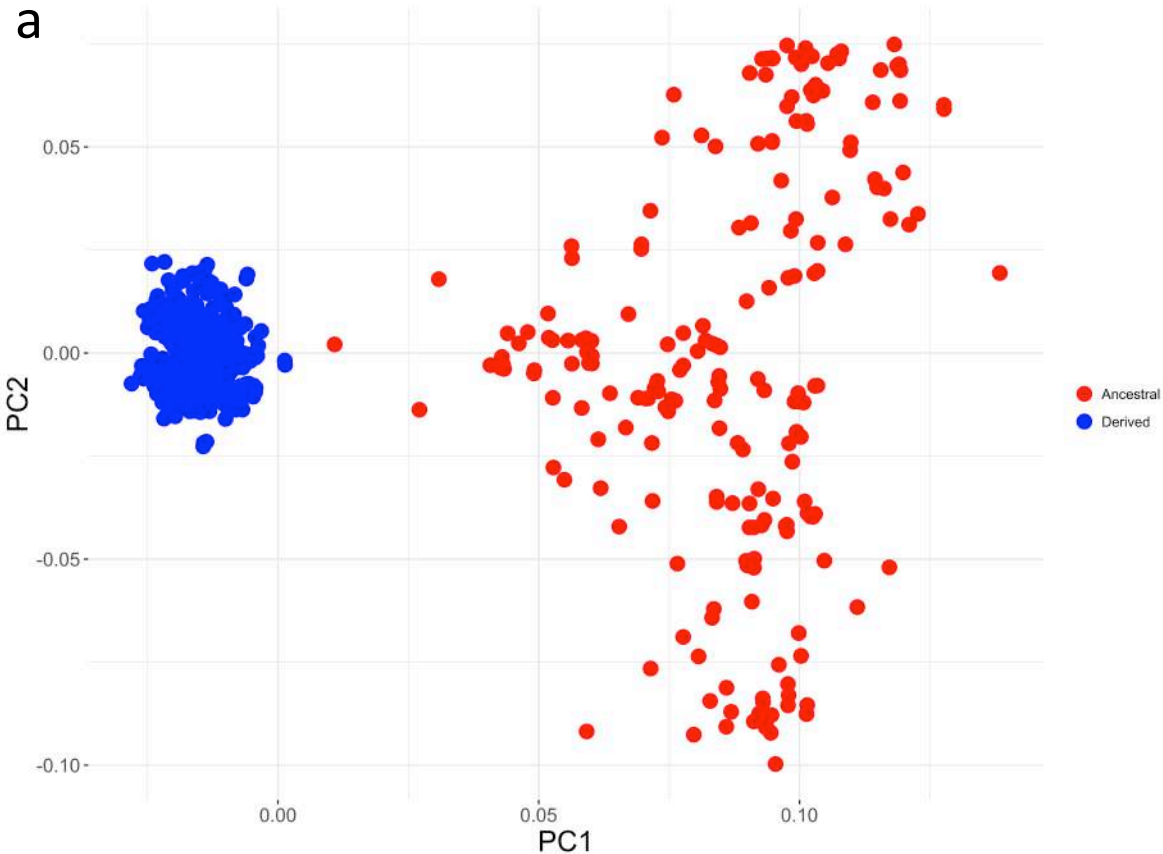◆ Moroccan

▲ Sub Saharan African

Supplementary Fig. 6

Supplementary Fig. 7
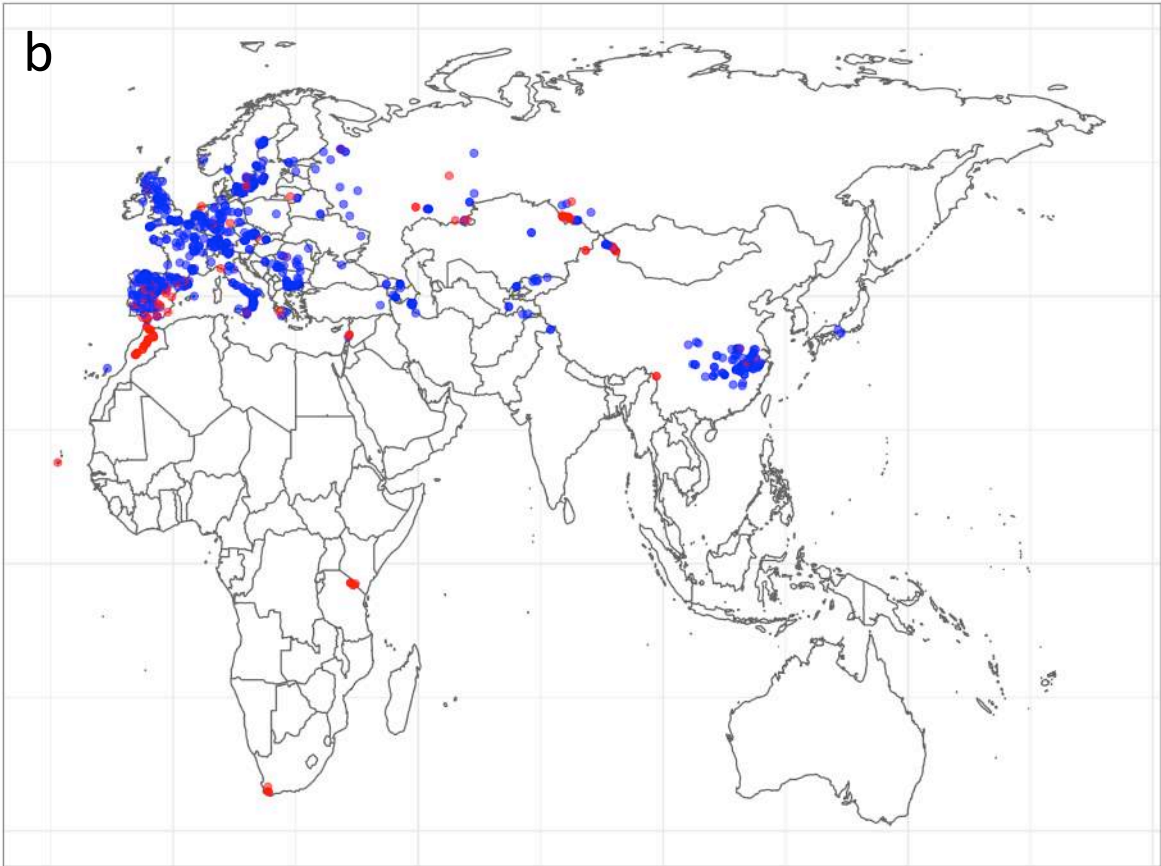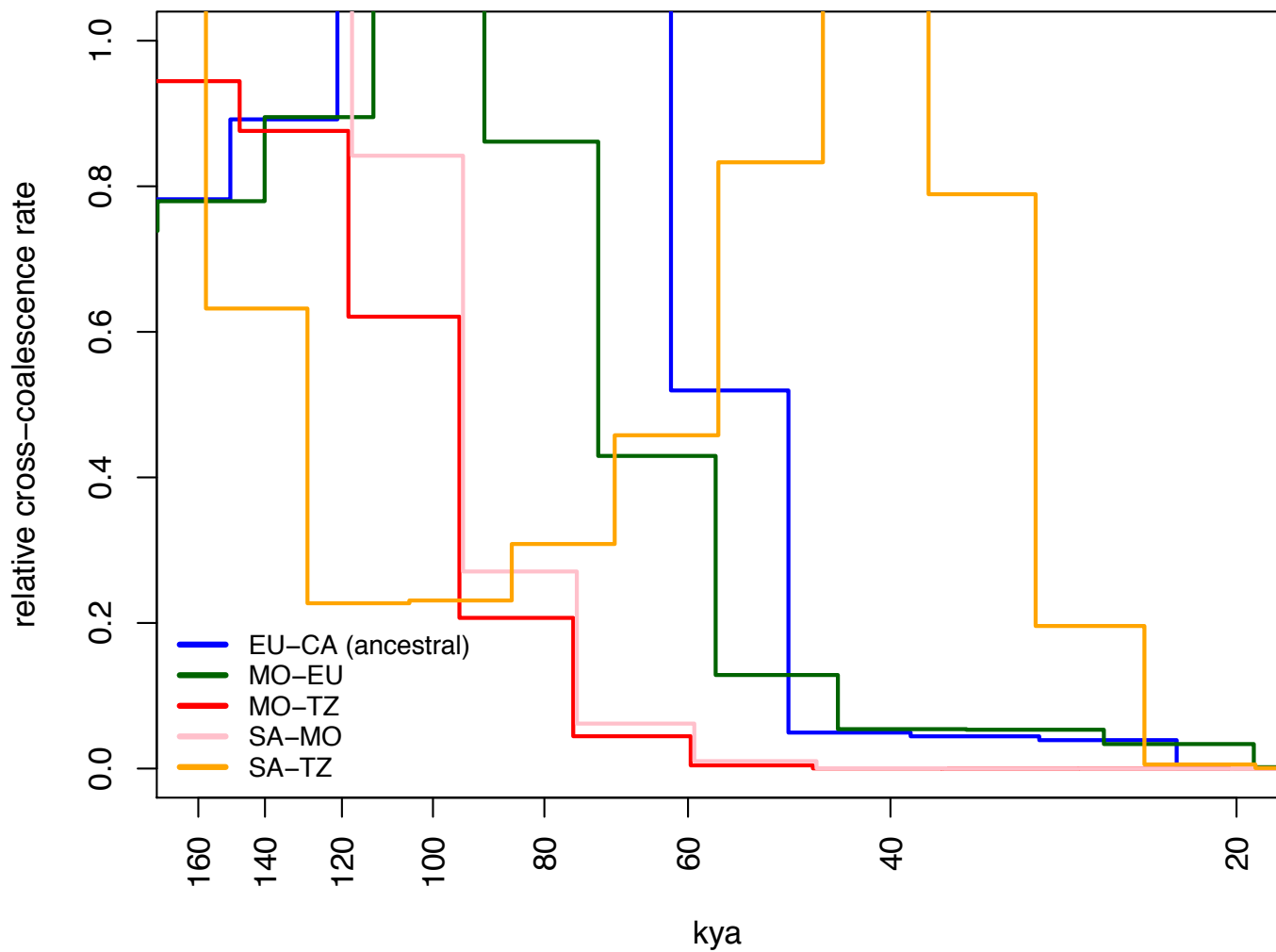
Supplementary Fig. 8

# Supplementary Fig. 9

Supplementary Fig. 10

a



b

# Supplementary Fig. 12



a

Europe, America

Canary island & Iberian relict

Moroccan, Cape Verde island, Italian & Levantine relict

*A. lyrata*

Yunnan, Tanzanian & South African relict

Central Asia

Yangtze

b

Europe, America

Canary island & Iberian relict

Moroccan, Cape Verde island, Italian & Levantine relict

*A. lyrata*

Yunnan, Tanzanian & South African relict

Central Asia

Yangtze

c

Europe, America

Canary island & Iberian relict

Moroccan, Cape Verde island, Italian & Levantine relict

*A. lyrata*

Yunnan, Tanzanian & South African relict

Yangtze

Central Asia