



HAL
open science

Ancient tree-topologies and gene-flow processes among human lineages in Africa

Gwenna Breton, Per Sjödin, Panagiotis I Zervakis, Romain Laurent, Alain Froment, Agnès E Sjöstrand, Barry S Hewlett, Luis B Barreiro, George H Perry, Himla Soodyall, et al.

► **To cite this version:**

Gwenna Breton, Per Sjödin, Panagiotis I Zervakis, Romain Laurent, Alain Froment, et al.. Ancient tree-topologies and gene-flow processes among human lineages in Africa. 2024. hal-04798019

HAL Id: hal-04798019

<https://hal.science/hal-04798019v1>

Preprint submitted on 22 Nov 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 **Title:**

2 **Ancient tree-topologies and gene-flow processes among human lineages**
3 **in Africa**

4
5 **Short Title:**

6 **Ancient demographic histories of African genetic lineages**

7
8 Gwenna Breton 1,2*

9 Per Sjödin 1

10 Panagiotis I. Zervakis 1

11 Romain Laurent 3

12 Alain Froment 4

13 Agnès E. Sjöstrand 1

14 Barry S. Hewlett 5

15 Luis B. Barreiro 6

16 George H. Perry 7,8

17 Himla Soodyall 9,10

18 Evelyne Heyer 3

19 Carina M. Schlebusch 1,11,12,*

20 Mattias Jakobsson 1,11,12,*

21 Paul Verdu 2,*

22

23 1 Human Evolution, Department of Organismal Biology, Evolutionary Biology Centre, Uppsala
24 University, Norbyvägen 18C, SE-752 36 Uppsala, Sweden

25 2 Centre for Medical Genomics, Department of Clinical Genetics and Genomics, Sahlgrenska
26 University Hospital, Gothenburg, Sweden

27 3 UMR7206 Eco-anthropology, CNRS-MNHN-Université Paris Cité, Paris, France

28 4 UMR 208 Patrimoines locaux, IRD-MNHN, Paris, France

29 5 Department of Anthropology, Washington State University, Vancouver, Washington, USA

30 6 Department of Medicine, Section of Genetic Medicine, University of Chicago, Chicago, IL, USA

31 7 Department of Anthropology, Pennsylvania State University, University Park, PA 16801, USA

32 8 Department of Biology, Pennsylvania State University, University Park, PA 16801, USA

33 9 Division of Human Genetics, School of Pathology, Faculty of Health Sciences, University of the
34 Witwatersrand and National Health Laboratory Service, Johannesburg, South Africa

35 10 Academy of Science of South Africa

36 11 Palaeo-Research Institute, University of Johannesburg, P.O. Box 524, Auckland Park, 2006, South
37 Africa

38 12 SciLifeLab, Uppsala, Sweden

39

40

41 * Authors for correspondence:

42 gwenna.breton@gu.se,

paul.verdu@mnhn.fr,

mattias.jakobsson@ebc.uu.se,

43 carina.schlebusch@ebc.uu.se

44

45

46

47 **Keywords:**

48 Africa; demographic history; admixture; migration; machine-learning; Approximate Bayesian
49 Computations; Whole Genome Sequences; Rainforest Hunter-Gatherers; Agriculturalists; Khoe and
50 San populations.

51

52

53

54 **Authors' contributions**

55 GB – Conceptualization, Data Curation, Formal Analysis, Investigation, Methodology, Software,
56 Supervision, Validation, Visualization, Writing – Original Draft Preparation.

57 PS – Conceptualization, Formal Analysis, Supervision, Validation, Writing – Review and Editing.

58 PLZ – Data Curation, Software, Writing – Review and Editing.

59 RL – Methodology, Software, Validation, Writing – Review and Editing.

60 AF – Investigation, Resources, Writing – Review and Editing.

61 AES – Investigation, Resources, Writing – Review and Editing.

62 BSH – Investigation, Resources, Writing – Review and Editing.

63 LBB – Investigation, Resources, Writing – Review and Editing.

64 GHP – Investigation, Resources, Writing – Review and Editing.

65 HS – Investigation, Resources, Writing – Review and Editing.

66 EH – Investigation, Resources, Writing – Review and Editing.

67 CMS – Conceptualization, Funding Acquisition, Investigation, Methodology, Project Administration,
68 Resources, Supervision, Validation, Writing – Review and Editing.

69 MJ – Conceptualization, Funding Acquisition, Investigation, Methodology, Project Administration,
70 Resources, Supervision, Validation, Writing – Review and Editing.

71 PV – Conceptualization, Formal Analysis, Funding Acquisition, Investigation, Methodology, Project
72 Administration, Resources, Software, Supervision, Validation, Visualization, Writing – Original Draft
73 Preparation.

74

75

76 **Abstract**

77

78 The deep history of humans in Africa and the complex divergences and migrations among
79 ancient human genetic lineages remain poorly understood and are the subject of ongoing
80 debate. We produced 73 high-quality whole genome sequences from 14 Central and Southern
81 African populations with diverse, well-documented, languages, subsistence strategies, and
82 socio-cultural practices, and jointly analyze this novel data with 104 African and non-African
83 previously-released whole genomes. We find vast genome-wide diversity and individual
84 pairwise differentiation within and among African populations at continental, regional, and
85 even local geographical scales, often uncorrelated with linguistic affiliations and cultural
86 practices. We combine populations in 54 different ways and, for each population combination
87 separately, we conduct extensive machine-learning Approximate Bayesian Computation
88 inferences relying on genome-wide simulations of 48 competing evolutionary scenarios. We
89 thus reconstruct jointly the tree-topologies and migration processes among ancient and recent
90 lineages best explaining the diversity of extant genomic patterns. Our results show the necessity
91 to explicitly consider the genomic diversity of African populations at a local scale, without
92 merging population samples indiscriminately into larger *a priori* categories based on
93 geography, subsistence-strategy, and/or linguistics criteria, in order to reconstruct the diverse
94 evolutionary histories of our species. We find that, for all different combinations of Central
95 and Southern African populations, a tree-like evolution with long periods of drift between short
96 periods of unidirectional gene-flow among pairs of ancient or recent lineages best explain
97 observed genomic patterns compared to recurring gene-flow processes among lineages.
98 Moreover, we find that, for 25 combinations of populations, the lineage ancestral to extant
99 Southern African Khoe-San populations diverged around 300,000 years ago from a lineage
100 ancestral to Rainforest Hunter-Gatherers and neighboring agriculturalist populations. We also
101 find that short periods of ancient or recent asymmetrical gene-flow among lineages often
102 coincided with epochs of major cultural and ecological changes previously identified by paleo-
103 climatologists and archaeologists in Sub-Saharan Africa.

104

105

Introduction

107 Unraveling genetic structure and gene flow among human lineages in Africa is crucial to the
108 understanding of the biological evolution and diversity of *Homo sapiens* throughout the continent
109 (Henn, Steele and Weaver, 2018; Schlebusch and Jakobsson, 2018; Pfennig *et al.*, 2023).

110 Population geneticists largely agree today that *Homo sapiens* spent a large part of its genetic
111 evolution within Africa only, between its gradual emergence from anatomically archaic forms 600,000-
112 200,000 years ago and the beginning of the Out-of-Africa expansions to the rest of the world around
113 100,000-50,000 years ago (Schlebusch and Jakobsson, 2018). For the past 20 years, the detailed
114 demographic and evolutionary history that shaped the genetic diversity of extant populations since the
115 Out-of-Africa has been, and is still, the subject of numerous investigations in various regions of Africa
116 (Verdu *et al.*, 2009, 2013; Schlebusch *et al.*, 2012; Breton *et al.*, 2014; Patin *et al.*, 2014, 2017; Perry
117 *et al.*, 2014; Busby *et al.*, 2016; Pierron *et al.*, 2017; Semo *et al.*, 2020; Lucas-Sánchez, Serradell and
118 Comas, 2021; Sengupta *et al.*, 2021; Fortes-Lima *et al.*, 2022, 2024; Laurent *et al.*, 2023; Pfennig *et al.*,
119 2023).

120 However, how ancient genetic divergences, and the dynamics of admixture and/or migration events
121 among ancient human lineages, influenced the genetic landscape observed today throughout Africa,
122 largely remains to be assessed (Bergström *et al.*, 2021). Classically tested demographic models
123 (Stringer, 2002; Henn, Steele and Weaver, 2018; Hollfelder *et al.*, 2021; Ragsdale *et al.*, 2023), range
124 from a unique ancestral *Homo sapiens* genetic population having recently and rapidly diverged into
125 extant African populations via series of founding events and/or multifurcations -often referred to as
126 tree-like models-, to models where extant populations diverged in a remote past and remained isolated
127 over long periods of time until relatively recently -often referred to as multiregional models-. Moreover,
128 each model may encompass possible gene exchanges among lineages via migration or admixture
129 processes. Importantly, the timing, duration, and/or intensity of such gene-flow events may
130 fundamentally change the most likely topology of bifurcating tree-like models compared to results
131 obtained without gene flows, and may also create reticulations among pairs of lineages throughout
132 history (Ragsdale *et al.*, 2023). Finally, these classical albeit highly complex models have recently been
133 enriched with possible introgression events from now extinct, unknown, or unsampled non-*Homo*
134 *sapiens* or ancient “ghost” *Homo-sapiens* lineages (Lachance *et al.*, 2012; Lorente-Galdos *et al.*, 2019;
135 Lipson *et al.*, 2022; Fan *et al.*, 2023; Pfennig *et al.*, 2023; Ragsdale *et al.*, 2023). Such latter models
136 became plausible in Africa in analogy to the likely events of introgressions from now extinct hominid
137 species into *Homo sapiens* lineages unveiled outside of Africa by the major advances of paleogenomics
138 (Meyer *et al.*, 2012; Prüfer *et al.*, 2014).

139 In this context, several recent investigations tested a variety of the above models with maximum-
140 likelihood approaches relying on whole genome sequences from several extant populations sampled
141 throughout the continent. They often reached highly contrasted results, with different ancient tree-
142 topologies between different numbers of ancient lineages (Schlebusch and Jakobsson, 2018; Lorente-
143 Galdos *et al.*, 2019; Lipson *et al.*, 2022; Fan *et al.*, 2023; Ragsdale *et al.*, 2023). Furthermore, they
144 found that reticulations among a limited number of ancient lineages explained observed genetic patterns
145 without the necessity of archaic introgressions (Ragsdale *et al.*, 2023); or, instead, that admixture with
146 other non-*Homo sapiens* or ancient *Homo sapiens* “ghost” populations most satisfactorily explained the
147 results (Lachance *et al.*, 2012; Lorente-Galdos *et al.*, 2019; Lipson *et al.*, 2022; Fan *et al.*, 2023).

148 A possible source of such vast differences among obtained results is likely the variety of population
149 sets considered in each study at a continental scale. Indeed, genetic diversity and differentiation among
150 populations is maximal in Africa (Tishkoff *et al.*, 2009; Skoglund and Mathieson, 2018; Pfennig *et al.*,
151 2023), and previously shown to be due to substantial differences in local populations’ demographic

152 histories (e.g. (Verdu *et al.*, 2009; Schlebusch *et al.*, 2012; Busby *et al.*, 2016; Patin *et al.*, 2017)).
153 Therefore, representing vast regions of the continent by a single population (Lorente-Galdos *et al.*,
154 2019; Ragsdale *et al.*, 2023), or even merging data from differentiated groups (Fan *et al.*, 2023), likely
155 influenced results obtained across studies and their interpretations.

156 In addition, and most importantly, all above models are highly nested, i.e. overlapping and
157 genetically largely indistinguishable for vast spaces of model-parameter values. Indeed, depending on
158 values of divergence times, effective population sizes, timing, duration and intensity of migration or
159 admixture events among lineages, tree-like, reticulated, and even multiregional models, with or without
160 ancient admixture, may be highly mimetic of one another. This essential methodological difficulty is
161 well illustrated by the often relatively similar values of posterior likelihoods of the various competing
162 models, thus indeed proved empirically hard to discriminate with the maximum-likelihood approaches
163 deployed, even more so when competing models differ in the number and specifications of parameters
164 explored with different statistics. These fundamental statistical issues remain despite the use of high-
165 quality whole genomes, an increasing number of populations considered at once, and increasingly
166 powerful methods based on innovative statistics (Gravel, 2012; Lorente-Galdos *et al.*, 2019; Ragsdale
167 and Gravel, 2019; Kamm *et al.*, 2020; Fan *et al.*, 2023; Ragsdale *et al.*, 2023). Finally, the lack of
168 ancient DNA data older than a few thousand years in Africa, strongly hampers the empirical testing of
169 the possible occurrence of ancient *Homo sapiens* “ghost” or non-*Homo sapiens* introgression events in
170 *Homo sapiens* lineages in Africa (Lachance *et al.*, 2012; Skoglund and Mathieson, 2018; Lorente-
171 Galdos *et al.*, 2019; Lipson *et al.*, 2022; Ragsdale *et al.*, 2023).

172 In this context, we investigated 74 high coverage (>30x) whole genome sequences from an
173 anthropologically well-characterized sample of hunter-gatherer, herder, and agricultural neighboring
174 populations from Central and Southern Africa (**FigureF1x, TableT1x**), merged together with 105
175 previously published high-quality whole genomes (SGDP-Mallick *et al.* 2016, HGDP-Meyer *et al.*
176 2012, 1KGP-Auton *et al.* 2015, Rasmussen *et al.* 2014, SAHGP-Choudhury *et al.* 2017). In particular,
177 we investigated genomic diversity patterns and aimed at reconstructing ancient demographic histories
178 among lineages having led to different groups of Khoe-San hunter-gatherer or herder populations from
179 Southern Africa (Schlebusch *et al.* 2012), and to different groups of both Eastern or Western Congo
180 Basin hunter-gatherer populations (often historically designated “Pygmies” by Europeans) and their
181 agriculturalist (“non-Pygmy”) neighbors with whom they share, nowadays, complex socio-economic
182 interactions (Verdu *et al.* 2013; Hewlett 2014). Investigating detailed ancient demography and gene-
183 flow among these three groups of populations, taking explicitly into account local differentiations and
184 possible gene-flow among and within groups, is of particular interest as they have previously been
185 identified to be among the most deeply diverged lineages in our species (Li *et al.*, 2008; Tishkoff *et al.*,
186 2009; Schlebusch and Jakobsson, 2018; Schlebusch *et al.*, 2020; Pfennig *et al.*, 2023).

187 We conducted formal statistical choice of competing scenarios and subsequent joint estimation of
188 parameters under the winning scenario using machine-learning Approximate Bayesian Computation
189 (ABC) (Tavaré *et al.*, 1997; Beaumont, Zhang and Balding, 2002; Blum and François, 2010; Csilléry,
190 François and Blum, 2012; Pudlo *et al.*, 2016), a methodological approach fundamentally differing from
191 all maximum-likelihood methods previously deployed, and only rarely explored in previous studies
192 reconstructing ancient demography in Africa (Lorente-Galdos *et al.*, 2019). ABC allows us to compare
193 a range of complex demographic scenarios presumed to have led to observed genetic patterns in
194 numerous population samples, and to estimate posterior parameter distributions best mimicking the data
195 under the winning scenario (e.g. (Verdu *et al.*, 2009; Lorente-Galdos *et al.*, 2019; Laurent *et al.*, 2023)).
196 Briefly, this is achieved based on informative summary-statistics computed on the observed data and
197 on numerous explicit genetic simulations for which scenario-parameters are drawn randomly in large
198 prior distributions set by the user. We thus aimed at inferring jointly competing tree-topologies,
199 divergence times, effective population size changes, and timing, duration and intensity of possible

200 asymmetric gene-flow events, to reconstruct the detailed evolutionary mechanisms underlying the
201 genomic diversity of extant Central and Southern African populations.

202

203 -----

204 **Figure 1x: Sampling location and number of whole-genome sequenced individuals.**

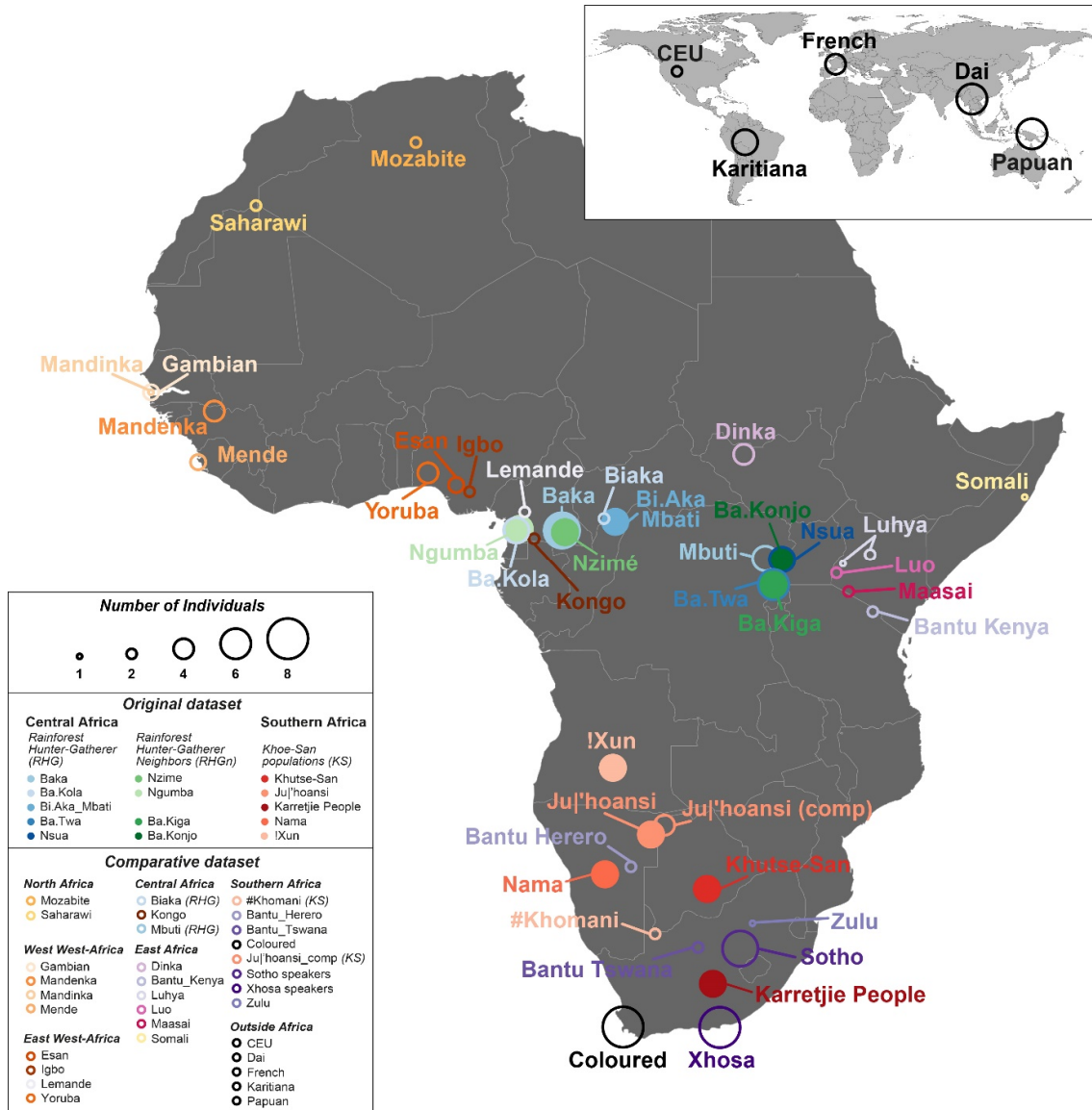
205 Individuals and populations originally sequenced in this work are indicated with filled circles. Individuals and populations
206 from previously published data merged with the original data set for comparison purposes are indicated with open circles.
207 Circles' diameters are proportional to the number of individuals as indicated in the top row of the legend in the bottom left
208 corner.

209 -----

210

211

212 **Figure F1x**
 213
 214
 215



216
217
218
219

TableT1x: Population table

Geographical location of population samples are indicated in **FigureF1x**.

Population Name ¹	N ²	Dataset ³	Sampling Location	Language Family (language)	Close socio-economic interactions with	Samples included in the ROH-ASD-ADMIXTURE analyses (FigureF2x, FigureF3x, FigureF4x)	Samples included in the ABC analyses (FigureF5x, FigureF6x, FigureF7x, FigureF8x, FigureF9x)
Original dataset							
Baka	7	this study	Bosquet (Cameroon)	Niger-Congo non Bantu Adamawa Ubangian Gbanzili (Baka)	Nzime	Yes	wRHG, 5 individuals with highest coverage
Nzime	5	this study	Messeca (Cameroon)	Bantu A842 (Nzime)	Baka	Yes	RHGn (West)
Ba.Kola	5	this study	Dispersed between Lolodorf and Kribi (Cameroon)	Bantu A80 (Kola)	Ngumba	Yes	wRHG
Ngumba	5	this study	Dispersed between Lolodorf and Kribi (Cameroon)	Bantu A80 (Ngumba)	Ba.Kola	Yes	RHGn (West)
Bi.Aka_Mbati	5	this study	Bombeketi section of Bagandou (Central African Republic)	Bantu C10 (Mbati)	Mbati farmers (not in the dataset)	Yes	wRHG
Ba.Kiga	5	this study	Mukono (Uganda)	Bantu J10 (Kiga)	Ba.Twa	Yes	RHGn (East)
Ba.Twa	6	this study	Kebiremu, Byumba, Kitariro, Mgungu, Nteko (Uganda)	Bantu J11 (Twa)	Ba.Kiga	Yes	eRHG, 5 individuals with highest coverage
Nsua	5	this study	Bundimassoli (Uganda)	Sudanic Mangbutu (Efe)	Ba.Konjo	Yes	eRHG
Ba.Konjo	5	this study	Mulimassenge (Uganda)	Bantu J40 (Konjo)	Nsua	Yes	RHGn (East)
!Xun	5	this study	Omega camp (Namibia) and Schmidtsdrift (Sout Africa) ⁴	Khoisan (Ju)	<i>Not applicable</i>	Yes	nKS
Ju 'hoansi	5	this study	Tsumkwe (Namibia)	Khoisan (Ju)	<i>Not applicable</i>	Yes	nKS
Nama	5	this study	Windhoek (Namibia)	Khoisan (Khoekhoe)	<i>Not applicable</i>	Yes	sKS
Karretjie People	5	this study	Colesberg (South Africa)	Khoisan (Tuu ancestral language) / Indo-European (Afrikaans current language)	<i>Not applicable</i>	Yes	sKS
Khutse-San	5	this study	Kutse Game reserve (Botswana)	Khoisan (Khoek Kalahari)	<i>Not applicable</i>	Yes	No
Comparative dataset							
Bantu_Herero	2	SGDP	Namibia	Bantu R30 (Herero)	<i>Not applicable</i>	Yes	No
Bantu_Kenya	2	SGDP	Kenya	<i>No data</i>	<i>Not applicable</i>	Yes	No
Bantu_Tswana	2	SGDP	South Africa	Bantu S31 (Tswana)	<i>Not applicable</i>	Yes	No
Biaka	2	SGDP	Central African Republic	Bantu C10 (Aka)	<i>Not applicable</i>	Yes	No
Coloured	8	SAHGP	South Africa	Indo-European (Afrikaans)	<i>Not applicable</i>	Yes	No
Dinka	4	SGDP, HGDP	Sudan	Nilo-Saharan (Dinka)	<i>Not applicable</i>	Yes	No
Esan	3	SGDP, KGP	Nigeria	Niger-Congo non Bantu (Esan)	<i>Not applicable</i>	Yes	No
Gambian	2	SGDP	Gambia	Niger-Congo non Bantu (Mandinka)	<i>Not applicable</i>	Yes	No
Igbo	2	SGDP	Nigeria	Niger-Congo non Bantu (Igbo)	<i>Not applicable</i>	Yes	No
Ju 'hoansi_com p.	4	SGDP, HGDP	Namibia	Khoisan (Ju)	<i>Not applicable</i>	Yes	No
#Khomani	2	SGDP	South Africa	Khoisan (Tuu)	<i>Not applicable</i>	Yes	No
Kongo	1	SGDP	Cameroon	<i>No data</i>	<i>Not applicable</i>	Yes	No
Lemande	2	SGDP	Cameroon	Bantu A46 (Lemande)	<i>Not applicable</i>	Yes	No
Luhya	3	SGDP, KGP	Kenya	Bantu JE32 (Luhya)	<i>Not applicable</i>	Yes	No
Luo	2	SGDP	Kenya	Nilo-Saharan (Dholuo)	<i>Not applicable</i>	Yes	No
Maasai	2	SGDP	Kenya	Nilo-Saharan (Maasai)	<i>Not applicable</i>	Yes	No

Mandenka	4	SGDP, HGDP	Senegal	Niger-Congo non Bantu (West Maninkakan)	<i>Not applicable</i>	Yes	No
Mandinka	1	KGP	Gambia	Niger-Congo non Bantu (Mandinka)	<i>Not applicable</i>	Yes	No
Mbuti	5	SGDP, HGDP	Democratic Republic of Congo	Bantu D30 (Bambuti)	<i>Not applicable</i>	Yes	eRHG
Mende	3	SGDP	Sierra Leone	Niger-Congo non Bantu (Mende)	<i>Not applicable</i>	Yes	No
Mozabite	2	SGDP	Algeria	Afro-Asiatic (Mozabite)	<i>Not applicable</i>	Yes	No
Saharawi	2	SGDP	Western Sahara	Afro-Asiatic (Saharawi)	<i>Not applicable</i>	Yes	No
Somali	1	SGDP	Kenya	Afro-Asiatic (Somali)	<i>Not applicable</i>	Yes	No
Sotho	7	SAHGP	South Africa	Bantu S30 (Sotho)	<i>Not applicable</i>	Yes	No
Xhosa	8	SAHGP	South Africa	Bantu S41 (Xhosa)	<i>Not applicable</i>	Yes	No
Yoruba	4	SGDP, HGDP	Nigeria	Niger-Congo non Bantu (Yoruba)	<i>Not applicable</i>	Yes	No
Zulu	1	SAHGP	South Africa	Bantu S42 (Zulu)	<i>Not applicable</i>	Yes	No
French	4	SGDP, HGDP	France	Indo-European (French)	<i>Not applicable</i>	Yes	No
CEU	2	KGP	United States of America	Indo-European (English)	<i>Not applicable</i>	Yes	No
Dai	6	SGDP, KGP	China	Tai (Dai)	<i>Not applicable</i>	Yes	No
Papuan	6	SGDP, HGDP	Papua New Guinea	Papuan (no data)	<i>Not applicable</i>	Yes	No
Karitiana	5	SGDP, HGDP	Brazil	Tupian (Karitiana)	<i>Not applicable</i>	Yes	No

220
221
222
223
224
225
226
227

¹Population Names are self-reported for the original dataset presented in this study

²Number of unrelated individuals considered in all analyses in this study (see **Material and Methods**)

³SGDP (Mallick *et al.*, 2016); KGP (The 1000 Genomes Project Consortium *et al.*, 2015); HGDP (Meyer *et al.*, 2012; Rasmussen *et al.*, 2014); SAHGP (Choudhury *et al.*, 2017).

⁴The place of origin of the !Xun is around Menongue in Angola.

Results

229 We generated high-coverage whole genome sequences (>30X) in 74 individuals from 14 Central and
230 Southern African populations for whom detailed ethno-anthropological information was gathered in the
231 field jointly with DNA samples (**FigureF1x** and **TableT1x**). We merged this original dataset with a
232 comparative dataset comprising 105 individuals from 27 African and five non-African populations for
233 whom similar high-quality raw whole genome sequencing data were made available to the community
234 (Meyer *et al.*, 2012; Rasmussen *et al.*, 2014; The 1000 Genomes Project Consortium *et al.*, 2015;
235 Mallick *et al.*, 2016; Choudhury *et al.*, 2017). After quality-control, variant-calling, and relatedness
236 filtering procedures conducted on all 179 individuals together, we retained 73 and 104 unrelated
237 individuals in our original and comparative datasets, respectively, for all subsequent analyses (see
238 **Material and Methods**).

239 Considering only the 73 Central and Southern African unrelated individuals newly sequenced here,
240 we identify a total of 26,780,319 biallelic SNPs (241,428 multiallelic SNPs), and 2,454,965 simple
241 insertions/deletions (969,025 complex indels), compared to the reference sequence of the human
242 genome GRCh38, out of which 854,114 (3.1893%) and 114,362 (4.6584%), respectively, were not
243 previously reported in dbSNP 156 (**SupplementaryTableST1x**). Among the 177 unrelated worldwide
244 individuals including our original dataset, we identify 36,272,545 biallelic SNPs (257,906 multiallelic
245 SNPs), and 3,159,306 simple insertions/deletions (977,542 complex indels), out of which 1,055,245
246 (2.9092%) and 189,124 (5.9863%), respectively, were not previously reported in dbSNP 156
247 (**SupplementaryTableST1x**).

248 Furthermore, we find a large variability in the mean number of biallelic SNPs across populations
249 (**FigureF2x-PanelA**, **SupplementaryTableST2x**). In particular, we find substantially more biallelic
250 SNPs within African populations and more variation of the mean number of SNPs across African
251 populations (from mean=3,726,018; SD=27,793 across 2 individuals in the Saharawi from Western
252 Sahara, to mean=4,628,957; SD=8,139 across 5 individuals in the Ju|'hoansi from Namibia), than within
253 and among non-African populations (from mean=3,214,086; SD=69,241 across 5 individuals in the
254 Karitiana from Brazil, to mean=3,551,792; SD=7,819 across 2 individuals from the USA CEU
255 population). Overall, Southern African Khoe-San populations exhibit the highest mean number of
256 biallelic SNPs as well as numbers of previously unreported SNPs, followed by Central African
257 Rainforest Hunter-Gatherer populations from the Congo Basin, and then by all other African
258 populations in our dataset. These results show that a substantial number of previously unknown variants
259 can still be found when investigating high-quality whole genome sequences from relatively under-
260 studied Sub-Saharan African populations, consistent with previous studies (Meyer *et al.*, 2012;
261 Rasmussen *et al.*, 2014; The 1000 Genomes Project Consortium *et al.*, 2015; Mallick *et al.*, 2016;
262 Choudhury *et al.*, 2017; Schlebusch *et al.*, 2020; Breton, Fortes-Lima and Schlebusch, 2021; Fan *et al.*,
263 2023; Ragsdale *et al.*, 2023).

264

265

266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282

FigureF2x: Whole autosomal genome bi-allelic SNPs counts, unbiased heterozygosities, and distributions of Runs of Homozygosity by bins of length, in 37 African and 5 non-African populations.

(A) Numbers of bi-allelic SNPs across individuals within populations of more than one individual, compared to the reference sequence of the human genome GCRh38 and to previously reported variants in dbSNP 156. Detailed variant-counts are provided in **SupplementaryTableST1x and SupplementaryTableST2x**. (B) Unbiased estimates of multi-locus heterozygosities (Nei, 1978), averaged across all variable (bi-allelic SNPs only) and non-variable autosomal sites with no-missing genotype within each population of more than one individual, corrected then for haploid population sample sizes. (C) Mean total ROH lengths in four bins of length categories for each population of more than one individual, separately. For each panel, we present a box zooming specifically on the results obtained for the original dataset of 73 unrelated individuals whole-genome sequenced from Central and Southern Africa.

See **Material and Methods** for filtering and calculation details and the software packages used. For (A) and (B), box-plots indicate the population median in between the first and third quartile of the box-limits, whiskers extending to data points no more than 1.5 times the interquartile range of the distribution, and empty circles for all more extreme points beyond this limit, if any. Population geographical location, categorization, and descriptions are provided in **FigureF1x and TableT1x**.

283 **FigureF2x**

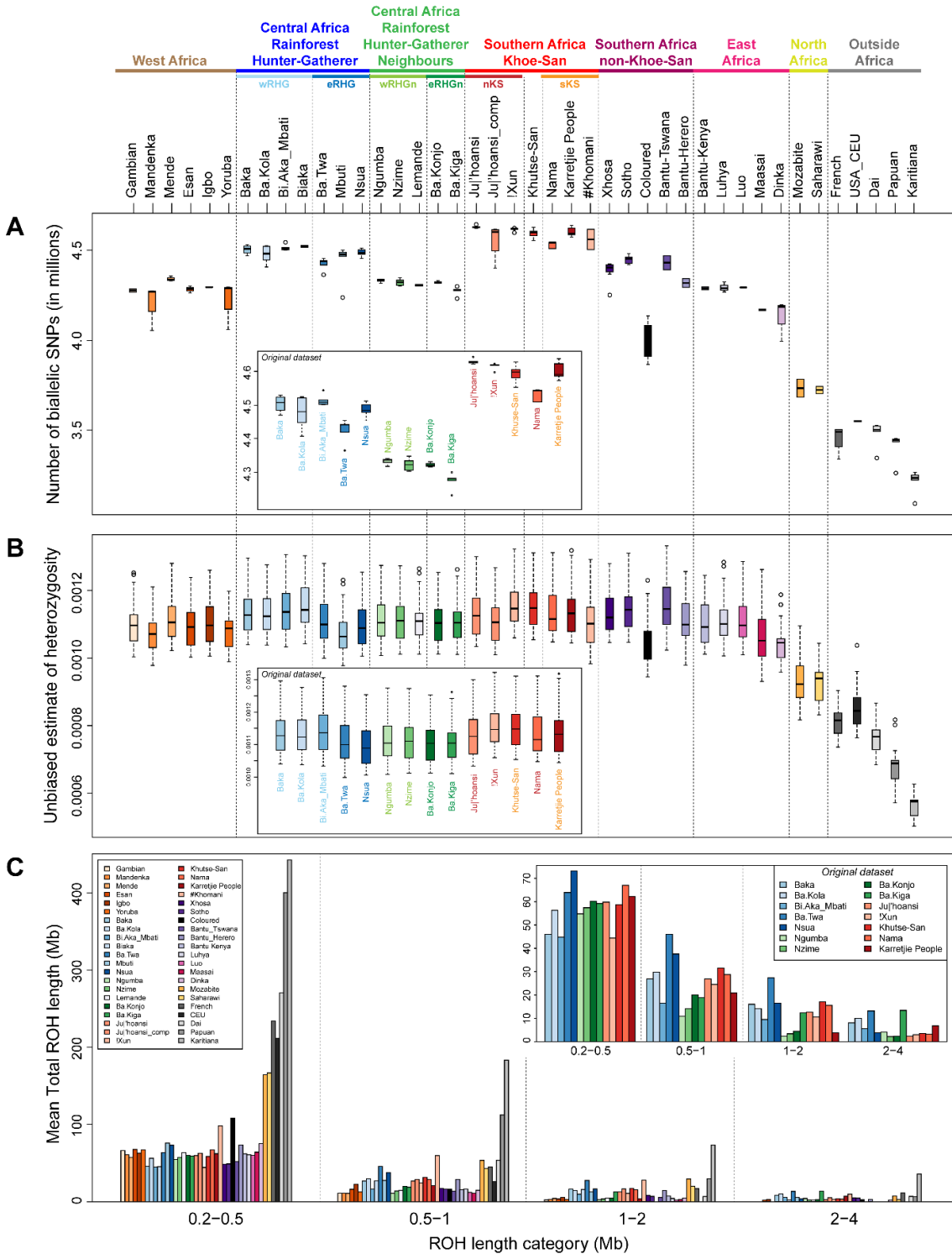
284

285

286

287

288



289 **1. Large genetic diversities and differentiations in Central and Southern Africa**

290 **1.A. Heterozygosities and Runs of Homozygosity**

291 We find very large differences across populations at a local scale in Central and Southern Africa in
292 unbiased heterozygosities (Nei, 1978), calculated on both varying and non-varying autosomal sites
293 (**FigureF2x-PanelA-B** and **SupplementaryTableST2x**). Indeed, we find varying genetic variation
294 across Northern Khoe-San populations (!Xun and Ju|'hoansi) and Southern KS populations (Nama,
295 Karretjie People and #Khomani). Furthermore, we find overall more variation in Rainforest Hunter-
296 Gatherer populations from the western part of the Congo Basin than in RHG populations from the east;
297 and western RHG populations exhibited much larger genetic diversities than all their respective RHG
298 neighbors.

299 In addition, we investigated the distributions of Runs of Homozygosity (ROH) across individuals
300 within each of the 46 populations in our dataset. In general (**FigureF2x-PanelC**), individuals from
301 Eastern RHG populations (Nsua, Mbuti, Ba.Twa), exhibit the longest total proportion of their autosomal
302 genome in ROH across Sub-Saharan African populations, for all ROH length-classes, with the
303 exception of #Khomani individuals from Southern Africa, and that of Coloured individuals for short
304 ROH only. Interestingly, we find that Western RHG populations (Ba.Kola, Baka, Biaka, Bi.Aka_Mbati)
305 have much less ROH of all classes than Eastern RHG. Furthermore, we find less short ROH, but more
306 long ROH, in Western RHG than in all RHG neighboring populations across the Congo Basin except
307 for the Ba.Kiga RHGn from Uganda. Finally, we find that all Northern and Southern Khoe-San
308 populations have relatively similar total lengths of short ROH compared to that of RHG neighbors, with
309 the notable exception of the !Xun, who have the smallest proportion of their genomes in short ROH,
310 worldwide. Nevertheless, we find substantially more ROH of longer class in KS populations than in
311 RHG neighbors, a relative pattern similar to what was observed for Western RHG. The distribution of
312 ROH sizes has been shown to be highly informative about important demographic processes such as
313 levels of inbreeding or endogamy within populations, and reproductive isolation or recent admixture
314 among populations (Pemberton *et al.*, 2012; Mooney *et al.*, 2018; Szpiech *et al.*, 2019; Laurent *et al.*,
315 2023). Thus, our ROH results highlight possibly vast differences in recent demographic processes and
316 levels of endogamy across Central and Southern African populations.

317 Altogether, our heterozygosity and ROH results show the vast genetic diversity of Sub-Saharan
318 African populations, as well as substantially diverging genomic patterns at a local geographical scale
319 among and within groups of Rainforest Hunter-Gatherers, RHG neighbors, and Khoe-San populations.
320

321 **1.B. Individual pairwise genetic differentiation**

322 We investigated individual pairwise genomic differentiation across the 177 worldwide unrelated
323 individuals in our dataset, including the 73 unrelated novel Central and Southern African individuals,
324 using Allele-Sharing Dissimilarities (ASD, (Bowcock *et al.*, 1991)). The Neighbor Joining Tree (NJT,
325 (Saitou and Nei, 1987; Gascuel, 1997)) representation of the pairwise ASD matrix shows three main
326 clusters among all African individuals, separated by a longer genomic distance from all non-African
327 individuals (**FigureF3x-PanelA**). One cluster corresponds to all Southern African Khoe-San
328 populations (represented by the pink and red symbols, see **FigureF3x-PanelB**), largely separated from
329 a cluster corresponding to Rainforest Hunter-Gatherer populations from the Congo Basin (represented
330 by the blue symbols); itself also largely separated from a third cluster grouping all other Western,
331 Eastern, Southern and Central African populations in our data set, the latter group including Rainforest
332 Hunter-Gatherer neighboring populations in green symbols (see, **FigureF3x-PanelB**).
333
334

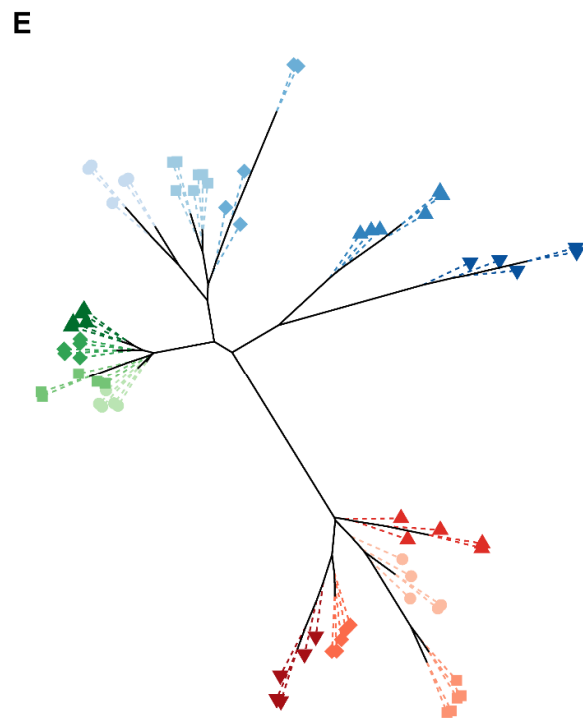
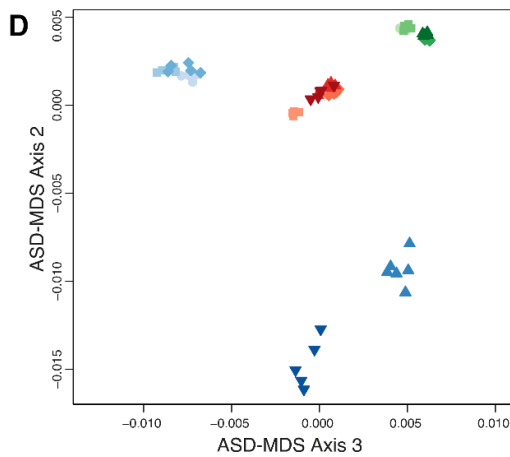
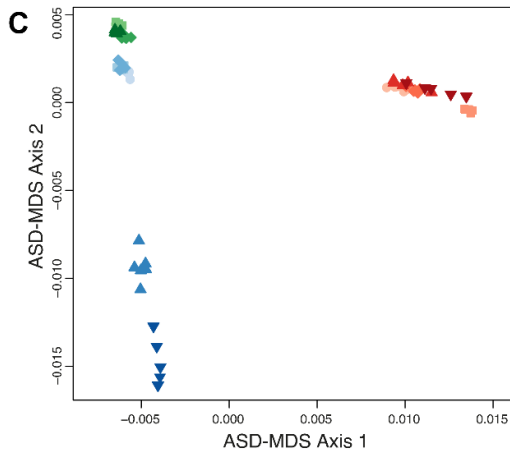
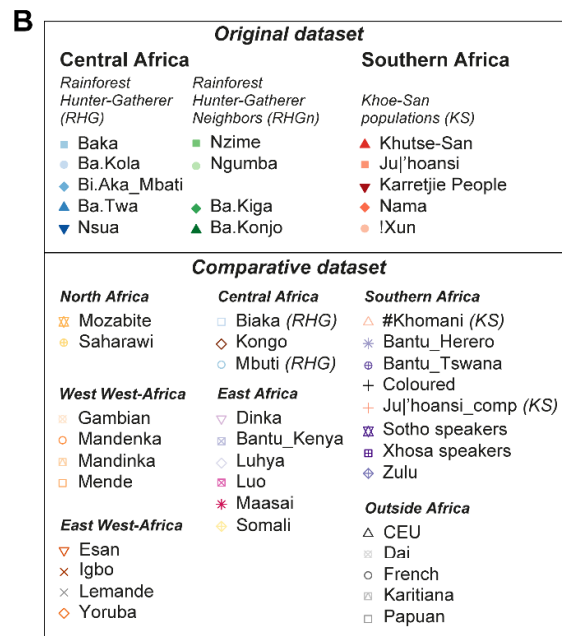
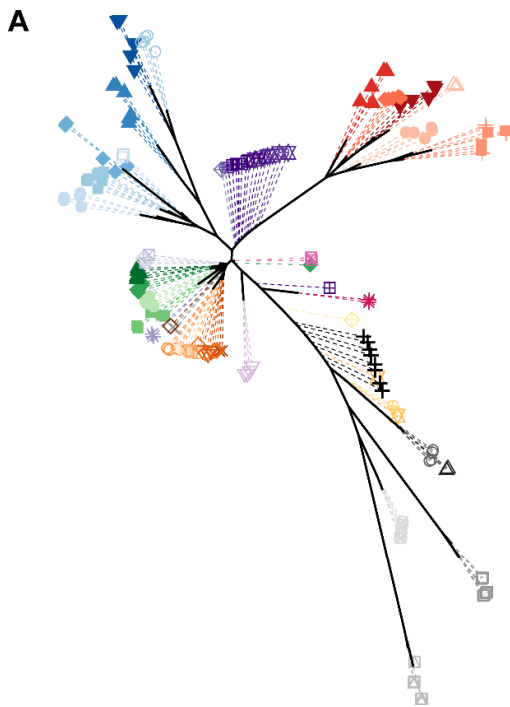
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350

FigureF3x: Individual pairwise genetic differentiation patterns.

Neighbor-Joining Tree NJT (Saitou and Nei, 1987; Gascuel, 1997) and Multi-Dimensional scaling representation of genome-wide Allele Sharing Dissimilarities (Bowcock *et al.*, 1991) among pairs of individuals at worldwide and regional African scales. We considered the 14,182,615 genome-wide autosomal SNPs pruned for low LD (r^2 threshold 0.1) to calculate ASD between all pairs of individuals. (A) NJT computed for all 177 worldwide individuals in the ASD matrix. For easing the visualization of the internal branches of the NJT, all terminal edges are represented in dotted lines each measuring 1/10th of their true size. (B) Individual symbols and colors identifying the 73 Central and Southern African individuals originally whole-genome sequenced with filled symbols, and the 104 individuals from worldwide populations, including Africa, merged with the original data with open symbols. (C) and (D) First three axes of ASD-MDS computed on a subset of the full ASD matrix comprising only the individuals sequenced anew, whose symbols are provided in (B). (E) NJT computed on the same individual subset of the ASD matrix used for (C) and (D). Terminal edges of this latter tree are represented with dotted lines each measuring 1/30th of their true size.

351
352
353
354

Figure F3x



355 Note importantly that individual pairwise dissimilarities at the genome-wide scale do not reflect
356 geography (**FigureF1x, TableT1x**), highlighting the large genetic differentiation found across African
357 populations and the complexity of its distribution at both a local and continental scale. From West to
358 East of the Congo Basin throughout Central Africa, RHG populations cluster separately from RHG
359 geographic neighbors with whom they nevertheless share close socio-economic interactions.
360 Furthermore, Central African RHGn are more genetically resembling certain other Western, Eastern
361 and Southern Africans than their immediate RHG geographical neighbors. Analogously, Southern
362 African Khoe-San populations are also clustering together and much more distant from other Southern
363 African geographically neighboring populations.

364 The Multidimensional Scaling (MDS) projection of the subsampled ASD matrix for only the 73
365 novel individuals from Central and Southern Africa, further shows substantial differentiation among
366 groups of individuals within each one of the three clusters identified at the continental scale. Indeed,
367 **FigureF3x-PanelC, D and E** show substantial differentiation between Western and Eastern Central
368 African Rainforest Hunter-Gatherer populations, and even substantial differentiation between Ba.Twa
369 and Nsua RHG very locally in Uganda. Conversely, we find relatively much shorter genetic
370 differentiation among pairs of RHG neighbors from West to East of the entire Congo Basin. Finally,
371 results also show substantial differentiation across Southern African Khoe-San populations, albeit these
372 populations are relatively closer from one another than the different RHG populations.

373 NJT and MDS based on the ASD pairwise matrix provide an important view of the major axes of
374 genetic differentiation and variation across samples, but do not easily allow to visually describe higher-
375 order axes of genomic variations. To do so, we conducted an ADMIXTURE analysis (Alexander et al.
376 2009), on the same entire worldwide dataset for increasing values of K (**FigureF4x**). This descriptive
377 method is known to capture the same information as ASD-MDS and ASD-NJT, but allows to explore
378 multiple axes of variation at the same time (Pritchard, Stephens and Donnelly, 2000; Rosenberg, 2002;
379 Falush, Stephens and Pritchard, 2003; Alexander, Novembre and Lange, 2009; Lawson, van Dorp and
380 Falush, 2018; Peter, 2022; Laurent *et al.*, 2023).

381 -----

382 **FigureF4x:**
383 **Worldwide interindividual genetic diversity patterns with ADMIXTURE**
384 Each individual is represented by a single vertical line divided in K colors each proportional to an individual's genotype
385 membership proportion assigned by ADMIXTURE (Alexander, Novembre and Lange, 2009) into the virtual cluster of that
386 color. Two individuals showing the same relative proportions of each color at a given value of K are genetically more
387 resembling one another than two individuals with different relative proportions of the same colors. Population and regional
388 geographic groupings are not considered for the calculations, individuals are *a posteriori* grouped by populations and ordered
389 by the user. Each population is separated by a thin vertical black line. Unsupervised clustering conducted with 840,031
390 genome-wide SNPs pruned for LD (r^2 threshold 0.1) and minor allele frequency (0.1) using ADMIXTURE for values of K
391 ranging from 2 to 10, considering 177 worldwide individuals. For each value of K separately, we performed 20 independent
392 runs and used PONG (Behr *et al.*, 2016), to identify groups of resembling results (called "modes") based on Symmetric
393 Similarity Coefficient measures of individual results. The number of runs among the 20 independent runs belonging to the
394 same mode result (i.e. that have pairwise SSC above 0.998), is indicated below the number of K on the left of each barplot
395 separately. Only those runs are averaged per individual to provide the barplot result for each value of K separately.
396 -----

397 -----
398
399

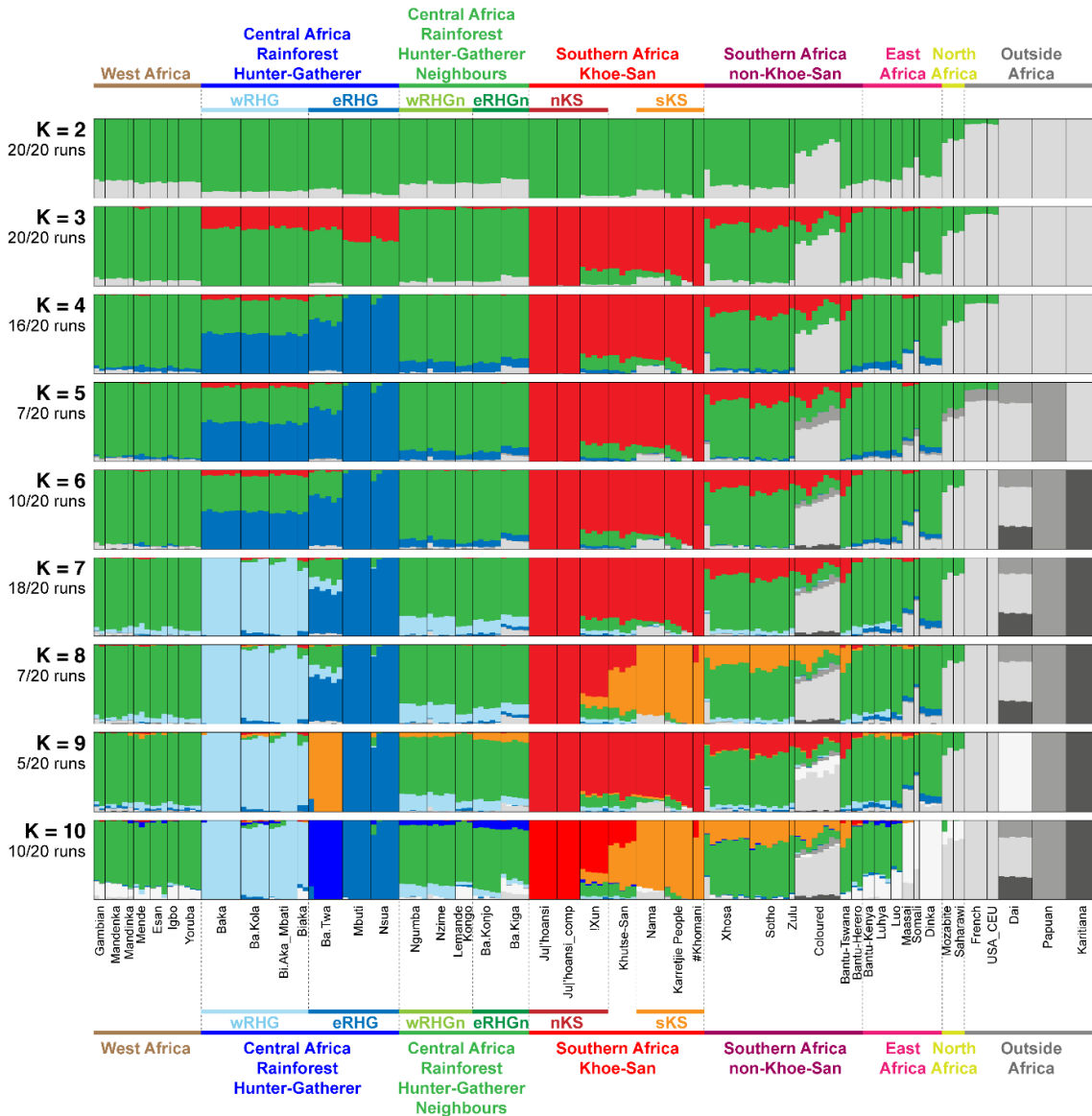
400 **FigureF4x**

401

402

403

404



405 At K=2 in **FigureF4x**, the green virtual genetic cluster is maximized in Northern KS populations
406 while the gray cluster is maximized in non-African populations, consistent with largest genetic distances
407 in the ASD-NJT in **FigureF3x**, with all other individuals in our dataset presenting intermediate, non-
408 resolved, genotype membership proportions to either cluster.

409 At K=3, a novel red cluster is maximized in all Southern KS populations separating them from the
410 rest of African individuals, also consistent with the ASD-NJT. At K=4, the novel blue cluster
411 differentiates Central African RHG populations, and in particular certain Eastern RHG populations who
412 maximize their membership proportions to this cluster.

413 The new dark and darker gray clusters respectively at K=5 and K=6 differentiate Oceanian Papuan
414 individuals and South American Karitiana individuals from one another and from all other non-African
415 individuals, without affecting substantially genotype membership proportions from African individuals.

416 At K=7, the novel light blue cluster differentiates Eastern and Western RHG populations, as
417 observed in **FigureF3x**.

418 At K=8, the novel orange cluster differentiates mainly among Northern and Southern Khoe-San
419 populations, with the Central Khutse San showing substantial membership proportions to both the red
420 and the orange clusters, respectively maximized in Northern and Southern KS.

421 At K=9, an alternative clustering solution is shown where East Asian Dai, Oceania Papuan, and
422 South American Karitiana individuals are all clustering separately in three fully resolved clusters,
423 affecting the African clustering patterns where the differentiation previously observed at K=8 among
424 Northern and Southern Khoe-San individuals disappears to the benefit of a novel clustering solution
425 where Ba.Twa cluster separately from all RHG and from Eastern RHG with whom they shared closer
426 genotype membership proportions to the medium blue clusters at previous values of K in particular.

427 Finally, at K=10, these alternative clustering solutions are resolved into seven different virtual
428 genetic clusters maximized in different groups of African individuals. The light blue cluster is
429 maximized in certain Western RHG populations. The dark blue cluster is maximized in the Ba.Twa
430 Eastern RHG individuals only, while the medium blue cluster is maximized in the two other Eastern
431 RHG populations mostly. The red cluster is mainly represented in the Northern KS individuals, while
432 the orange cluster is maximized in Southern KS. Finally, the green cluster is maximized in all other
433 African populations, albeit note that it is not fully resolved at this value of K since no individual exhibits
434 100% genotype membership to this cluster.

435 Most importantly, note that numerous Central and Southern African individuals retain, from K=2
436 to 10, intermediate genotype membership proportions in between certain clusters. This should be
437 interpreted primarily as these individuals being at intermediate genetic distance between individuals
438 presenting 100% membership to either cluster, which is consistent with ASD-MDS and ASD-NJT
439 projections (**FigureF3x**). In turn, such intermediate distances may be either due to yet-unresolved
440 ADMIXTURE clustering, which may appear at higher values of K, or can be due to admixture having
441 occurred between ancestors respective to each cluster, the latter interpretation being likely but not
442 formally tested by the ADMIXTURE method. The possible occurrence of gene-flow events across pairs
443 of lineages and their influence on tree-topologies are incorporated explicitly in our Approximate
444 Bayesian Computation inferences detailed below (Tavaré *et al.*, 1997; Pritchard *et al.*, 1999; Beaumont,
445 Zhang and Balding, 2002; Blum and François, 2010; Csilléry, François and Blum, 2012; Pudlo *et al.*,
446 2016; Raynal *et al.*, 2019).

447
448 Altogether, our descriptive results highlight the vast genetic diversity and differentiation of African
449 populations, maximized in certain groups of Central and Southern African populations, as previously
450 reported (e.g. (Meyer *et al.*, 2012; Rasmussen *et al.*, 2014; Mallick *et al.*, 2016; Choudhury *et al.*, 2017;
451 Fan *et al.*, 2023; Ragsdale *et al.*, 2023; Fortes-Lima *et al.*, 2024)). Finally, and most importantly, our
452 results highlight in particular the vast genetic diversity and differentiation of populations at a very local

453 scale in Africa, including among immediate neighbors, not trivially correlated with geographical
454 distances among populations, linguistic classification, nor mode of subsistence. In fact, these results
455 anticipate that historical and demographic inferences in Africa may substantially differ when
456 considering different population samples across the continent and even at a local scale, as well as
457 advocate for extreme caution when grouping individuals and populations samples into larger categories.
458

459 **2. Demographic and migration histories in Central and Southern African**

460 **2.A. Summary of the ABC inference design**

461 We aimed at reconstructing the demographic and migration history that produced extant genomic
462 patterns observed within and among Central and Southern African populations (**FigureF2x, FigureF3x,**
463 **FigureF4x**), with machine-learning ABC scenario-choices followed by posterior-parameter
464 estimations. To do so, we considered eight competing possible tree-topologies among five extant
465 Northern and Southern Khoe-San, Western and Eastern Rainforest Hunter-Gatherer, and Rainforest
466 Hunter-Gatherer neighboring populations (**FigureF5x**). Furthermore, for each topology, we considered
467 two possible gene-flow processes among recent and ancient genetic lineages: instantaneous asymmetric
468 gene-flow corresponding to instantaneous unidirectional introgression events between pairs of lineages;
469 and recurring asymmetric gene-flow corresponding to unidirectional recurring migrations among pairs
470 of lineages (scenarios “i” and “r” respectively, **FigureF5x**). Finally, we considered, for each topology
471 and each gene-flow process, three nested gene-flow intensities for each event among pairs of lineages
472 separately: no to very high gene-flow rates (scenarios “i1” and “r1”); no to moderate rates (scenarios
473 “i2” and “r2”); and no to limited rates (scenarios “i3” and “r3”).

474 This led to $8 \times 2 \times 3 = 48$ competing scenarios in total which can be grouped in different ways to
475 address formally two nested major questions:

- 476
- 477 i) which ancestral lineages to extant populations diverged first from the others and when did they do so,
478 when considering complex gene-flow events among recent and ancient lineages?
479
 - 480 ii) did gene-flow events occur recurrently or more instantaneously among pairs of lineages during the
481 evolutionary history of Central and Southern African populations? When did these events occur? How
482 intense were they?
483

484 Importantly, we aimed at considering the vast genetic diversity and differentiation within and
485 among populations and groups of populations at a local and regional scale (**FigureF2x, FigureF3x,**
486 **FigureF4x**, and (Tishkoff *et al.*, 2009; Skoglund and Mathieson, 2018; Pfennig *et al.*, 2023)).
487 Therefore, we replicated the same ABC scenario-choice and posterior-parameter estimation procedures
488 considering, in turn, 54 different combinations of five sampled populations of five individuals’ whole
489 autosomal genomes each (**FigureF5x**).

490 We thus conducted a total of 240,000 coalescent simulations under these 48 competing scenarios
491 (5000 simulations per scenario), by drawing parameter values from prior distributions set by the user
492 (**FigureF5x, TableT2x**). For each simulation, we then calculated a vector of 337 summary-statistics
493 (**TableT3x**), within and among the five simulated populations; each vector of summary-statistics thus
494 corresponds to a vector of parameter values drawn randomly from prior distributions and used for one
495 simulation. We then deployed Random Forest ABC procedures (Pudlo *et al.*, 2016; Estoup *et al.*, 2018),
496 to identify the winning scenario or group of scenarios for each 54 combinations of five sampled “real”
497 populations separately. Finally, under the winning scenario, we produced 100,000 simulations,
498 computed 202 summary-statistics for each simulation (**TableT3x**), and performed Neural Network

499 ABC posterior parameter joint-estimation (Blum and François, 2010; Csilléry, François and Blum,
500 2012), providing posterior-parameter distributions most likely to have produced observed genomic data,
501 for each scenario-parameter.

502

503

504 **FigureF5x: 48 competing scenarios for the history of Central and Southern African populations.**

505 (A) Eight competing topologies for the demographic history of five Central and Southern African extant lineages. The Northern
506 Khoe-San lineage (nKS) is represented by either the Ju|'hoansi or the !Xun population. The Southern Khoe-San lineage (sKS)
507 is represented by either the Karretjie People or the Nama population. The Rainforest Hunter-Gatherer Neighbors lineage
508 (RHGn) is represented either by the Western Congo Basin Nzime or Ngumba populations, or the Eastern Ba.Konjo or Ba.Kiga
509 populations. The Western Rainforest Hunter-Gatherer lineage (wRHG) is represented by either the Baka, the Ba.Kola, or the
510 Bi.Aka_Mbati population. The Eastern Rainforest Hunter-Gatherer lineage (eRHG) is represented by either the Nsua, the
511 Ba.Twa, or the Mbuti population. As indicated in legend in the bottom right corner of the panel, the eight topologies can be
512 grouped according to i) which ancestral lineage diverged first from the two others, or ii) whether the nKS and sKS lineages
513 split earlier or later than the split of the wRHG and eRHG lineages. For each topology, we consider possible changes in
514 lineages' constant effective population size, N_e , after each divergence event, t . For each topology, we considered possible
515 gene-flow events among ancestral lineages and among recent lineages, represented as horizontal uni-directional arrows. (B)
516 Example for a given topology (1a), of the fact that each one of the eight topologies are considered under two alternative
517 competing "instantaneous" (Scenario i-1a) or "recurring" (Scenario r-1a) gene-flow processes. Instantaneous gene-flow
518 processes consider that genes can be exchanged among pairs of ancestral or recent lineages at a single parameterized time, tad ,
519 with independent parameters of gene-flow intensity, m , from lineage "A" to "B" and from lineage "B" to "A". Recurring gene-
520 flow processes consider instead that gene-flow occur at each generation between pairs of lineages with a constant independent
521 rate, m , from lineage "A" to "B" and from lineage "B" to "A". Note that for all eight topologies and for both instantaneous
522 and recurring gene-flow processes, we parameterized separately the gene-flow from lineage "A" into lineage "B", and the
523 gene-flow from lineage "B" into "A", in order to allow for possibly asymmetrical gene-flow events among pairs of lineages.
524 For both instantaneous and recurring gene-flow processes, we considered three nested intensities of gene-flow parameters
525 indicated at the bottom of (B). Scenarios i1 or r1 consider that the gene-flow parameters m in each topology are independently
526 drawn in Uniform distributions between 0 and 1 or 0 and 0.05, respectively, thus allowing for the possibility of no to high
527 gene-flows from one lineage to the other. Scenarios i2 or r2 consider that the gene-flow parameters m in each topology are
528 independently drawn in Uniform distributions between 0 and 0.125 or 0 and 0.0125, respectively, thus allowing only for the
529 possibility of no to intermediate gene-flow from one lineage to the other. Finally, scenarios i3 or r3 consider that the gene-
530 flow parameters m in each topology are independently drawn in Uniform distributions between 0 and 0.01 or 0 and 0.001,
531 respectively, thus only allowing for the possibility of no to reduced gene-flow from one lineage to the other. Altogether we
532 thus considered $8 \times 2 \times 3 = 48$ competing scenarios for the demographic and migration history of Central and Southern African
533 populations inferred with machine-learning Approximate Bayesian Computation. We conducted Random-Forest ABC
534 scenario choice procedures among groups of these 48 scenarios in turn for 54 different sets of five observed populations each,
535 represented as grey lines between population names, each population comprising five whole-genome sequenced individuals.
536 Sampled-population names considered in the 54 combinations are given below each tree-leave, and combinations are limited
537 to those pairing one RHGn with the specific eastern or western RHG population with whom they share complex socio-
538 economic interactions, as indicated in **TableT1x**. See the detailed description of scenarios and their grouping, their parameters
539 and their respective prior distributions used for ABC inference in **Material and Methods** and **TableT2x**.

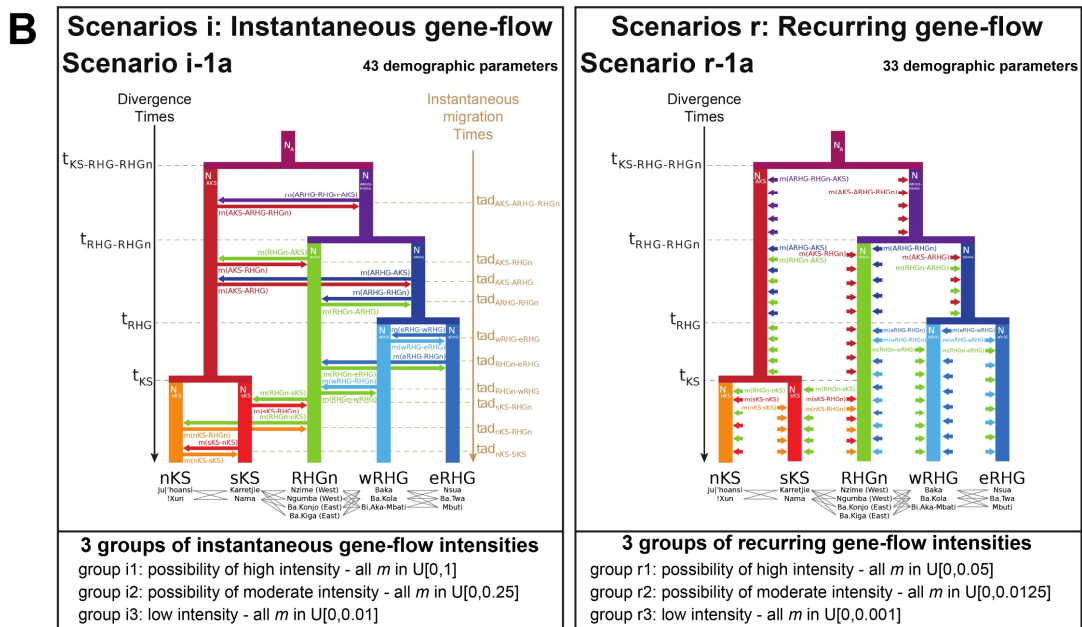
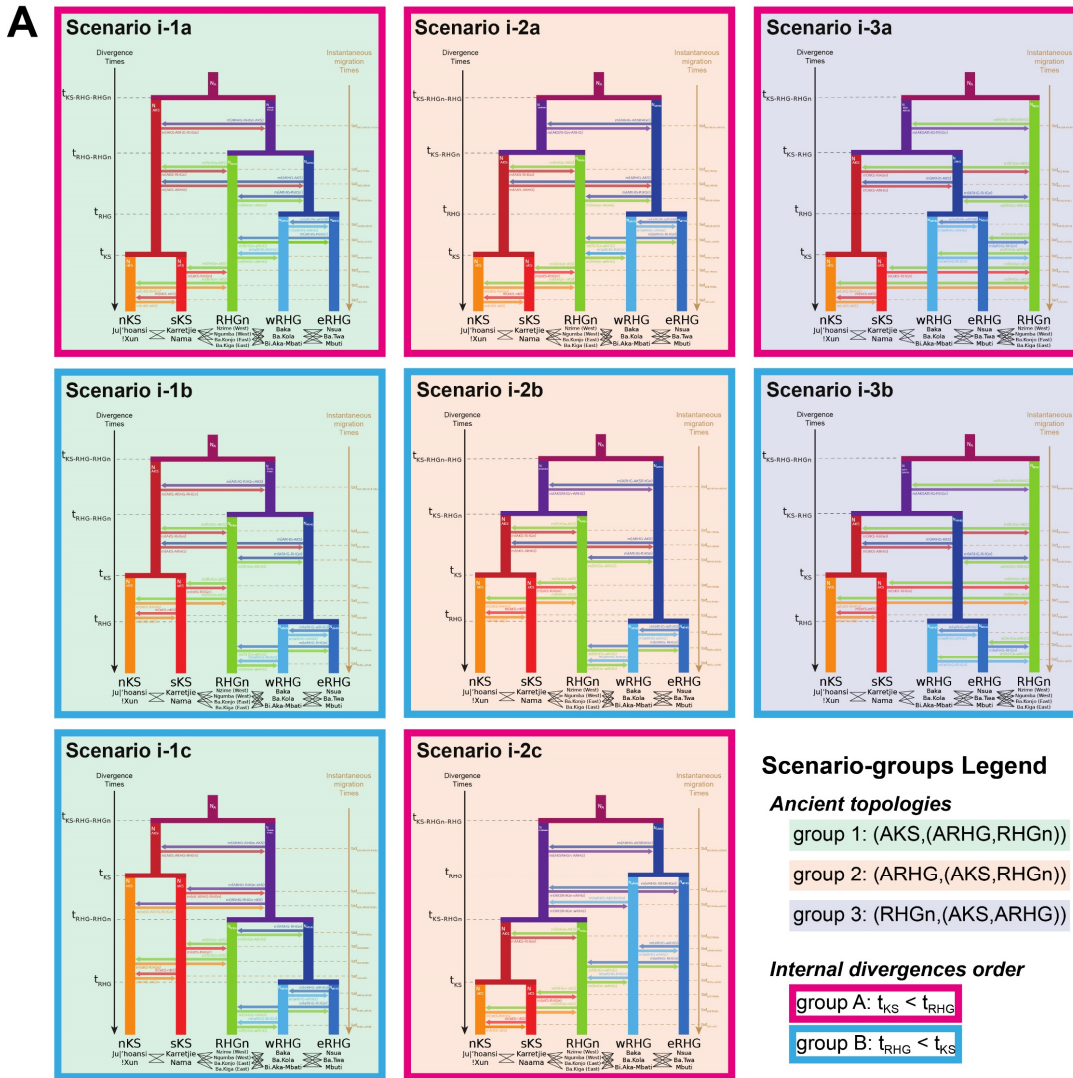
540

541

542

543
544
545
546

Figure F5x



547
548
549
550

TableT2x: 48 competing scenarios parameters' description and prior distributions
Codes for scenarios topologies and gene-flow processes and all scenario parameters are detailed in **FigureF5x** and in **Material and Methods**.

Parameter	Description	Prior	Scenario topology	Scenario gene-flow	Condition
t_{RHG}	Time of divergence between wRHG and eRHG	Uniform[1,15000]	1a-2a-2c-3a 1b-1c-2b-3b	i1-i2-i3 and r1-r2-r3	$t_{RHG} > t_{KS}$ $t_{KS} > t_{RHG}$
t_{KS}	Time of divergence between nKS and sKS	Uniform[1,15000]	1a-2a-2c-3a 1b-1c-2b-3b	i1-i2-i3 and r1-r2-r3	$t_{RHG} > t_{KS}$ $t_{KS} > t_{RHG}$
$t_{RHG-RHGn}$	Time of divergence between the RHGn lineage and the lineage ancestral to all RHG lineages	Uniform[1,15000]	1a-1b-1c	i1-i2-i3 and r1-r2-r3	$t_{RHG-RHGn} > t_{RHG}$
$t_{KS-RHGn}$	Time of divergence between the RHGn lineage and the lineage ancestral to all KS lineages	Uniform[1,15000]	2a-2b-2c	i1-i2-i3 and r1-r2-r3	$t_{KS-RHGn} > t_{KS}$
t_{KS-RHG}	Time of divergence between the KS and RHG ancestral lineages	Uniform[1,15000]	3a-3b	i1-i2-i3 and r1-r2-r3	$t_{KS-RHG} > t_{KS}$ and $t_{KS-RHG} > t_{RHG}$
$t_{KS-RHG-RHGn}$	Time of the divergence event ancestral to all lineages	Uniform[1,15000]	all 8 topol.	and r1-r2-r3	$t_{KS-RHG-RHGn} >$ all other times
$tad_{AKS-sKS}$	Time for both the unidirectional instantaneous gene-flow events between nKS and sKS	Uniform[1,15000] Uniform[0,1]	all 8 topol.	i1-i2-i3	$t_{KS} > tad_{AKS-sKS}$
$m(nKS-sKS)$ or $m(sKS-nKS)$	Instantaneous gene-flow intensity from sKS into nKS, or from nKS into sKS independently	Uniform[0,0.125] Uniform[0,0.01] Uniform[0,0.05]	all 8 topol.	i2 i3	x
$m(nKS-sKS)$ or $m(sKS-nKS)$	Recurring gene-flow intensity from sKS into nKS, or from nKS into sKS independently	Uniform[0,0.0125] Uniform[0,0.001]	all 8 topol.	r1 r2 r3	x
$tad_{AKS-RHGn}$	Time for both the unidirectional instantaneous gene-flow events between nKS and RHGn	Uniform[1,15000] Uniform[0,1]	all 8 topol.	i1-i2-i3	$t_{KS} > tad_{AKS-RHGn}$ and $t_{RHG-RHGn} > tad_{AKS-RHGn}$
$m(nKS-RHGn)$ or $m(RHGn-nKS)$	Instantaneous gene-flow intensity from RHGn into nKS, or from nKS into RHGn independently	Uniform[0,0.125] Uniform[0,0.01] Uniform[0,0.05]	all 8 topol.	i2 i3	x
$m(nKS-RHGn)$ or $m(RHGn-nKS)$	Recurring gene-flow intensity from RHGn into nKS, or from nKS into RHGn independently	Uniform[0,0.0125] Uniform[0,0.001]	all 8 topol.	r1 r2 r3	x
$tad_{AKS-RHGn}$	Time for both the unidirectional instantaneous gene-flow events between sKS and RHGn lineages	Uniform[1,15000] Uniform[0,1]	all 8 topol.	i1-i2-i3	$t_{KS} > tad_{AKS-RHGn}$ and $t_{RHG-RHGn} > tad_{AKS-RHGn}$
$m(sKS-RHGn)$ or $m(RHGn-sKS)$	Instantaneous gene-flow intensity from RHGn into sKS, or from sKS into RHGn independently	Uniform[0,0.125] Uniform[0,0.01] Uniform[0,0.05]	all 8 topol.	i2 i3	x
$m(sKS-RHGn)$ or $m(RHGn-sKS)$	Recurring gene-flow intensity from RHGn into sKS, or from sKS into RHGn independently	Uniform[0,0.0125] Uniform[0,0.001]	all 8 topol.	r1 r2 r3	x
$tad_{wRHG-eRHG}$	Time for both the unidirectional instantaneous gene-flow events between wRHG and eRHG	Uniform[1,15000] Uniform[0,1]	all 8 topol.	i1-i2-i3	$t_{RHG} > tad_{wRHG-eRHG}$
$m(wRHG-eRHG)$ or $m(eRHG-wRHG)$	Instantaneous gene-flow intensity from eRHG into wRHG, or from wRHG into eRHG independently	Uniform[0,0.125] Uniform[0,0.01] Uniform[0,0.05]	all 8 topol.	i2 i3	x
$m(wRHG-eRHG)$ or $m(eRHG-wRHG)$	Recurring gene-flow intensity from eRHG into wRHG, or from wRHG into eRHG independently	Uniform[0,0.0125] Uniform[0,0.001]	all 8 topol.	r1 r2 r3	x
$tad_{wRHG-RHGn}$	Time for both the unidirectional instantaneous gene-flow events between wRHG and RHGn	Uniform[1,15000] Uniform[0,1]	all 8 topol.	i1-i2-i3	$t_{RHG} > tad_{wRHG-RHGn}$ and $t_{KS-RHGn} > tad_{wRHG-RHGn}$
$m(wRHG-RHGn)$ or $m(RHGn-wRHG)$	Instantaneous gene-flow intensity from RHGn into wRHG, or from wRHG into RHGn independently	Uniform[0,0.125] Uniform[0,0.01] Uniform[0,0.05]	all 8 topol.	i2 i3	x
$m(wRHG-RHGn)$ or $m(RHGn-wRHG)$	Recurring gene-flow intensity from RHGn into wRHG, or from wRHG into RHGn independently	Uniform[0,0.0125] Uniform[0,0.001]	all 8 topol.	r1 r2 r3	x
$tad_{eRHG-RHGn}$	Time for both the unidirectional instantaneous gene-flow events between eRHG and RHGn	Uniform[1,15000] Uniform[0,1]	all 8 topol.	i1-i2-i3	$t_{RHG} > tad_{eRHG-RHGn}$ and $t_{KS-RHGn} > tad_{eRHG-RHGn}$
$m(eRHG-RHGn)$ or $m(RHGn-eRHG)$	Instantaneous gene-flow intensity from RHGn into eRHG, or from eRHG into RHGn independently	Uniform[0,0.125] Uniform[0,0.01] Uniform[0,0.05]	all 8 topol.	i2 i3	x
$m(eRHG-RHGn)$ or $m(RHGn-eRHG)$	Recurring gene-flow intensity from RHGn into eRHG, or from eRHG into RHGn independently	Uniform[0,0.0125] Uniform[0,0.001]	all 8 topol.	r1 r2 r3	x
$tad_{ARHG-RHGn}$	Time for both the unidirectional instantaneous gene-flow events between the lineage ancestral to nKS and sKS, and the lineage ancestral to wRHG and eRHG lineages	Uniform[1,15000] Uniform[0,1]	all 8 topol. except 2c	i1-i2-i3	$tad_{RHGn-ARHG} > t_{RHG}$
$m(RHGn-ARHG)$ or $m(ARHG-RHGn)$	Instantaneous gene-flow intensity from the ancestral lineage ARHG into RHGn, or from RHGn into ARHG independently	Uniform[0,0.125] Uniform[0,0.01] Uniform[0,0.05]	all 8 topol. except 2c	i2 i3	x
$m(RHGn-ARHG)$ or $m(ARHG-RHGn)$	Recurring gene-flow intensity from the ancestral lineage ARHG into RHGn, or from RHGn into ARHG independently	Uniform[0,0.0125] Uniform[0,0.001]	all 8 topol. except 2c	r1 r2 r3	x
$tad_{AKS-RHGn}$	Time for both the unidirectional instantaneous gene-flow events between the lineage ancestral to nKS and sKS, and the lineage ancestral to wRHG and eRHG lineages	Uniform[1,15000] Uniform[0,1]	all 8 topol. except 1c	i1-i2-i3	$tad_{AKS-RHGn} > t_{KS}$
$m(RHGn-AKS)$ or $m(AKS-RHGn)$	Instantaneous gene-flow intensity from the ancestral lineage AKS into RHGn, or from RHGn into AKS independently	Uniform[0,0.125] Uniform[0,0.01] Uniform[0,0.05]	all 8 topol. except 1c	i2 i3	x
$m(RHGn-AKS)$ or $m(AKS-RHGn)$	Recurring gene-flow intensity from the ancestral lineage AKS into RHGn, or from RHGn into AKS independently	Uniform[0,0.0125] Uniform[0,0.001]	all 8 topol. except 1c	r1 r2 r3	x
$tad_{AKS-ARHG}$	Time for both the unidirectional instantaneous gene-flow events between the lineage ancestral to nKS and sKS, and the lineage ancestral to wRHG and eRHG lineages	Uniform[1,15000] Uniform[0,1]	all 8 topol. except 1c-2c	i1-i2-i3	$tad_{AKS-ARHG} > t_{KS}$ and $tad_{AKS-ARHG} > t_{RHG}$
$m(AKS-ARHG)$ or $m(ARHG-AKS)$	Instantaneous gene-flow intensity from the ancestral lineage ARHG into the ancestral lineage AKS, or from AKS into ARHG independently	Uniform[0,0.125] Uniform[0,0.01]	all 8 topol. except 1c-2c	i2 i3	x

$m(\text{AKS-ARHG})$ or $m(\text{ARHG-AKS})$	Recurring gene-flow intensity from the ancestral lineage ARHG into the ancestral lineage AKS, or from AKS into ARHG independently	Uniform[0,0.05] Uniform[0,0.0125] Uniform[0,0.001]	all 8 topol. except 1c-2c	r1 r2 r3	x
$tad_{\text{AKS-ARHG-RHGn}}$	Time for both the unidirectional instantaneous gene-flow events between the nKS and the lineage ancestral to ARHG and RHGn lineages	Uniform[1,15000]	1c	i1-i2-i3	$t_{\text{KS}} > tad_{\text{AKS-ARHG-RHGn}}$ and $tad_{\text{AKS-ARHG-RHGn}} > t_{\text{RHG-RHGn}}$
$m(\text{nKS-ARHG-RHGn})$ or $m(\text{ARHG-RHGn-nKS})$	Instantaneous gene-flow intensity from the lineage ARHG-RHGn ancestral to lineages RHGn and ARHG into nKS, or from nKS into ARHG-RHGn independently	Uniform[0,1] Uniform[0,0.125] Uniform[0,0.01]	1c	i1 i2 i3	x
$m(\text{nKS-ARHG-RHGn})$ or $m(\text{ARHG-RHGn-nKS})$	Recurring gene-flow intensity from the lineage ARHG-RHGn ancestral to lineages RHGn and ARHG into nKS, or from nKS into ARHG-RHGn independently	Uniform[0,0.05] Uniform[0,0.0125] Uniform[0,0.001]	1c	r1 r2 r3	x
$tad_{\text{AKS-ARHG-RHGn}}$	Time for both the unidirectional instantaneous gene-flow events between the sKS and the lineage ancestral to ARHG and RHGn lineages	Uniform[1,15000]	1c	i1-i2-i3	$t_{\text{KS}} > tad_{\text{AKS-ARHG-RHGn}}$ and $tad_{\text{AKS-ARHG-RHGn}} > t_{\text{RHG-RHGn}}$
$m(\text{sKS-ARHG-RHGn})$ or $m(\text{ARHG-RHGn-sKS})$	Instantaneous gene-flow intensity from the lineage ARHG-RHGn ancestral to lineages RHGn and ARHG into sKS, or from sKS into ARHG-RHGn independently	Uniform[0,1] Uniform[0,0.125] Uniform[0,0.01]	1c	i1 i2 i3	x
$m(\text{sKS-ARHG-RHGn})$ or $m(\text{ARHG-RHGn-sKS})$	Recurring gene-flow intensity from the lineage ARHG-RHGn ancestral to lineages RHGn and ARHG into sKS, or from sKS into ARHG-RHGn independently	Uniform[0,0.05] Uniform[0,0.0125] Uniform[0,0.001]	1c	r1 r2 r3	x
$tad_{\text{AKS-RHGn-wRHG}}$	Time for both the unidirectional instantaneous gene-flow events between the wRHG and the lineage ancestral to AKS and RHGn lineages	Uniform[1,15000]	2c	i1-i2-i3	$t_{\text{RHG}} > tad_{\text{AKS-RHGn-wRHG}}$ and $tad_{\text{AKS-RHGn-wRHG}} > t_{\text{KS-RHGn}}$
$m(\text{wRHG-AKS-RHGn})$ or $m(\text{AKS-RHGn-wRHG})$	Instantaneous gene-flow intensity from the lineage AKS-RHGn ancestral to lineages AKS and RHGn into wRHG, or from wRHG into AKS-RHGn independently	Uniform[0,1] Uniform[0,0.125] Uniform[0,0.01]	2c	i1 i2 i3	x
$m(\text{wRHG-AKS-RHGn})$ or $m(\text{AKS-RHGn-wRHG})$	Recurring gene-flow intensity from the lineage AKS-RHGn ancestral to lineages AKS and RHGn into wRHG, or from wRHG into AKS-RHGn independently	Uniform[0,0.05] Uniform[0,0.0125] Uniform[0,0.001]	2c	r1 r2 r3	x
$tad_{\text{AKS-RHGn-eRHG}}$	Time for both the unidirectional instantaneous gene-flow events between the eRHG and the lineage ancestral to AKS and RHGn lineages	Uniform[1,15000]	2c	i1-i2-i3	$t_{\text{RHG}} > tad_{\text{AKS-RHGn-eRHG}}$ and $tad_{\text{AKS-RHGn-eRHG}} > t_{\text{KS-RHGn}}$
$m(\text{eRHG-AKS-RHGn})$ or $m(\text{AKS-RHGn-eRHG})$	Instantaneous gene-flow intensity from the lineage AKS-RHGn ancestral to lineages AKS and RHGn into eRHG, or from eRHG into AKS-RHGn independently	Uniform[0,1] Uniform[0,0.125] Uniform[0,0.01]	2c	i1 i2 i3	x
$m(\text{eRHG-AKS-RHGn})$ or $m(\text{AKS-RHGn-eRHG})$	Recurring gene-flow intensity from the lineage AKS-RHGn ancestral to lineages AKS and RHGn into eRHG, or from eRHG into AKS-RHGn independently	Uniform[0,0.05] Uniform[0,0.0125] Uniform[0,0.001]	2c	r1 r2 r3	x
$tad_{\text{AKS-ARHG-RHGn}}$	Time for both the unidirectional instantaneous gene-flow events between the lineage ancestral to KS and the lineage ancestral to RHG and RHGn lineages	Uniform[1,15000]	1a-1b-1c	i1-i2-i3	$tad_{\text{AKS-ARHG-RHGn}} > t_{\text{RHG-RHGn}}$ and $tad_{\text{AKS-ARHG-RHGn}} > t_{\text{KS-RHGn}}$ and $t_{\text{KS-RHGn}} > tad_{\text{AKS-ARHG-RHGn}}$
$m(\text{AKS-ARHG-RHGn})$ or $m(\text{ARHG-RHGn-AKS})$	Instantaneous gene-flow intensity from the lineage ARHG-RHGn ancestral to lineages ARHG and RHGn into the ancestral lineage AKS, or from AKS into ARHG-RHGn independently	Uniform[0,1] Uniform[0,0.125] Uniform[0,0.01]	1a-1b-1c	i1 i2 i3	x
$m(\text{AKS-ARHG-RHGn})$ or $m(\text{ARHG-RHGn-AKS})$	Recurring gene-flow intensity from the lineage ARHG-RHGn ancestral to lineages ARHG and RHGn into the ancestral lineage AKS, or from AKS into ARHG-RHGn independently	Uniform[0,0.05] Uniform[0,0.0125] Uniform[0,0.001]	1a-1b-1c	r1 r2 r3	x
$tad_{\text{AKSRHGn-ARHG}}$	Time for both the unidirectional instantaneous gene-flow events between the lineage ancestral to KS and RHGn, and the lineage ancestral to RHG lineages	Uniform[1,15000]	2a-2b-2c	i1-i2-i3	$tad_{\text{AKSRHGn-ARHG}} > t_{\text{KS-RHGn}}$ and $tad_{\text{AKSRHGn-ARHG}} > t_{\text{RHG}}$ and $t_{\text{KS-RHGn}} > tad_{\text{AKSRHGn-ARHG}}$
$m(\text{AKSRHGn-ARHG})$ or $m(\text{ARHG-AKSRHGn})$	Instantaneous gene-flow intensity from the ancestral lineage ARHG into the lineage AKSRHGn ancestral to AKS and RHGn, or from AKSRHGn into ARHG independently	Uniform[0,1] Uniform[0,0.125] Uniform[0,0.01]	2a-2b-2c	i1 i2 i3	x
$m(\text{AKSRHGn-ARHG})$ or $m(\text{ARHG-AKSRHGn})$	Recurring gene-flow intensity from the ancestral lineage ARHG into the lineage AKSRHGn ancestral to AKS and RHGn, or from AKSRHGn into ARHG independently	Uniform[0,0.05] Uniform[0,0.0125] Uniform[0,0.001]	2a-2b-2c	r1 r2 r3	x
$tad_{\text{AKSARHG-RHGn}}$	Time for both the unidirectional instantaneous gene-flow events between the lineage ancestral to KS and ARHG, and the RHGn lineage	Uniform[1,15000]	3a-3b	i1-i2-i3	$tad_{\text{AKSARHG-RHGn}} > t_{\text{KS-RHGn}}$ and $t_{\text{KS-RHGn}} > tad_{\text{AKSARHG-RHGn}}$
$m(\text{AKSARHG-RHGn})$ or $m(\text{RHGn-AKSARHG})$	Instantaneous gene-flow intensity from the lineage RHGn into the lineage AKSARHG ancestral to AKS and ARHG, or from AKSARHG into RHGn independently	Uniform[0,1] Uniform[0,0.125] Uniform[0,0.01]	3a-3b	i1 i2 i3	x
$m(\text{AKSARHG-RHGn})$ or $m(\text{RHGn-AKSARHG})$	Recurring gene-flow intensity from the lineage RHGn into the lineage AKSARHG ancestral to AKS and ARHG, or from AKSARHG into RHGn independently	Uniform[0,0.05] Uniform[0,0.0125] Uniform[0,0.001]	3a-3b	r1 r2 r3	x
N_{nKS}	Effective population size for the nKS lineage		all 8 topol.	i1-i2-i3 and r1-r2-r3	x
N_{sKS}	Effective population size for the sKS lineage		all 8 topol.	i1-i2-i3 and r1-r2-r3	x
N_{AKS}	Effective population size for the lineage ancestral to the nKS and the sKS		all 8 topol.	i1-i2-i3 and r1-r2-r3	x
N_{RHGn}	Effective population size for the RHGn lineage		all 8 topol.	i1-i2-i3 and r1-r2-r3	x
N_{wRHG}	Effective population size for the wRHG lineage	Uniform[10,100000]	all 8 topol.	i1-i2-i3 and r1-r2-r3	x
N_{eRHG}	Effective population size for the eRHG lineage		all 8 topol.	i1-i2-i3 and r1-r2-r3	x
N_{ARHG}	Effective population size for the lineage ancestral to the wRHG and the eRHG		all 8 topol.	i1-i2-i3 and r1-r2-r3	x
$N_{\text{ARHG-RHGn}}$	Effective population size for the lineage ancestral to the RHGn and the ARHG		1a-1b-1c	i1-i2-i3 and r1-r2-r3	x
$N_{\text{AKS-RHGn}}$	Effective population size for the lineage ancestral to the RHGn and the AKS		2a-2b-2c	i1-i2-i3 and r1-r2-r3	x

$N_{AKS-ARHG}$	Effective population size for the lineage ancestral to the AKS and the ARHG		3a-3b	i1-i2-i3 and r1-r2-r3	x
N_A	Effective population size for the lineage ancestral to all lineages		all 8 topol.	i1-i2-i3 and r1-r2-r3	x
Mutation rate	Fixed mutation rate	$1.25 \cdot 10^{-8}$	all 8 topol.	i1-i2-i3 and r1-r2-r3	x
Recombination rate	Fixed recombination rate	$1 \cdot 10^{-8}$	all 8 topol.	i1-i2-i3 and r1-r2-r3	x
Transition to transversion ratio	Fixed transition to transversion rate	0.33	all 8 topol.	i1-i2-i3 and r1-r2-r3	x

551
552
553

554
555
556

TableT3x: Summary-statistics used in Random-Forest ABC scenario-choice and Neural-Network ABC posterior-parameter estimation.

Name of statistic (or group of statistics)	What was computed	Number of values	RF-ABC scenario-choice	NN-ABC posterior-parameter estimation
Total number of biallelic sites among the five populations	sum	1	yes	yes
Total number of multiallelic sites among the five populations	sum	1	yes	yes
ASD distance within population for each five population separately	average and variance	10	yes	no
First five dimensions of the projection of the ASD matrix within population for each five population separately	average and variance	50	yes	yes
Proportion of biallelic sites within each five population separately	sum in a population divided by total number of biallelic sites in the dataset (no missing data)	5	yes	yes
Proportion of homozygous sites within each five population separately	sum in a population divided by total number of biallelic sites in the dataset (no missing data)	5	yes	yes
Proportion of homozygous ancestral sites within each five population separately	sum in a population divided by total number of homozygous sites in the population (no missing data)	5	yes	no
Proportion of private biallelic sites within each five population separately	sum in a population divided by total number of biallelic sites in the dataset (no missing data)	5	yes	yes
Frequency of the minor allele of the private biallelic sites within each five population separately	average and variance	10	yes	no
Expected heterozygosity at a segregating site within each five population separately	average and variance	10	yes	no
Tajima's D within each five population separately	average and variance for chromosomes with at least one variant	10	yes	yes
Unfolded site frequency spectrum (SFS) within each five population separately	proportion for each of 9 classes	45	yes	no
Watterson's theta within each five population separately	sum of biallelic sites divided by the harmonic number of the sample size (n=10)	5	yes	no
Total number of runs of homozygosity (ROH) within each five population separately	sum	5	yes	no
Total number of ROH by length class (<200 kbp, 200-500 kbp, >500 kbp) within each five population separately	sum	15	yes	yes
ROH length for each length class within each five population separately	average	15	yes	no
Total ROH length within each five population separately	average and variance	10	yes	no
ASD distance between each pair of five populations	average and variance	20	yes	no
First five dimensions of the projection of the ASD matrix between each pair of five populations	average and variance	100	yes	yes
Pairwise FST between each pair of five population	Pairwise FST (Weir and Cockerham 1984)	10	yes	yes
TOTAL number of statistics used in ABC inference			337	202

557
558
559
560

561 **2.B. Are we able to mimic observed data with simulations under the 48 competing scenarios?**

562 In order to conduct ABC inferences, we first checked whether we were able to simulate data, under the
563 48 competing scenarios, for which vectors of summary statistics could mimic those obtained separately
564 from the 54 different combinations of five sampled populations each.

565 We first conducted goodness-of-fit permutation tests, separately for the 54 combinations of five
566 populations, and found that the observed vectors of 337 summary statistics were never significantly
567 different (54 goodness-of-fit permutation p-values > 0.56), from the 240,000 vectors of 337 summary
568 statistics computed from simulations under the 48 competing scenarios. Second, we computed a two-
569 dimensional PCA on the vectors of summary-statistics obtained from simulations on which we
570 projected, separately, the 54 vectors of these summary-statistics computed on observed data, and found
571 that observations each fell well within the space of simulations (**SupplementaryFigureSF1x**).

572 Both results show that we empirically were able to simulate data for which summary-statistics
573 were reasonably close to the observed ones, at least in parts of the parameter space used for simulations
574 under the 48 competing scenarios, for each one of the 54 separate combinations of five sampled
575 populations each. We could therefore reasonably proceed with Random Forest ABC scenario-choice
576 and Neural Network ABC posterior-parameter estimation procedures.

577

578 **2.C. Which historic scenario best explains genomic patterns in extant African populations?**

579 2.C.1. Instantaneous or recurring gene-flows of which intensity among African lineages?

580 **Panel A** in **FigureF6x** shows that gene-flow processes among recent or ancient genetic lineages most
581 likely occurred during very limited “instantaneous” periods of time, rather than recurrently, throughout
582 the entire evolutionary history of Central and Southern African populations. Indeed, for each one of the
583 54 different combinations of five sampled populations, the group of 24 scenarios considering
584 instantaneous gene-flow processes provides vectors of summary statistics systematically closest to the
585 observed ones, whichever the tree-topology or intensity of the gene-flow considered.

586 Furthermore, we find that we imperatively need these gene-flow processes to be potentially high
587 in order to best explain the observed data for each one of the 54 combinations of population samples,
588 whichever the tree-topology and gene-flow process (**Panel B** in **FigureF6x**), whichever the tree-
589 topology but considering gene-flow processes separately (**Panel C** in **FigureF6x**), or even considering
590 24 competing scenarios under instantaneous gene-flow processes only, whichever the tree-topology
591 (**SupplementaryFigureSF6xPart1**). Finally, note that when considering all 48 competing-scenarios
592 separately, a very challenging task *a priori* given the number of competing scenarios and the high level
593 of nestedness among scenarios, the winning scenarios predicted by the RF-ABC procedure is, among
594 the 54 tests conducted, systematically found for scenarios considering an instantaneous gene-flow
595 process allowing for possibly highly intense gene-flow rates (**SupplementaryFigureSF6xPart2-3**).

596 Overall, for all the above analyses, the *a priori* discriminatory powers of the RF, estimated as
597 cross-validation procedures based on simulations used as pseudo-observed data, were overall good
598 across groups of scenarios, despite increased confusion among scenarios based on gene-flow intensity
599 within groups of instantaneous or recurring processes, as expected due to increased scenario-nestedness
600 among classes of no-to-low, no-to-moderate, or no-to-high intensities
601 (**SupplementaryFigureSF6xPart0-PanelA-B-C, SupplementaryFigureSF6xPart1-2-3**).

602

603

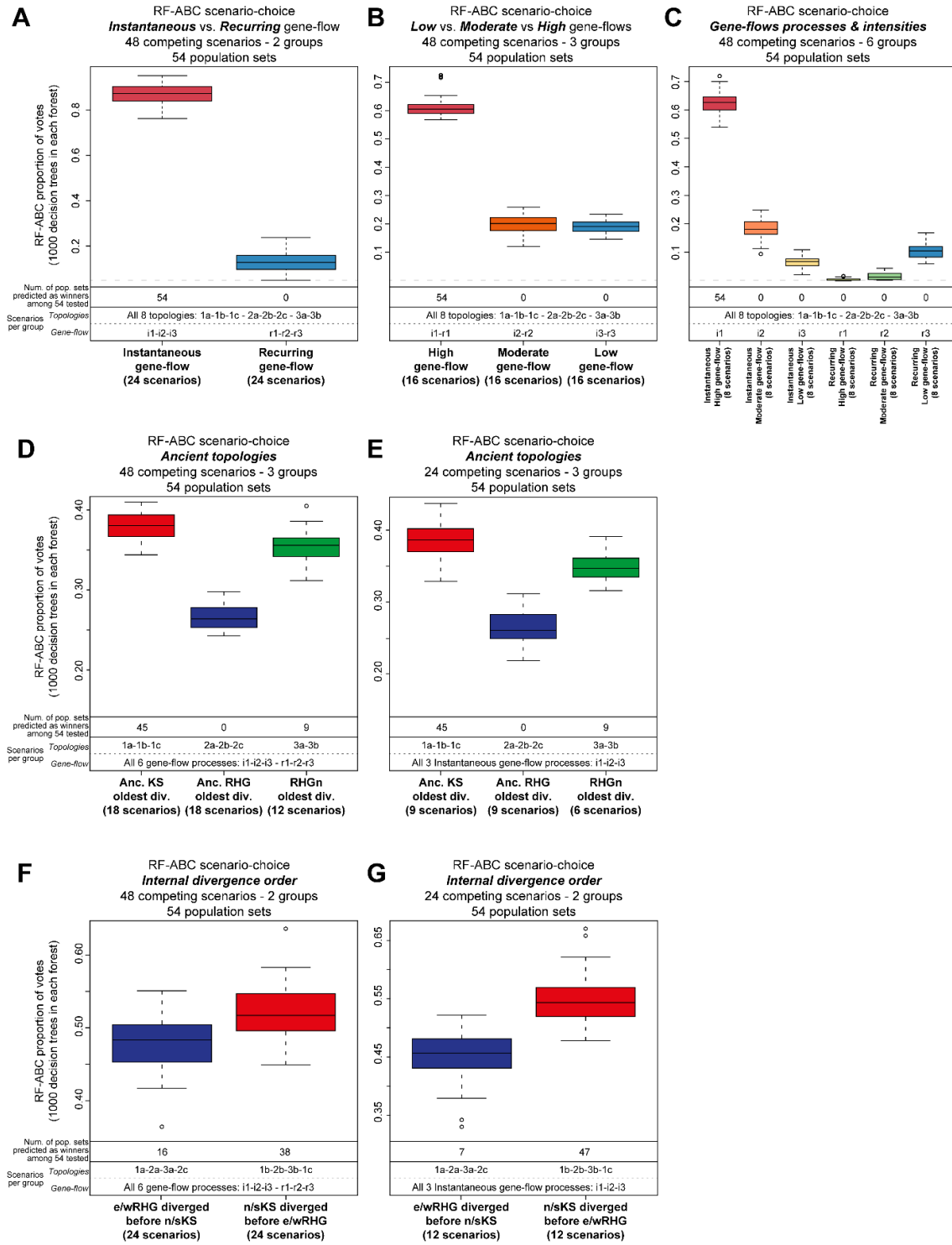
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637

FigureF6x: Random-Forest ABC scenario-choices among 48 competing scenarios.

Random Forest Approximate Bayesian Computation scenario-choice results (Pudlo *et al.*, 2016; Estoup *et al.*, 2018), are conducted, for each analysis in each panel, separately for 54 different combinations of five Central and Southern African sampled populations each. Posterior proportions of votes obtained for each 54 sampled-populations combinations, and for RF-ABC analysis in each panel respectively, are provided as box-plots indicating the median in between the first and third quartile of the box-limits, whiskers extending to data points no more than 1.5 times the interquartile range of the distribution, and empty circles for all more extreme points beyond this limit, if any. In each RF-ABC analysis in each panel, proportion of votes (indicated in the y-axis) are obtained using 1,000 decision trees in the random forest for each group of competing scenarios (indicated in the x-axis table and labels), using 5,000 simulations per scenario and 337 summary-statistics (see **Material and Methods** and **TableT3x** for details). Detailed competing-scenarios grouped in each group are indicated in the two bottom lines of the x-axis table, with topology and gene-flow processes codes provided in **FigureF5x** and explained in detail in **Material and Methods**. The top-line of the x-axis table indicates, for each group of competing scenarios in each analysis, the number of times the corresponding group of scenarios was predicted as the winning one among the 54 different combinations of five sampled-populations. (A) Scenario-choice results for the two groups of instantaneous or recurring gene-flow processes (24 scenarios and 120,000 simulations in each two competing groups), all gene-flow intensities and all topologies “being equal”. (B) Scenario-choice results for the three groups of intensities for gene-flow processes (16 scenarios and 80,000 simulations in each three competing groups), all topologies and both instantaneous and recurring gene-flow processes “being equal”. (C) Scenario-choice results for the six groups of gene-flow processes and intensities separately (8 scenarios and 40,000 simulations in each six competing groups), all topologies “being equal”. (D) Scenario-choice results for the three groups of scenario topologies differing in the order of ancient lineages divergence (differing number of scenarios in each three competing groups are evened by randomly sampling the same number of simulations, equal to 60,000, in each group, see **Material and Methods**), all gene-flow processes and intensities “being equal”. (E) Scenario-choice results for the same test as (D) restricted to the 24 competing scenarios considering instantaneous gene-flow processes only (differing number of scenarios in each three competing groups are evened by randomly sampling the same number of simulations, equal to 30,000, in each group, see **Material and Methods**), all gene-flow intensities “being equal”. (F) Scenario-choice results for the two groups of scenario topologies differing in the relative order of divergences between Northern and Southern KS, and Western and Eastern RHG lineages, respectively (24 competing scenarios and 120,000 simulations in each two groups), all gene-flow processes and intensities “being equal”. (G) Scenario-choice results for the same test as (F) restricted to the 24 competing scenarios considering instantaneous gene-flow processes only (12 competing scenarios and 60,000 simulations in each two groups), all gene-flow intensities “being equal”.

638
639
640
641
642
643

Figure F6x



644 2.C.2. Ancient tree-topology among African lineages?

645 **FigureF6x** shows that, overall, the RF-ABC scenario-choice procedures vote in favor of an ancient
646 tree-topology where the ancestral Khoe-San lineage diverged first followed by the divergence between
647 the ancient Rainforest Hunter-Gatherer lineage and that of their neighbors (Scenarios tree-topologies
648 1a-1b-1c, **FigureF5x**), in a large majority (45/54) of the 54 combinations of observed populations tested
649 here, whichever the gene-flow process and intensity (**FigureF6x-PanelD**), or when considering
650 instantaneous gene-flow processes, only, whichever their intensities (**FigureF6x-PanelE**). Note that
651 the tree-topologies where the RHGn lineage diverge first from the lineage ancestral to that of the RHG
652 and the KS (Scenarios tree-topologies 3a-3b, **FigureF5x**), is winning in a minority of the tests (9 out of
653 54 combinations of sampled populations), whichever the gene-flow process and/or intensity.
654 Interestingly, we found that tree-topologies where the ancient lineage for RHG populations diverges
655 first from the two others (Scenarios tree-topologies 2a-2b-2c, **FigureF5x**), is never favored in our
656 analyses, whichever the gene-flow process and/or intensity.

657

658 2.C.3. Recent tree-topology among African lineages?

659 Finally, **FigureF6x** shows that Northern and Southern Khoe-San extant lineages likely diverged from
660 one-another before the divergence between Eastern and Western Rainforest Hunter-Gatherer lineages
661 (Scenarios tree-topologies 1b-1c-2b-3b, **FigureF5x**), whichever the ancient tree-topology and gene-
662 flow process and intensity (**FigureF6x-PanelF**), and whichever the ancient tree-topology for
663 instantaneous gene-flow processes only, all classes of intensities “being equal” (**FigureF6x-PanelG**).
664 Interestingly, while the “b” group of tree-topologies wins for 38 out of 54 combinations of sampled
665 populations, this majority of RF-ABC votes is larger (47/54) when considering only the group of 24
666 competing scenarios for which gene-flows are instantaneous rather than recurring.

667

668 2.C.4. Conclusion: intersecting scenario-choice results for the history of Africa.

669 Altogether, when intersecting RF-ABC scenario-choice inferences in groups of gene-flow processes
670 and tree-topologies (**FigureF6x**), with inferences considering all competing-scenarios separately
671 (**SupplementaryFigureSF6x-Part2-3**), and among 54 combinations of five sampled populations, we
672 find that Scenario i1-1b (**FigureF5x**), systematically wins in all conformations of scenario-choice tests
673 for 25 combinations of populations (**SupplementaryTableST3x**), while the second best scenario
674 (Scenario i1-3b) wins consistently across tests for only four combinations of population out of 54.

675

676 Therefore, our results point to a scenario for the history of Central and Southern African
677 populations where the ancestral KS lineage diverged first, followed by the divergence between the
678 ancestral RHG and the RHGn lineages, with a divergence between the extant Northern and Southern
679 KS lineages occurring independently before that of the extant Eastern and Western RHG lineages. Most
680 importantly, the winning scenario necessarily involves gene-flow events among all pairs of lineages
681 throughout history which occurred during relatively short “instantaneous” periods of time, rather than
682 recurrently over larger periods of time. Furthermore, these instantaneous gene-flow events must each
683 have been possibly of high intensity, rather than all limited to moderate or low intensities, and possibly
684 asymmetrical between pairs of lineages, to best explain the observed data. Finally, our results also
685 highlight that considering different combinations of extant populations separately in the analyses may
686 provide contrasted results where alternative evolutionary scenarios may best explain observations.

686

687 For conservativeness and to further consider the genetic diversity of extant Central and Southern
688 African populations at a local or regional scale, we will henceforth conduct all Neural-Network ABC
689 posterior-parameter estimations separately for the 25 combinations of sampled populations
690 systematically providing the Scenario i1-1b as the winner.

690

691

692 **2.D. Which scenario-parameters best explain genomic patterns in extant African populations?**

693 2.D.1. Divergence times.

694 Divergence times are most often well estimated by our Neural Network ABC posterior parameter
695 inferences (**FigureF7x** and **SupplementaryTableST4x**), with relatively narrow 90% Credibility
696 Intervals (CI) and posterior distributions very often largely departing from the priors. Interestingly, we
697 find relatively consistent posterior estimates across the 25 different combinations of five populations
698 for, separately, the Eastern and Western RHG divergence between 383 generations before present (gbp)
699 (mode point estimate, 90%CI=[220-934]) and 2174 gbp (90%CI=[1,226-4,373],
700 **SupplementaryTableST4x**), and for the Northern and Southern KS divergence between 1,486 gbp
701 (mode point estimate, 90%CI=[844-3,798]) and 6,201 gbp (90%CI=[4,465-8,126]). Combining the
702 posterior distributions for these two parameters (**FigureF7x-PanelA-B**, **TableT4x**), respectively,
703 provides, synthetically, a modal point-estimate divergence time between eRHG and wRHG of 892 gbp
704 (90%CI=[422-2,531]), and one between nKS and sKS, largely more ancient, having occurred some
705 2,623 gbp (90%CI=[1,481-5,763]).

706 In a more remote past, the divergence time between the lineage ancestral to all RHG populations
707 and the RHG neighbors' lineage is also most often well estimated, based on CI-width and departure
708 from prior distributions, for each 25 combinations of populations respectively, albeit results are more
709 variable across sets of population combinations than for more recent divergence times (**FigureF7x-
710 PanelC** and **SupplementaryTableST4x**). Indeed, ancient RHG and RHGn lineages diverged between
711 3,904 gbp (mode point estimate, 90%CI=[2,828-6,636]) and 10,727 gbp (90%CI=[9,268-12,736],
712 **SupplementaryTableST4x**) across population combinations, and synthetically combining results
713 together 6,726 gbp (90%CI=[3,876-10,748], **FigureF7x-PanelC** and **TableT4x**).

714 Finally, we also find relatively variable posterior estimates for the most ancient divergence time in
715 our tree-topology, between the ancestral lineages to all extant KS lineages and the ancestral lineages to
716 all RHG and the RHG neighbors; with, again, posterior-parameter distributions across 25 combinations
717 of sampled populations very often satisfactorily estimated, with relatively reduced 90%CI and
718 substantial departure from the priors (**FigureF7x-PanelD**, **SupplementaryTableT4x**). Overall, we find
719 that the original most ancient divergences in our tree-topology occurred some 12,117 gbp
720 (90%CI=[8,875-14,538]), when all results are combined together synthetically (**FigureF7x-PanelD** and
721 **TableT4x**).

722 Altogether, the relatively large variation in divergence times posterior estimates further back in
723 time across combinations of sampled populations empirically highlights that demographic and historical
724 reconstructions largely depend on the specific samples considered when investigating highly
725 differentiated Sub-Saharan populations. Our results thus explicitly advocate for caution when
726 summarizing results obtained across population sets and further likely explain the, sometimes, apparent
727 discrepancies that arise across studies considering different sample sets and/or artificially merging
728 several samples from relatively differentiated populations (e.g. (Fan *et al.*, 2023; Ragsdale *et al.*, 2023)).

729
730
731

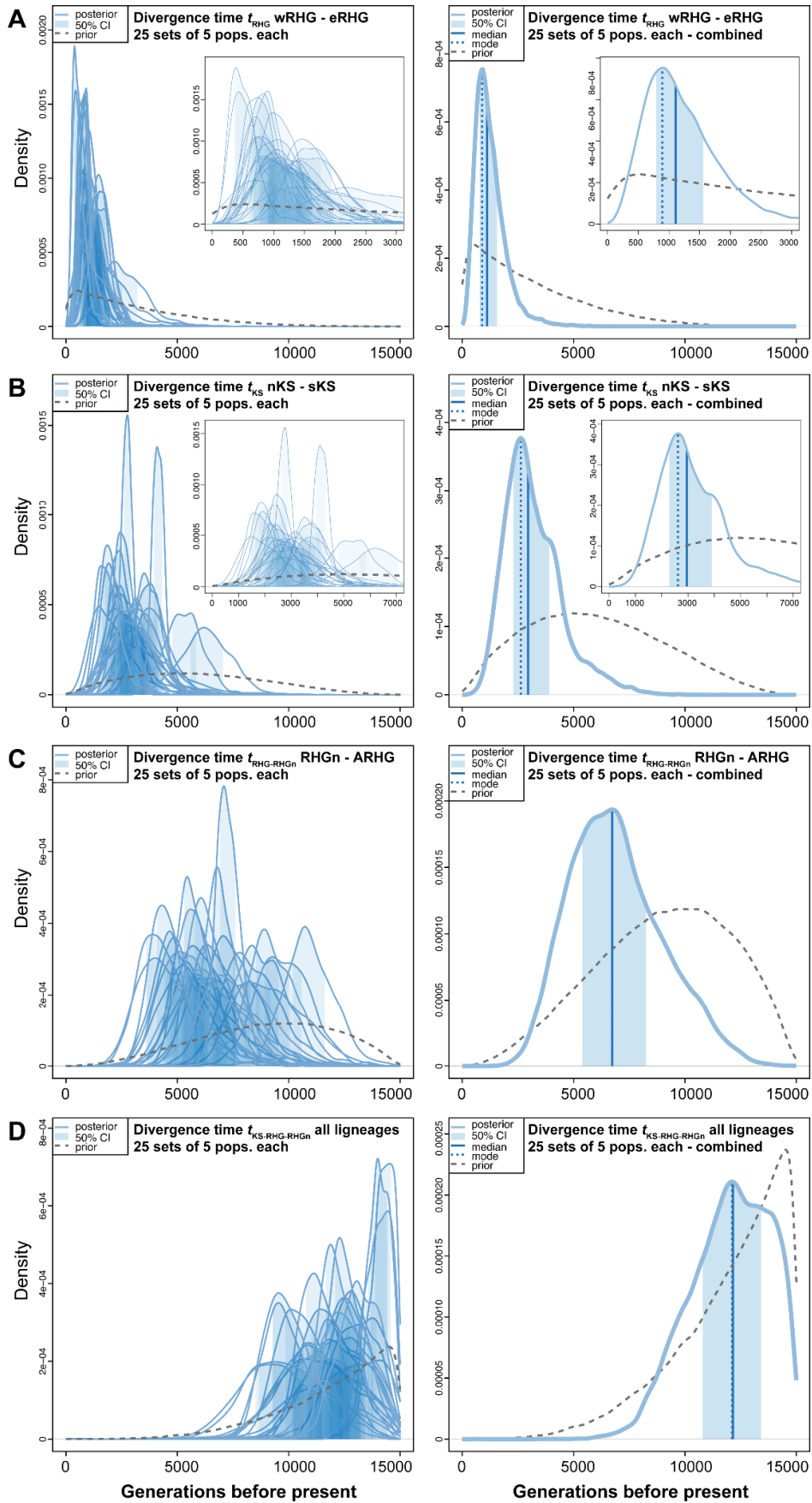
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758

FigureF7x: Neural Network ABC posterior parameter distributions of divergence times.

Neural Network Approximate Bayesian Computation posterior parameter joint estimations (Blum and François, 2010; Csilléry, François and Blum, 2012), of topological divergence-times (in generations before present) for 25 sets of five Central and Southern African populations for which the winning scenario identified by RF-ABC was Scenario i1-1b (**FigureF5x**). NN ABC posterior parameter estimation procedures were conducted using 100,000 simulations under Scenario i1-1b, each simulation corresponding to a single vector of parameter values drawn randomly from prior distributions provided in **TableT2x**. We considered 42 neurons in the hidden layer of the NN and a tolerance level of 0.01, corresponding to the 1,000 simulations providing summary-statistics closest to the observed one. NN posterior estimates are based on the logit transformation of parameter values using an Epanechnikov kernel between the corresponding parameter's prior bounds (see **Material and Methods** and **TableT2x**). Posterior parameter densities are represented with solid blue lines. 50% Credibility Intervals are represented as the light blue area under the density. The median and mode values are represented as a solid and dotted blue vertical line, respectively. Parameter prior distributions are represented as dotted gray lines. For all panels, the left plots represent the NN-ABC posterior parameter distributions for each 25 sets of five Central and Southern African populations winning under Scenario i1-1b, separately (**SupplementaryTableST3x** and **SupplementaryTableST4x**). For all panels, the right plots represent a single parameter posterior distribution obtained from combining the 25 posterior distributions together. (A) Results for parameter t_{RHG} corresponding to the split time between the Western and Eastern Rainforest Hunter-Gatherer (RHG) lineages (**FigureF5x**). (B) Results for parameter t_{KS} corresponding to the split time between Northern Khoe-San (nKS) and Southern Khoe-San (sKS) lineages (**FigureF5x**). (C) Results for parameter $t_{RHG-RHGn}$ corresponding to the split time between the Rainforest Hunter-Gatherer neighboring population lineage (RHGn) and the lineage ancestral to Western and Eastern RHG (**FigureF5x**). (D) Results for parameter $t_{KS-RHG-RHGn}$ corresponding to the split time between the lineage ancestral to KS populations and the lineage ancestral to all RHG and RHGn lineages; this event thus corresponding to the split time of all lineages in the history of Central and Southern African populations (**FigureF5x**). Results for all left panels are summarized in **TableT4x**.

759
760
761
762

Figure F7x



763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779

TableT4x: NN-ABC posterior parameter estimation of all parameters in Scenario i1-1b for results from 25 sets of five populations each, combined altogether.

Entire posterior distributions for each 25 analyses separately, for which summary statistics are presented in **SupplementaryTableT2x**, are combined together before re-computing the Mode, Median, Mean, 50%CI and 90%CI presented in this table. Therefore, values in the table correspond to posterior densities plotted in the right panels of **FigureF7x**, **FigureF8x**, and **SupplementaryFigureF7xPart1-Part2**. Parameter definitions and priors are provided in **TableT2x** and represented graphically in the corresponding scenario panel of **FigureF5x**. Values in this table are used for the synthetic schematic representation of the demographic and migration history of Central and Southern African populations presented in **FigureF9x**. Values in italic are not satisfactorily departing from the priors for the results from 25 sets of five populations each when combined together. In the case of the parameters N_{AKS} , N_{ARHG} , and $N_{ARHG-RHGn}$, this is due to parameter un-identifiability for all 25 analyses separately as almost none of these posterior distributions depart from their respective priors (**SupplementaryFigureF7xPart1**). In the case of all the migration rates parameters, **SupplementaryFigureF7xPart2** clearly shows both an overall lack of identifiability of parameters in the majority of the 25 analyses, separately, and a large variability of posterior estimates among those sets of analyses for which posterior distributions satisfactorily depart from the priors.

Parameter	Mode	Median	Mean	50% CI	90% CI
N _{nKS}	14184	23881	29908	14721-39894	7825-73423
N _{sKS}	11042	19672	26050	11059-35737	4623-68767
N _{wRHG}	9461	20586	26551	11506-36934	5074-67317
N _{eRHG}	12592	15735	19968	10059-24368	3993-53065
N _{RHGn}	18022	28739	33965	18076-45474	9150-77157
N _{AKS}	86572	61359	59978	41254-80294	19523-94743
N _{ARHG}	57397	55236	54989	36771-73803	16922-91847
N _{ARHG-RHGn}	46968	48438	48585	30042-66886	9894-88177
N _A	18507	24995	28406	17135-35537	9399-60616
t _{KS}	2623	2960	3199	2287-3921	1481-5763
t _{RHG}	892	1111	1255	773-1574	422-2531
t _{RHG-RHGn}	6726	6729	6931	5393-8271	3876-10748
t _{KS-RHG-RHGn}	12117	12152	12001	10797-13422	8875-14538
tad _{nKS-sKS}	395	532	597	347-786	147-1232
tad _{nKS-RHGn}	1383	1659	1835	1140-2333	591-3710
tad _{sKS-RHGn}	1142	1378	1533	928-1970	418-3165
tad _{wRHG-RHGn}	302	348	418	227-550	97-944
tad _{eRHG-RHGn}	393	520	616	342-801	171-1358
tad _{wRHG-eRHG}	411	555	653	367-851	157-1429
tad _{RHGn-ARHG}	2756	3193	3523	2352-4347	1369-6828
tad _{AKS-RHGn}	4490	4588	4831	3731-5590	2828-7892
tad _{AKS-ARHG}	4408	4458	4730	3629-5498	2658-7963
tad _{AKS-ARHG-RHGn}	9679	9532	9524	7888-11066	6358-12919
m(nKS-sKS)	0.6727	0.6018	0.5685	0.3770-0.7830	0.0892-0.9385
m(sKS-nKS)	0.7330	0.6314	0.6072	0.4481-0.7911	0.1851-0.9346
m(nKS-RHGn)	0.4827	0.5336	0.5324	0.3595-0.7206	0.1277-0.9082
m(RHGn-nKS)	0.3668	0.4842	0.4966	0.2930-0.7092	0.0875-0.9153
m(sKS-RHGn)	0.4810	0.5093	0.5082	0.3190-0.6977	0.1062-0.9089
m(RHGn-sKS)	0.4920	0.4525	0.4587	0.2549-0.6432	0.0771-0.8880
m(RHGn-wRHG)	0.6461	0.5674	0.5482	0.3609-0.7438	0.1142-0.9258
m(wRHG-RHGn)	0.3925	0.4341	0.4447	0.2729-0.6058	0.0873-0.8413
m(RHGn-eRHG)	0.3305	0.4140	0.4365	0.2319-0.6291	0.0618-0.8750
m(eRHG-RHGn)	0.6072	0.5726	0.5507	0.3583-0.7577	0.1073-0.9243
m(wRHG-eRHG)	0.5034	0.4724	0.4766	0.2736-0.6758	0.0802-0.8891
m(eRHG-wRHG)	0.5453	0.5396	0.5340	0.3291-0.7465	0.1069-0.9388
m(RHGn-ARHG)	0.5802	0.4897	0.4851	0.2797-0.6870	0.0781-0.8904
m(ARHG-RHGn)	0.4081	0.5080	0.5110	0.2993-0.7299	0.0853-0.9345
m(ARHG-AKS)	0.4734	0.5082	0.5056	0.3033-0.7097	0.0850-0.9212
m(AKS-ARHG)	0.5824	0.5220	0.5121	0.3036-0.7210	0.0773-0.9255
m(RHGn-AKS)	0.7032	0.5235	0.5139	0.3053-0.7261	0.0887-0.9158
m(AKS-RHGn)	0.4066	0.5027	0.5066	0.3136-0.7197	0.0929-0.8935
m(AKS-ARHG-RHGn)	0.5357	0.5008	0.4916	0.3017-0.6756	0.0834-0.8904
m(ARHG-RHGn-AKS)	0.3359	0.4501	0.4680	0.2424-0.6856	0.0663-0.9148

781
782
783
784

785 2.D.2. Effective population sizes.

786 Effective population sizes (N_e) for all recent Northern and Southern KS, Western and Eastern RHG,
787 and RHG neighbors' lineages are often reasonably well estimated with relatively reduced 90%CI and
788 substantial departure from their priors, for each 25 sets of sampled population combinations
789 (**SupplementaryFigureSF7xPart1, SupplementaryTableST4x**). Combining posterior distributions
790 among the 25 separate tests, we find N_e posterior estimates ranging from modal point-estimates of
791 9,461 diploid effective individuals (90%CI=[5,074-67,317]) in the wRHG to 18,022 (90%CI=[9,150-
792 77,157]) in the RHGn extant lineages (**SupplementaryFigureSF7xPart1, TableT4x**). Importantly, we
793 were unable to satisfactorily estimate, for almost all 25 combinations of sampled populations, the three
794 effective population sizes for, respectively, the ancient KS lineage, the ancient RHG lineage, and the
795 lineage ancestral to all RHG and RHG neighbors' extant lineages (**TableT4x,**
796 **SupplementaryFigureSF7xPart1, SupplementaryTableST4x**).

797 Conversely, our posterior estimates of the effective size of the lineage most ancestral to all our
798 populations was satisfactorily estimated, with relatively narrow 90%CI and substantial departure from
799 the prior distributions, for almost all 25 sets of population combinations. Combining posterior
800 distributions, despite a noteworthy variation across the 25 tests (**SupplementaryFigureSF7xPart1,**
801 **SupplementaryTableST4x**), we find an ancestral effective size for the lineage ancestral to all our
802 Central and Southern African extant populations of 18,507 diploid effective individuals
803 (90%CI=[9,399-60,616]), the largest posterior estimate across all recent and ancient lineages in our
804 analyses (**TableT4x**). Note that we considered large priors (Uniform[100-100,000] diploid individuals)
805 for constant effective population sizes in all lineages with possible changes at each divergence time, for
806 simplicity. Therefore, the posterior estimates here found may be different in future procedures
807 considering more complex effective demographic regimes likely to have occurred in certain lineages,
808 such as possible bottlenecks and/or population expansions (e.g. (Patin *et al.*, 2014; Schlebusch *et al.*,
809 2020; Seidensticker *et al.*, 2021)).

810

811 2.D.3. Instantaneous gene-flow times.

812 In between each lineage divergence times, we estimated the time of occurrence of instantaneous and
813 potentially asymmetric gene-flow exchanges across pairs of lineages (**FigureF5x**), for each 25
814 combinations of sampled populations separately. For four out of the six recent gene-flow times between
815 pairs of extant lineages ($tad_{nKS-sKS}$, $tad_{wRHG-eRHG}$, $tad_{wRHG-RHGn}$, $tad_{eRHG-RHGn}$, **FigureF5x, TableT2x**), we
816 find relatively consistent posterior estimates across the 25 tests, almost all with relatively reduced
817 90%CI and substantial departure from the priors (**FigureF8x-PanelA-D, SupplementaryTableST4x,**
818 **TableT4x**). Note that, while we also obtain similarly satisfactory posterior distributions among the 25
819 tests for gene-flows among Northern KS and RHGn lineages ($tad_{nKS-RHGn}$), and among Southern KS and
820 RHGn lineages ($tad_{sKS-RHGn}$), respectively, we found much larger variance across the 25 different
821 population sets (**FigureF8x-PanelE-F, SupplementaryTableST4x**).

822 Interestingly, we also obtain satisfactory posterior estimates of gene-flow timing among ancient
823 lineages for almost all 25 sets of population combinations ($tad_{ARHG-RHGn}$, $tad_{AKS-RHGn}$, $tad_{AKS-ARHG}$, tad_{AKS-}
824 $ARHG-RHGn$, **FigureF5x, TableT2x**). Nevertheless, note that posterior estimates are increasingly variable
825 from one set of sampled populations to the other as the estimated parameter is further back in time in
826 the tree-topology (**FigureF8x-PanelG-J, TableT4x**).

827 These results further illustrate that different sample sets used in genetic inferences can provide
828 substantially differing demographic inferences, sometimes for even the most recent gene-flow events,
829 but most often for more ancient gene-flow events across ancestral lineages.

830

831

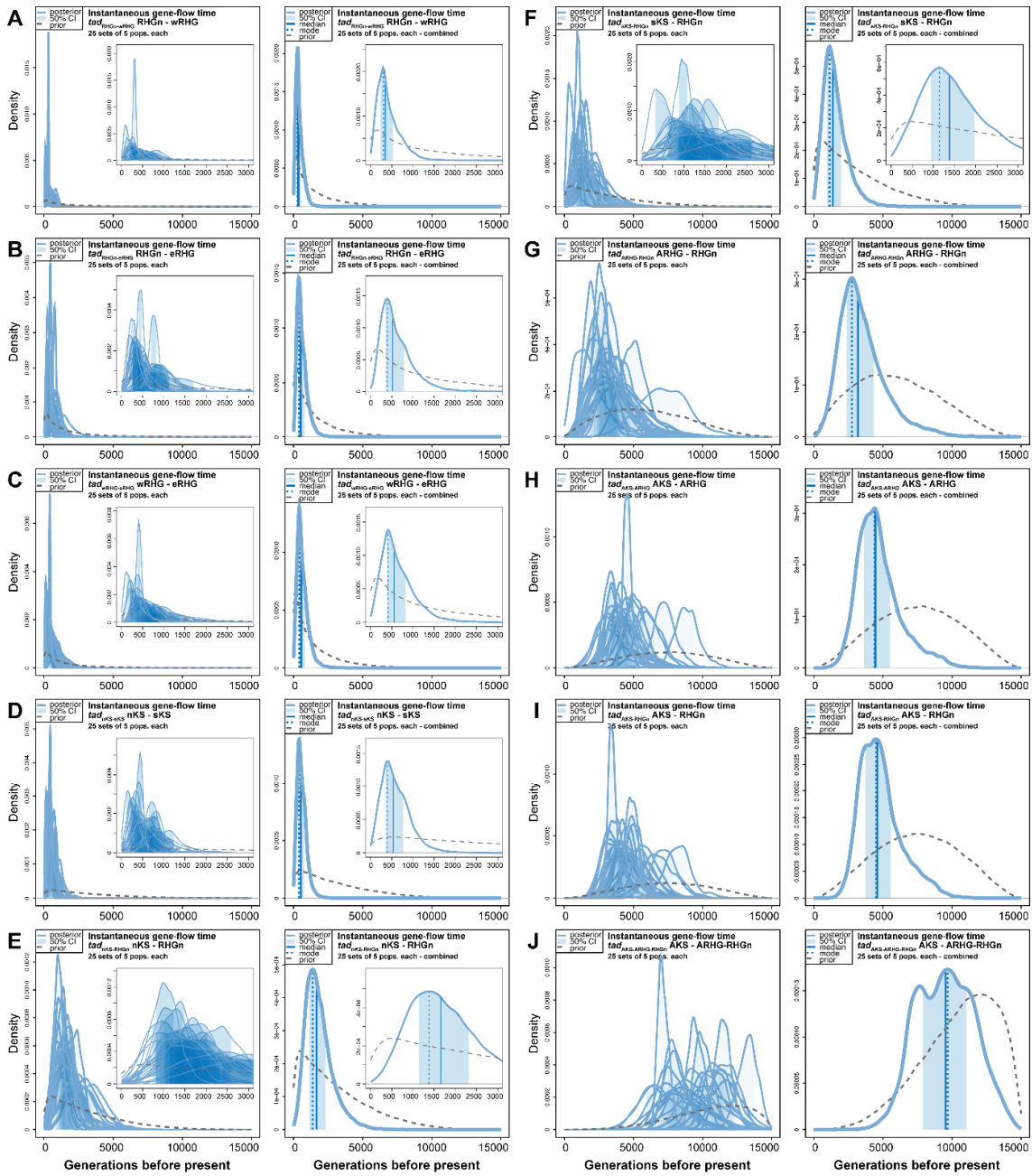
832
833
834
835
836
837
838
839
840
841
842

FigureF8x: ABC posterior parameter distribution of gene-flow times.

Neural Network Approximate Bayesian Computation (Blum and François, 2010; Csilléry, François and Blum, 2012), posterior parameter joint estimations of gene-flow instantaneous times *tad* (in generations before present) for 25 sets of five Central and Southern African populations for which the winning scenario identified by RF-ABC was Scenario i1-1b (**FigureF5x, SupplementaryTableST3x**). Methodological details are provided in **FigureF7x** caption and detailed in **Material and Methods**.

843
844
845
846

Figure F8x



847 2.D.4. Gene-flow intensities.

848 We obtain overall unsatisfactory posterior estimates of the 20 separate unidirectional gene-flow rates
849 for the 10 separate instantaneous gene-flow events, for a majority of the 25 sets of population
850 combinations (**SupplementaryFigureSF7xPart2, SupplementaryTableST4x**), with relatively large
851 90%CI and reduced posterior distributions' departure from priors. Considering only the several
852 satisfactory posterior estimates for each gene-flow parameter separately, we notice their very large
853 variability across the sets of population combinations for each parameter, making it unreasonable to
854 combine posterior distributions to obtain a synthetic value (**TableT4x**). The strong parameter-
855 estimation limitation encountered here is extensively discussed in the **Discussion** section below.

856

857 2.D.5. A synthesis of the demographic history of Central and Southern African populations.

858 We propose, in **FigureF9x**, a schematic synthesis of all the above inference results for the demographic
859 history of 25 out of 54 combinations of five Central and Southern African populations, for which all
860 our RF-ABC scenario-choice procedures provided systematically consistent results.

861

862

863

864

865

866

867 **FigureF9x: Schematic inferred demographic and migration history of Central and Southern African populations.**

868 Schematic representation of the winning Scenario i1-1b (**FigureF5x**), and Neural Network ABC posterior parameter mode
869 estimates summarizing results obtained separately for 25 sets of five Central and Southern African populations
870 (**SupplementaryTableST3x**), represented by the gray lines in between population names. For the time of each divergence and
871 gene-flow event, mode point estimates are provided in generations before present (gbp) in bold, and 90% credibility intervals
872 are provided between parentheses (**TableT4x**). We provide two estimates of the divergence times estimates in years before
873 present (ybp), one (upper) corresponding to 30 years per generation and the other (lower) to 20 years per generation. Mode
874 point estimates of effective population sizes N_e are provided in numbers of diploid effective individuals and width of lineages
875 are proportional to the estimated N_e (**TableT4x, SupplementaryFigureSF7xPart1**). Note that NN-ABC posterior
876 distributions for the effective population sizes of the ancestral Khoe-San lineage (AKS), for the ancestral RHG lineage
877 (ARHG), and for the lineage ancestral to RHGn and ARHG, were all three poorly distinguished from their respective prior
878 distributions; they were therefore considered them to be un-estimated (**TableT4x, SupplementaryFigureSF7xPart1**). Note
879 that posterior distributions for instantaneous asymmetric gene-flow rates were overall poorly departing from their priors, or
880 highly variable when substantially departing from their priors, among the 25 sets of population combinations
881 (**SupplementaryFigureSF7xPart2, SupplementaryTableST4x**). Instantaneous asymmetric gene-flow intensities are
882 therefore considered to be un-estimated. All posterior distributions are shown graphically in **FigureF7x, FigureF8x,**
883 **SupplementaryFiguresSF7xPart1-2,** and detailed in **TableT4x** and **SupplementaryTablesST4x**.

884

885

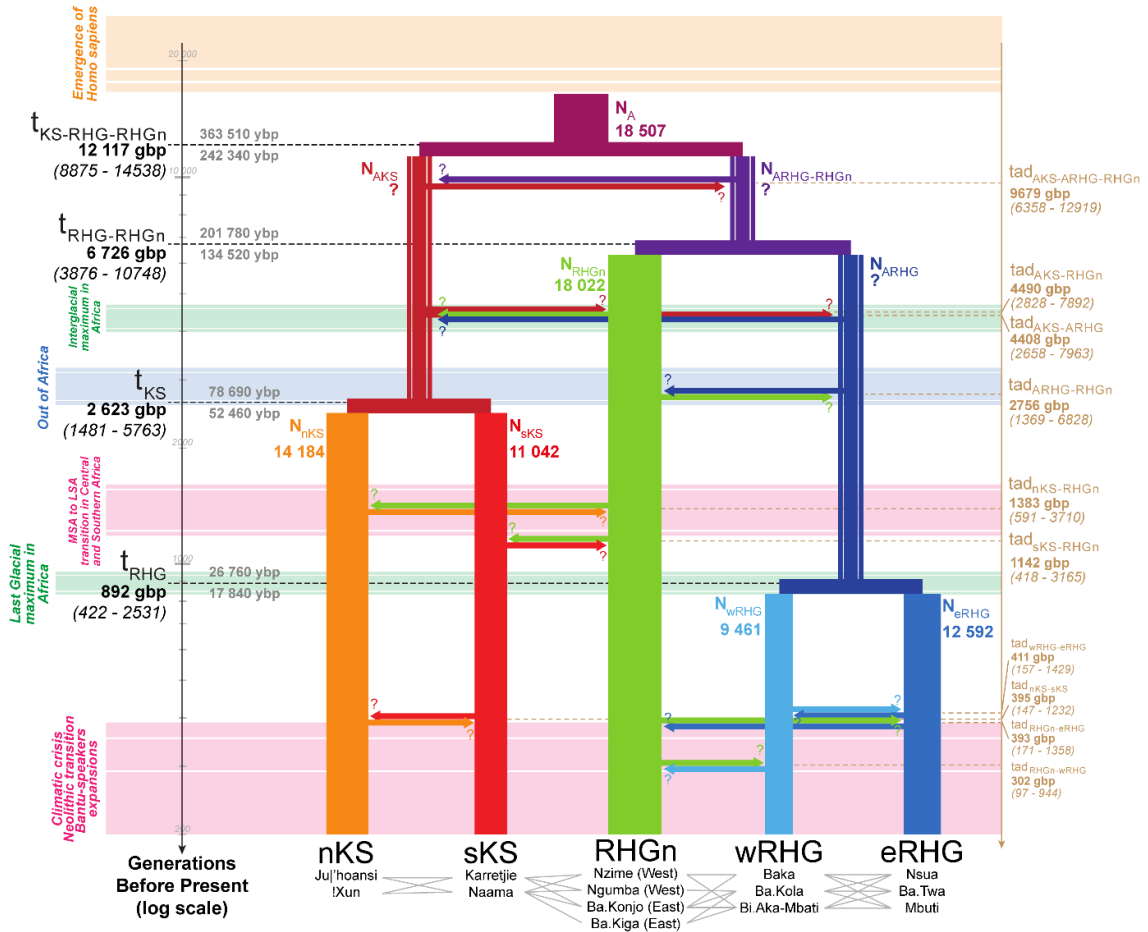
886

887

888 **FigureF9x**

889

890



Discussion

892 We aimed at inferring jointly the tree-topology, instantaneous or recurring gene-flow processes,
893 and asymmetric gene-flow intensities across pairs of recent and ancient lineages that most likely
894 produced genome-wide genetic patterns in 54 different combinations of five Central and Southern
895 African populations. Therefore, we explicitly considered the large genetic differentiation observed
896 across African populations at regional and local scales, and the possible confounding effects of complex
897 gene-flow processes on tree-topologies' predictions among 48 complex scenarios in formal statistical
898 competition.

899 First, our results altogether show that, in fact, considering different sets of populations at a regional
900 and local scale in Africa highlights different aspects of the demographic histories of African
901 populations, sometimes substantially divergent. This demonstrates in practice that apparently discrepant
902 results obtained across previous studies may not necessarily be absolutely reconciled, but rather
903 interpreted as separate illustrations of the large diversity of demographic histories experienced by our
904 species in Africa since its emergence in a remote past. This further strongly advocates for avoiding
905 artificially merging individual samples in un-discriminated *a priori* categories based on geography,
906 subsistence-strategy, and/or linguistics criteria. Instead, we recommend developing novel approaches,
907 probably at a large computational cost as experienced here, to explicitly take into account African
908 genetic diversity at the continental, regional, and/or local scales, and *a posteriori* interpreting the
909 diversity of results obtained separately.

910

Short periods of gene-flow rather than recurring migrations in the ancient history of Africa

912 As pointed out here and in numerous previous theoretical or empirical studies (e.g. (Verdu and
913 Rosenberg, 2011; Gravel, 2012; Lachance *et al.*, 2012; Harris and Nielsen, 2016; Lorente-Galdos *et al.*,
914 2019; Fan *et al.*, 2023; Ragsdale *et al.*, 2023)), whether ancient gene-flow processes occurred
915 recurrently or more instantaneously across lineages has fundamental consequences on the biologic and
916 cultural evolution of our species. Indeed, recurrent gene-flow processes allow for continuous allelic
917 exchanges across populations -thus never reproductively isolated throughout entire periods of
918 evolution-, which may strongly influence the relative influence of drift and selection across populations
919 and the notion of ancient tree-like evolution in human populations in Africa. Conversely, we found that
920 instantaneous gene-flow processes systematically vastly outperformed recurring migration processes,
921 whichever the tree-topology for ancient and recent Central and Southern African lineages divergences,
922 and for all combinations of five observed populations. Therefore, our results unambiguously favor an
923 evolutionary history of African lineages ancestral to a variety of Central and Southern African
924 populations where *Homo sapiens* populations experienced long periods of isolation and drift, followed
925 by short periods of possibly asymmetric gene-flow, which may in turn have induced some reticulations
926 among lineages (Mazet *et al.*, 2016; Ragsdale *et al.*, 2023). Notably, this scenario further allows for
927 differential selection and adaptation processes across lineages (e.g. (Lachance *et al.*, 2012; Schlebusch
928 *et al.*, 2012; Perry *et al.*, 2014)), and for ancient admixture-related selection processes within Africa,
929 analogous to previous findings of such processes having occurred more recently throughout the
930 continent (Breton *et al.*, 2014; Patin *et al.*, 2017; Hamid *et al.*, 2021; Cuadros-Espinoza *et al.*, 2022)).

931 Nevertheless, it is certain that *Homo sapiens* evolutionary history of ancient and recent gene-flow
932 in Africa has involved more complex processes among pairs of lineages than the ones here identified.
933 Previous studies unquestionably considered arguably more realistic models to explain their data, with
934 either a conjunction of recurring symmetric migrations among pairs of lineages with possible
935 instantaneous introgressions among certain pre-specified recent and/or ancient lineages (Lorente-
936 Galdos *et al.*, 2019; Ragsdale *et al.*, 2023), or only pre-specified instantaneous introgressions without

937 recurring gene-flows (Lachance *et al.*, 2012; Fan *et al.*, 2023), including specific events involving the
938 contributions from *Homo sapiens* or non-*Homo sapiens* unsampled lineages. In this context, our results
939 highlight that complex instantaneous introgression scenarios should probably be preferred as a starting-
940 point from where one could then increase complexity, for instance by considering more than a single
941 instantaneous gene-flow event among pairs of ancient lineages.

942 Moreover, our results advocate for the possibility of at least some such introgressions to be
943 relatively intense, as scenarios considering low or moderate gene-flows, whether instantaneous or
944 recurring, are systematically poorly mimetic of observed genomic patterns. Echoing this result, a
945 previous study identified such high levels of instantaneous introgressions across African lineages (Fan
946 *et al.*, 2023), albeit the authors considered very different maximum-likelihood approaches based on
947 different statistics and scenario specifications, as well as different population samples, sometimes
948 artificially merged in larger categories at a continental scale.

949 Nevertheless, our approach failed to provide satisfactory posterior estimates of gene-flow
950 intensities, a serious limitation which we discuss extensively in the following methodological limits and
951 perspectives section below. Importantly, previous studies investigating conjunctions of recurring and
952 more instantaneous migrations overall estimated that symmetric recurring gene-flows were weak,
953 whether recent or ancient (e.g. (Ragsdale *et al.*, 2023)). Although our results are difficult to compare
954 with previous studies due to strong differences in model specifications, statistics used and
955 methodological approaches as well as population samples, this latter previous result may be reflected
956 in ours that disfavor recurring migration models as being poorly explicative of our data.

957

958 **Inferred demographic history in Central and Southern Africa**

959 Extensive previous work agrees that divergences among lineages ancestral respectively to extant Khoe-
960 San, Rainforest Hunter-Gatherers, and Rainforest Hunter-Gatherer neighbors were among the most
961 ancient in *Homo sapiens* evolution (e.g. reviewed in (Schlebusch and Jakobsson, 2018; Pfennig *et al.*,
962 2023)). However, the relative order of their divergences as well as their timing has been a matter of
963 extensive debate (Schlebusch *et al.*, 2020; Lipson *et al.*, 2022; Fan *et al.*, 2023). This is due in part to
964 split-time inferences conducted on tree-like topologies considering gene-flow or no gene-flow, using
965 ancient DNA data or not, and to differing methods, statistics, and population samples.

966

967 Most ancient divergences in Africa

968 Here, we formally tested, with Random-Forest ABC scenario-choice, which ancient tree-topology best
969 explained extant genomic patterns, whichever the gene-flow processes and intensities across pairs of
970 recent and ancient lineages, for 54 different combinations of five different Khoe-San, Rainforest
971 Hunter-Gatherer, and Rainforest Hunter-Gatherer neighboring populations. It is the first time to our
972 knowledge that such formal comparison of numerous competing-scenarios is conducted systematically
973 for a variety of sets of population samples at a regional and local scale, rather than comparing
974 maximum-likelihood values from vastly differing models, each obtained separately using differing
975 population samples (Lipson *et al.*, 2022; Fan *et al.*, 2023; Ragsdale *et al.*, 2023). Importantly, note that
976 (Lorente-Galdos *et al.*, 2019) also explored highly complex demographic scenarios in Africa with ABC,
977 focusing on varying unidirectional introgression processes among certain lineages for a fixed tree-shape
978 among, in particular, three extant African populations at the continental scale.

979 Our results show that whichever the gene-flow processes considered, the lineage ancestral to
980 extant KS populations diverged first from the lineage ancestral to RHG and RHGn, for almost all
981 combinations of sampled populations. Furthermore, we estimated with Neural Network ABC posterior-
982 parameter inferences that this original divergence likely occurred in a remote past ~300,000 years ago,
983 followed by the divergence between ancestral RHG and RHGn lineages ~165,000 years ago,
984 considering between 20 and 30 years per generations.

985 Both estimates fell within the upper bound of previously obtained results considering a variety of
986 similar population samples, genetic data, inference methods, and statistics, albeit scenarios may not
987 have always been specified in analogous ways (Patin *et al.*, 2009, 2014; Verdu *et al.*, 2009, 2013;
988 Schlebusch *et al.*, 2012, 2020; Lipson *et al.*, 2022; Fan *et al.*, 2023). Note however, that considering
989 different combinations of populations at the regional and local scales allowed us to identify a small
990 minority (4/54) of such combinations for which an alternative scenario (Scenario i1-3b), where RHGn
991 lineages diverged first from a lineage ancestral to all RHG and KS extant populations, somewhat
992 analogous to those proposed in (Fan *et al.*, 2023), would better explain the data consistently across our
993 analyses, again whichever the gene-flow processes and range of intensities considered. This latter result
994 may explain apparently discrepant results across some previous studies, which would then be due to
995 differing population samples used for demographic inferences and/or to artificial merging of
996 differentiated populations into larger groups (Lorente-Galdos *et al.*, 2019; Lipson *et al.*, 2022; Fan *et al.*,
997 2023; Ragsdale *et al.*, 2023), thus advocating for further accounting for the vast genetic diversity
998 among African populations even at a local scale in future studies.

999 Note that, in (Ragsdale *et al.*, 2023), the fundamental divergence between their Nama Khoe-San
1000 sample and other Sub-Saharan populations was dated to ~110,000-135,000 years ago under their two
1001 best-fitting models, thus strongly discrepant with our findings. However, these important results are
1002 very difficult to compare with ours, since the authors did not consider any Central African Rainforest
1003 Hunter-Gatherers nor their neighbors in their inferences, and since they investigated highly complex
1004 models specified completely differently from those here envisioned. In particular, Ragsdale and
1005 colleagues included in their models possible very ancient genetic structures, long before *Homo sapiens*
1006 emergence, a feature that is unspecified in our scenarios which considered simply a single ancestral
1007 population in which all extant lineages ultimately coalesce. Nevertheless, note that substructure and
1008 reticulation within the ancestral population is not *per se* incompatible with our scenarios. In fact, it may
1009 be compatible with our posterior estimates of a large effective population ancestral to all extant
1010 populations here investigated, the largest among all inferred ancient and recent Central and Southern
1011 African effective population sizes. Therefore, it will be reasonable in future work to complexify the
1012 scenarios here proposed to evaluate whether very ancient substructures and reticulations within our
1013 ancestral population, prior to the original divergence between Southern and Central African
1014 populations, may improve the fit to the observed genomic data, as proposed by (Ragsdale *et al.*, 2023).

1015 More recent divergences in Central and Southern Africa

1016 More recently during the evolutionary history of Sub-Saharan Africa, we found that the divergence
1017 among Northern and Southern Khoe-San populations largely pre-dated the divergence of Western and
1018 Eastern Congo Basin Rainforest Hunter-Gatherers, a question rarely addressed to our knowledge. First,
1019 we found that KS divergence dated sometime between 50,000 and 80,000 years ago, thus substantially
1020 more recently than estimates previously proposed (Schlebusch and Jakobsson, 2018). Beyond vast
1021 differences in models, methods, and statistics used to provide either inferences, note that our results for
1022 this divergence time were substantially variable across pairs of sampled populations used in each
1023 analysis, some specific sets of populations providing posterior estimates consistent with previous result.
1024 It thus further highlights that complex inferences in Africa imperatively need to explicitly consider
1025 population variation at a local scale.

1026
1027 Interestingly, our synthesized estimates for the Northern and Southern Khoe-San population
1028 divergence were relatively synchronic to the genetic onset of the Out-of-Africa (e.g. (Schlebusch and
1029 Jakobsson, 2018)). Population genetics inferences only provide possible mechanisms to explain
1030 observed genetic patterns and are, in essence, not addressing the possible causes underlying the inferred
1031 mechanisms. In this context, we may hypothesize that the global climatic shifts inducing massive
1032 ecological changes that have occurred in Africa at that time (e.g. (Beyer *et al.*, 2021)), sometimes

1033 proposed to have triggered ancient *Homo sapiens* movements Out-of-Africa, may also have triggered,
1034 independently, the genetic isolation among ancestral Khoe-San populations. Nevertheless, where the
1035 ancestors of extant Khoe-San populations lived at that time remains unknown and is nevertheless crucial
1036 to further elaborate possible scenarios for the causes of the genetic divergence here inferred.

1037 Long after this divergence, we found that Rainforest Hunter-Gatherer populations across the
1038 Congo Basin diverged roughly between 17,000 and 27,000 years ago, relatively consistently across
1039 pairs of sampled populations used for inferences; estimates highly consistent with previous studies
1040 (Patin *et al.*, 2009, 2014; Lopez *et al.*, 2018), despite major differences in gene-flow specifications
1041 across RHG groups between studies. Interestingly, this divergence time is relatively synchronic with
1042 absolute estimates for the Last-Glacial Maximum in Sub-Saharan Africa (e.g. (Bartlein *et al.*, 2011)).
1043 The fragmentation of the rainforest massif during this period in the Congo Basin may have induced
1044 isolation between Eastern and Western RHG extant populations, as plausibly previously proposed
1045 (Patin *et al.*, 2009). However, similarly as above for the Northern and Southern Khoe-San populations
1046 divergence, where the ancestors of extant Eastern and Western Rainforest Hunter-Gatherers lived
1047 remains unknown, which prevents us from formally testing this hypothesis (Perry and Verdu, 2017).

1048 Altogether, these results show that Central African Rainforest Hunter-Gatherer and Southern
1049 African Khoe-San populations have had, respectively, extensive time for selection processes, including
1050 adaptive introgression processes, to have influenced independently both groups of populations as well
1051 as populations within each group separately (e.g. (Schlebusch *et al.*, 2012; Breton *et al.*, 2014; Perry *et al.*,
1052 2014; Patin *et al.*, 2017)).

1053 Ancient and recent instantaneous gene-flow times in Africa

1054 We obtained reasonably well estimated instantaneous asymmetric gene-flow times among almost all
1055 pairs of ancient and recent lineages, albeit the variation of estimates across sets of Central and Southern
1056 African sampled populations increased substantially with most ancient times. In this context, we deem
1057 it hard to confidently try to interpret the most ancient event of instantaneous gene-flow between the two
1058 most ancestral lineages in our tree-topology. However, other instantaneous gene-flow time estimates
1059 throughout the topology were more consistently estimated overall, and showed relative synchronicity
1060 in some cases, which has never been reported before to our knowledge, even if they were specified
1061 independently in our models and drawn from large distribution *a priori*.

1062 Interestingly, we found strong indications for almost synchronic events of introgressions having
1063 occurred during the Last Interglacial Maximum in Africa (Mazet *et al.*, 2016), between ~90,000 and
1064 ~135,000 years ago (when considering 20 or 30 years per generation). They involved gene-flow
1065 between lineages ancestral to Khoe-San populations and ancestors of Rainforest Hunter-Gatherer
1066 neighbors on the one hand and, on the other hand, between lineages ancestral to Khoe-San populations
1067 and the lineage ancestral to all Rainforest Hunter-Gatherers. An increase in material-based culture
1068 diversification and innovation, possibly linked to climatic and environmental changes locally, has
1069 previously been observed during this period of the Middle Stone Age in diverse regions of continental
1070 Africa; prompting a long-standing debate as to its causes if human populations were subdivided and
1071 isolated biologically and culturally at the time (Ziegler *et al.*, 2013; Scerri *et al.*, 2018; Gosling, Scerri
1072 and Kaboth-Bahr, 2022; Thomas *et al.*, 2022).

1073 In this context, our results instead may suggest that population movements at that time among
1074 previously isolated populations may itself have triggered the observed increased cultural diversification
1075 locally, even if obvious signs of pan-African cultural spread at the time are difficult to assess (Gosling,
1076 Scerri and Kaboth-Bahr, 2022). In turn, it would interestingly echo the known effect of within-
1077 population genetic diversity increase induced by genetic admixture between previously isolated
1078 populations (e.g. (Long, 1991; Verdu and Rosenberg, 2011; Gravel, 2012; Laurent *et al.*, 2023)).
1079

1080 Then, note that we estimated that the instantaneous gene-flow event between the ancestral
1081 Rainforest Hunter-Gatherers lineage and that of their extant neighbors seemingly occurred
1082 synchronically to the genetic Out-of-Africa ((Beyer *et al.*, 2021); see above). This would imply that
1083 possible climatic and ecological shifts at that time may not have only induced population divergences
1084 and displacement, but may also have triggered population gene-flow.

1085 Relatively more recently, around 30,000 years ago, we found two loosely synchronic gene-flow
1086 events between ancestors to extant Central African Rainforest Hunter-Gatherer neighbors' lineages and,
1087 separately, Northern and Southern Khoe-San lineages. This corresponds to the end of the Interglacial
1088 Maximum and a period of major cultural changes and innovations during the complex transition from
1089 Middle Stone Age to Late Stone Age in Central and Southern Africa (Cornelissen, 2002; Ziegler *et al.*,
1090 2013; Mesfin, Oslisly and Forestier, 2021; Bader *et al.*, 2022; Thomas *et al.*, 2022). Nevertheless,
1091 connecting the two lines of genetic and archaeological evidence to conclude for increased population
1092 movements at the time and their possible causes should be considered with caution. Indeed, in addition
1093 to genetic-dating credibility-intervals being inherently much larger than archaeological dating, this
1094 period remains highly debated in paleoanthropology mainly due to the scarcity and complexity of the
1095 material-based culture records, and that of climatic and ecological changes locally, across vast regions
1096 going from the Congo Basin to the Cape of Good Hope (Cornelissen, 2002; Ziegler *et al.*, 2013; Mesfin,
1097 Oslisly and Forestier, 2021; Bader *et al.*, 2022; Thomas *et al.*, 2022).

1098 Finally, we found strong signals for multiple instantaneous gene-flow events having occurred
1099 between almost all five recent Central and Southern African lineages between 6000 and 12,000 years
1100 ago, during the onset of the Holocene in that region, shortly before or during the beginning of the last
1101 Post Glacial Maximum climatic crisis in Western Central Africa (Lézine *et al.*, 2019), the emergence
1102 and spread of agricultural techniques (Phillipson, 2005), and the demic expansion of now-Bantu-
1103 speaking populations from West Central Africa into the rest of Central and Southern Africa (Bostoen
1104 *et al.*, 2015; Patin *et al.*, 2017; Fortes-Lima *et al.*, 2024). These results are consistent with previous
1105 investigations that demonstrated the determining influence of Rainforest Hunter-Gatherer neighboring
1106 populations' migrations through the Congo Basin in shaping complex socio-culturally determined
1107 admixture patterns (Patin *et al.*, 2009, 2014; Verdu *et al.*, 2009, 2013), including admixture-related
1108 natural selection processes (Perry *et al.*, 2014; Patin *et al.*, 2017; Lopez *et al.*, 2018, 2019). As our
1109 estimates for introgression events are in the upper bound of previous estimates for the onset of the so-
1110 called "Bantu expansion" throughout Central and Southern Africa, we may hypothesize here that major
1111 climatic and ecological changes that have occurred at that time may have triggered increased population
1112 mobility and gene-flow events between previously isolated populations, rather than consider that the
1113 Bantu-expansions themselves were the cause for all the gene-flow events here identified.

1114 Finally, we did not find signals of more recent introgression events from Bantu-speaking
1115 agriculturalists populations into Northern or Southern Khoe-San populations, in particular among the
1116 !Xun, albeit such events have been identified in several previous studies (see (Schlebusch and
1117 Jakobsson, 2018)). This is likely due to the fact that we considered only a limited number of individual
1118 samples from each population, and therefore may lack power to detect these very recent events with our
1119 data and approach.

1120

1121 **Conceptual, methodological, and empirical limitations and perspectives for inferring ancient** 1122 **histories from observed genomic data**

1123 All previous attempts at reconstructing the human evolutionary histories which led to genomic patterns
1124 observed today across African populations have faced major conceptual, methodological, and empirical
1125 challenges. Conceptually, large amounts of more-or-less nested scenarios can be envisioned *a priori*
1126 to explain extant genetic diversity, based on previous results from paleo-anthropology, population
1127 genetics, and paleogenomics. These scenarios may range from tree-like models without gene-flow

1128 events among ancient or recent lineages to complex networks of weakly differentiated populations
1129 exchanging migrants over large periods of time, with or without the contribution of ancient *Homo* or
1130 non-*Homo* now extinct or unsampled lineages. Systematically exploring all possible models is often
1131 methodologically out of reach due to differing fundamental scenario-specifications or, when scenarios
1132 are specified and parameterized in analogous ways, due to scenarios' nestedness and un-identifiability
1133 in certain parts of their parameter spaces (e.g. (Robert, Mengersen and Chen, 2010)). Empirically,
1134 formal scenario comparisons are first hampered by necessarily limited amounts of genomic data
1135 representative, at continental, regional, and local scales, of the known diversity and differentiation of
1136 human populations in Africa; in addition to yet limited amounts of ancient DNA data throughout the
1137 continent. Finally, empirical limitations also emerge from the use of different statistics to explore
1138 genomic diversity patterns, which possibly each capture different facets of human evolutionary
1139 histories, thus providing discrepant results and interpretations only in appearances. While machine-
1140 learning ABC procedures provide significant advantages over other maximum-likelihood approaches,
1141 in particular concerning the formal exploration of competing scenarios' fit to observed data across
1142 numerous highly complex sometimes nested scenarios and using numerous summary-statistics (Blum
1143 and François, 2010; Robert, Mengersen and Chen, 2010; Pudlo *et al.*, 2016; Estoup *et al.*, 2018), the
1144 above challenges are also largely faced in this study.

1145 While most divergence and gene-flow times, as well as effective sizes, were inferred satisfactorily
1146 in our results, the lack of satisfactory posterior estimates for all 20 gene-flow rates parameters consistent
1147 among the 25 sets of population combinations likely stems from different limitations. First, considering
1148 only five individual genomes per population is inherently limiting when trying to estimate gene-flow
1149 rates from inter-population summary-statistics in ABC (Fortes-Lima *et al.*, 2021). In future studies,
1150 increasing sample sizes and adding statistics based on the inter and intra-individual distribution of
1151 admixture fractions (Verdu and Rosenberg, 2011; Gravel, 2012; Ragsdale and Gravel, 2019; Fortes-
1152 Lima *et al.*, 2021), will likely improve the posterior estimation of these parameters using machine-
1153 learning ABC approaches (Lorente-Galdos *et al.*, 2019). Note, however, that these statistics require
1154 non-trivial computation time in ABC frameworks comprising hundreds of thousands of simulations
1155 (Boitard *et al.*, 2016; Jay, Boitard and Austerlitz, 2019). Furthermore, in all cases, considering multiple
1156 sets of different and substantially genetically diverse populations as the ones here considered from sub-
1157 Saharan Africa may inevitably leads to large variation across gene-flow rates' posterior estimates
1158 depending on the specific population-combinations.

1159 Second, for simplicity, we considered a single instantaneous time for the two separate
1160 unidirectional gene-flow events, for each gene-flow event between pairs of lineages separately.
1161 Therefore, while our RF-ABC scenario-choice results strongly support instantaneous asymmetric gene-
1162 flow events rather than recurring asymmetric gene-flows to best explain extant genomic patterns, it is
1163 plausible that scenarios where each unidirectional gene-flow event may occur at a different time
1164 (Lorente-Galdos *et al.*, 2019), will more realistically explain the observed data. If this is the case, it is
1165 possible that the choice of a unique time for two separate gene-flow events rendered the corresponding
1166 gene-flow rates harder to identify with our joint NN-ABC posterior parameter estimation procedure.
1167 Furthermore, we explored two extreme gene-flow processes, establishing an open competition between
1168 only instantaneous gene-flow and only recurring gene-flow processes. While we demonstrated that only
1169 instantaneous gene-flow processes vastly outperformed only recurring gene-flow processes, whichever
1170 the range of intensities for each event, we only established here a starting point for the future
1171 complexification of evolutionary scenarios to be tested. For instance, it will be of natural interest to
1172 consider, next, more than a single "pulse" of gene-flow between any two ancient or recent lineage.

1173 Third, the neural-networks used here may have been unable to satisfactorily identify the 20 gene-
1174 flow parameters among the 43 jointly-estimated parameters, due to their lack of complexity when
1175 considering a unique layer of hidden neurons. For instance, and as a first step, considering Multilayer-

1176 Perceptrons in the future, *i.e.* more complex neural-networks with additional layers of hidden neurons,
1177 may help posterior estimations of these parameters (e.g. (Wang, Czerminski and Jamieson, 2021)), at
1178 the known cost of non-trivial parameterization of the neural-networks themselves (e.g. (Leung *et al.*,
1179 2003; Jay, Boitard and Austerlitz, 2019; Huang *et al.*, 2024)).

1180 Altogether, joint posterior estimation of numerous gene-flow rates parameters and the timing of
1181 their occurrence under highly complex demographic scenarios remains one of the most challenging
1182 tasks in population genetics. It will unquestionably benefit from future analytical theoretical
1183 developments (Mooney *et al.*, 2023; Agranat-Tamir, Mooney and Rosenberg, 2024), and the
1184 improvement of machine-learning-based inference procedures (e.g. (Murtagh, 1991; Chen *et al.*, 2020;
1185 Yelmen and Jay, 2023; Huang *et al.*, 2024)).

1186

1187 **Ancient admixture with *Homo sapiens* or non-*Homo sapiens* unsampled lineages?**

1188 We did not explore possible contributions from unsampled lineages, whether from non-*Homo sapiens*
1189 or from ancient “ghost” human populations, and therefore cannot formally evaluate the likeliness of the
1190 occurrence of such events to explain observed data. In all cases, our results demonstrate that explicitly
1191 considering ancient admixture from unsampled populations is not a necessity to explain satisfactorily
1192 large parts of the observed genomic diversity of extant Central and Southern African populations,
1193 consistently with a previous study (Ragsdale *et al.*, 2023), and conversely to others (Lipson *et al.*, 2022;
1194 Fan *et al.*, 2023; Pfennig *et al.*, 2023); at least when considering jointly the 337 relatively classical
1195 population genetics summary-statistics used here for demographic inferences. As discussed above, our
1196 results formally comparing competing-scenarios rather than comparing posterior likelihoods of highly
1197 complex yet vastly differing models, provide a clear and reasonable starting point for future
1198 complexification of scenarios comprising possible contributions from ancient or ghost unsampled
1199 populations, which will unquestionably benefit from the explicit use of additional novel summary-
1200 statistics ((Ragsdale and Gravel, 2019; Fan *et al.*, 2023; Ragsdale *et al.*, 2023); see also above).

1201 In any case, the complexification of scenario-specifications to account for possible past “archaic”
1202 or “ancient” introgressions will not fundamentally solve the issue of the current lack of reliable ancient
1203 genomic data older than a few hundreds or thousands of years from Sub-Saharan Africa (Skoglund *et al.*,
1204 2017; Vicente and Schlebusch, 2020; Pfennig *et al.*, 2023). Indeed, analogously to archaic
1205 admixture signals that were unambiguously identified outside Africa only when ancient DNA data were
1206 made available for Neanderthals and Denisovans (e.g. (Meyer *et al.*, 2012; Prüfer *et al.*, 2014)), we
1207 imperatively need to overcome this lack of empirical ancient DNA data in Africa to formally test
1208 whether, or not, ancient human or non-human now extinct lineages have contributed to shaping extant
1209 African diversity.

1210

1211

1212

1213

Material and Methods

1214

Population samples

1215

Central and Southern Africa dataset

1216

We investigated high-coverage whole genomes newly generated for 74 individual samples (73 after relatedness filtering, see below), from 14 Central and Southern African populations (**FigureF1x, TableT1x**). Based on extensive ethno-anthropological data collected from semi-directed interviews with donors and their respective communities, we grouped *a posteriori* the 73 individual samples from 14 populations in three larger categories.

1221

The Baka, Ba.Kola, Bi.Aka_Mbati, Ba.Twa, and Nsua individuals were categorized as Rainforest Hunter-Gatherers (RHG), based on several criteria including self-identification and relationships with other-than-self, socio-economic practices and ecology, mobility behavior, and musical practices (Verdu *et al.*, 2009; Hewlett, 2017).

1225

The Nzime, Ngumba, Ba.Kiga, and Ba.Konjo individuals were categorized as Rainforest Hunter-Gatherer neighbors (RHGn), based on the same set of categorization-criteria as RHG. Indeed, RHG and RHGn populations in Central Africa are known to identify themselves separately and to live in different ways in the Central African rainforest, whilst often sharing languages locally as well as complex socio-economic relationships and interactions that may include intermarriages (Verdu *et al.*, 2013). Historically, RHG populations have been designated as “Pygmies” by European colonists, an exogenous term from ancient Greek that is sometimes used derogatorily by the RHGn.

1232

Finally, the Nama, Ju!’hoansi, Karretjie People, !Xun, and Khutse_San individuals were categorized as Khoe-San populations (KS) based primarily on self-identification and relationships with other-than-self locally, lifestyle, and languages (Schlebusch, 2010). Note that research was approved by the South African San Council. Individuals from these populations were whole-genome sequenced anew with improved methods (see below), compared to previous work (Schlebusch *et al.*, 2020).

1237

Note that the three groups can be further subdivided based on geography into Western and Eastern RHG and RHGn, and into Northern, Central, and Southern KS (**FigureF1x, TableT1x**).

1239

Comparative dataset

1240

The 74 genomes were analyzed together with 105 (104 after relatedness filtering, see below) high-coverage genomes from 32 worldwide populations (**FigureF1x, TableT1x**). 64 samples from the Simon's Genome Diversity Project (SGDP) (Mallick *et al.*, 2016); 9 HGDP samples (Meyer *et al.*, 2012); 7 samples from the 1000 Genomes project (KGP) (The 1000 Genomes Project Consortium *et al.*, 2015); one Karitiana sample (Rasmussen *et al.*, 2014); 24 samples from the South African Human Genome Project (SAHGP) (Choudhury *et al.*, 2017). Selection criteria were: Illumina paired-end reads, high-coverage (>30X), access to raw data and no (known) population substructure.

1247

Note that, although samples from these datasets have not been collected by us for the purpose of this study, based on previous knowledge and publications, the Biaka and Mbuti HGDP samples can be categorized as RHG and the #Khomani and Ju!’hoansi_comp samples as KS (**FigureF1x, TableT1x**).

1248

1249

1250 **Data generation**

1251 **DNA extraction**

1252 DNA was extracted from saliva with the DNA Genotek OG-250 kit for the Baka, Nzime, Bi.Aka Mbatu,
1253 Nsua, and Ba.Konjo samples. DNA was extracted from buffy coats with DNeasy Blood&Tissue spin-
1254 column Qiagen™ kits for the Ba.Kola, Ngumba and Ba.Kiga samples. Both extraction methods were
1255 used for the Ba.Twa samples. For the Nama, Ju|'hoansi, Karretjie People, !Xun, and Khutse_San
1256 samples, DNA was extracted from EDTA-blood using the salting-out method (Miller, Dykes and
1257 Polesky, 1988).

1258 **Library preparation and sequencing**

1259 Library preparation and sequencing was performed by the SciLifeLab SNP&SEQ Technology platform
1260 in Uppsala, Sweden. Libraries were prepared with TruSeq DNA preparation kits, and paired-end
1261 sequencing (150 bp read length), was performed on Illumina HiSeqX machines with v2 sequencing
1262 chemistry, to a coverage of 30X or more. For the Southern African samples, we included paired-end
1263 data (100 bp read length) sequenced on Illumina HiSeq2000 machines obtained previously for the same
1264 libraries (Schlebusch *et al.*, 2020).

1265 **Sequencing data quality-control processing and relatedness filtering**

1266 The processing pipeline used in the present study is adapted from the Genome Analyses Toolkit
1267 (GATK) “Germline short variant discovery (SNPs + Indels)” Best Practices workflow (McKenna *et al.*,
1268 2010; DePristo *et al.*, 2011; Van Der Auwera *et al.*, 2013). It is described and compared to the original
1269 GATK Best Practices workflow in (Breton *et al.*, 2021). **SupplementaryFigureSF8x** gives an
1270 overview of the pipelines used for the processing; all template codes with accompanying detailed
1271 explanations for all processing steps are provided below and in the corresponding GitHub repository
1272 (<https://github.com/Gwennid/africa-wgs-descriptive>).

1273 Reads were mapped to a decoy version of the human reference genome, GRCh38 (1000 Genomes
1274 Project version), with BWA-MEM (Li and Durbin, 2009) from BWAKIT v0.7.12 using a strategy
1275 appropriate for ALT contigs. The input was either FASTQ or BAM files that were reverted to unmapped
1276 BAM or to FASTQ prior to mapping. Mapped reads were selected with SAMtools version 1.1 and
1277 Picard version 1.126 (<https://broadinstitute.github.io/picard/>).

1278 Duplicate reads were marked (at the lane level) with Picard version 1.126 MarkDuplicates.
1279 Realignment around indels was performed (at the lane level) with GATK version 3.5.0
1280 RealignerTargetCreator and IndelRealigner. We then performed a “triple mask base quality score
1281 recalibration (BQSR)”, where the sample’s variation is used together with dbSNP, as described in
1282 (Schlebusch *et al.*, 2020; Breton *et al.*, 2021). This step includes merging the reads from a given sample
1283 with SAMtools version 1.1; calling variants with GATK version 3.5.0 HaplotypeCaller with “--
1284 genotyping_mode DISCOVERY”; and (in a parallel track) standard BQSR with GATK version 3.5.0
1285 BaseRecalibrator and PrintReads, and dbSNP version 144. Following triple mask BQSR (performed at
1286 the lane level), the recalibrated reads from a given sample are merged with SAMtools version 1.1, sorted
1287 and indexed with Picard version 1.126 SortSam. Duplicate marking and realignment around indels were
1288 performed again, at the sample level.

1289 Variant calling was performed independently for the autosomes and the X chromosomes, due to
1290 the difference in ploidy. SNP and short indels were first called with GATK version 3.7 HaplotypeCaller
1291 in each genome, with the “--emitRefConfidence BP_RESOLUTION” option. Multi-sample GVCFs

1292 were then generated with CombineGVCF, and variants were finally jointly called with GATK
1293 GenotypeGVCFs. Genotypes were emitted for each site (option “allSites”).

1294 Variants were filtered with GATK version 3.7 Variant Quality Score Recalibration (VQSR) using
1295 recommended resources from the GATK bundle: HapMap (version 3.3), 1000 Genomes Project
1296 Illumina Omni 2.5M SNP array, 1000 Genomes Project phase 1 high confidence SNPs, and dbSNP
1297 version 151 for the SNPs; and (Mills *et al.*, 2011) gold standard indels and dbSNP version 151 for the
1298 indels. For the X chromosome, autosomal variants were included for training the VQSR model
1299 (VariantRecalibrator step).

1300 We controlled the dataset for related individuals with KING (Manichaikul *et al.* 2010) “--kinship”
1301 and plink version 1.90b4.9 (Purcell *et al.*, 2007) “--genome --ppc-gap 100”. The dataset was pruned for
1302 linkage disequilibrium (LD) before relatedness estimation with plink, using plink “--indep-pairwise”
1303 with sliding windows of 50 SNPs, shifting by five SNPs, and a r^2 threshold of 0.5. Both methods
1304 identified two pairs of first-degree relatives. We excluded the sample with greatest missingness from
1305 each pair to produce the family unrelated working datasets henceforth used.

1306 The callset was then further refined; two samples with a first degree relative in the dataset were
1307 excluded with GATK version 3.7 SelectVariants; a “FAIL” filter status was set on sites that are
1308 ambiguous in the reference genome (base “N”) or had greater than 10% missingness, identified with
1309 VCFtools version 0.1.13 “--missing-site” (Danecek *et al.*, 2011). Moreover, for the autosomes, variants
1310 heterozygous in all samples were marked as failed (Hardy Weinberg Equilibrium -HWE- filter,
1311 identified with VCFtools version 0.1.13 “--hardy”).

1312 Coverage was computed with QualiMap version 2.2 (Okonechnikov, Conesa and García-Alcalde,
1313 2016) including or not duplicates (“-sd” option for the latter). The average coverage (without duplicates)
1314 across individuals within populations is provided in **SupplementaryTableST2x**.

1315 **Descriptive analyses**

1316 **Variant counts**

1317 We obtained per-sample and aggregated metrics of variant counts (**FigureF2x**), with the tool
1318 *CollectVariantCallingMetrics* of the software package Picard v2.10.3
1319 (<https://broadinstitute.github.io/picard/>), applied to the full 177 worldwide individuals dataset after
1320 VQSR, relatedness, HWE and 10% site missingness filtering. We conducted variant counts procedures
1321 for chromosomes 1 to 22. We used dbSNP 156 as a reference; we downloaded
1322 “GCF_000001405.40.gz” and modified the contig names to match contig names in the VCFs.
1323 Separately, we applied the same variant counts pipeline as above to the Central and Southern African
1324 73 individuals’ original dataset extracted from the full dataset with BCFtools version 1.17
1325 (<https://github.com/samtools/bcftools>), using view with the option “-S”. All scripts are provided in the
1326 corresponding GitHub repository (<https://github.com/Gwennid/africa-wgs-descriptive>).

1327 **Genome-wide Heterozygosity**

1328 We calculated observed and expected heterozygosities, for the autosomes and for the X-chromosome
1329 separately using custom Python, Bash, and R scripts (full pipeline available at GitHub
1330 <https://github.com/Gwennid/africa-wgs-descriptive>), for the 177 worldwide individuals after variant-
1331 counts filtering described above. For autosomes, we considered the main contigs (“chr1” to “chr22”).
1332 For the X chromosome, we excluded the Pseudo-Autosomal Regions with coordinates from GRCh38
1333 (PAR1: 10,000 to 2,781,480, PAR2: 155,701,382 to 156,030,896).

1334 For each population with more than one individual and for each autosome and X-chromosome
1335 separately, we counted the number of variable and non-variable sites and excluded multi-allelic sites,
1336 indels, and sites with missing genotypes in at least one individual in the population for simplicity and
1337 conservativeness. For each autosome and each individual, we then counted separately the numbers of
1338 observed homozygous and heterozygous sites of each configuration compared to the reference (0/0;
1339 0/1; 1/0; 1/1). Based on these counts, for each population and each chromosome separately, we first
1340 calculated observed heterozygosities simply as the average proportion of heterozygous individuals per
1341 variable sites with no-missing data in the population sample. We also computed unbiased expected
1342 multi-locus heterozygosities for all variable loci with no missing genotypes in the population as in
1343 equations (2) and (3) in (Nei, 1978). We averaged this value across all sites with no missing information
1344 in the population sample including both variable and non-variable sites, and finally corrected it for
1345 haploid population sample sizes (**FigureF2x**). All scripts are provided in the corresponding GitHub
1346 repository (<https://github.com/Gwennid/africa-wgs-descriptive>).

1347 **Runs of homozygosity**

1348 We identified runs of homozygosity (ROH) using the *homozyg* tool from PLINK version 1.90b4.9
1349 (Purcell *et al.*, 2007). We selected autosomal biallelic SNPs from the variant-counts pipeline described
1350 above with GATK version 3.7 SelectVariants with options “-selectType SNP -restrictAllelesTo
1351 BIALLELIC -excludeFiltered”. We converted the VCF to TPED and then binary plink fileset with
1352 VCFtools version 0.1.13 “-plink-tped” and plink. For the 177 individuals worldwide, we considered
1353 only ROHs measuring more than 200 Kb (*--homozyg-kb 200*), containing at least 200 SNPs (*--homozyg-
1354 snp 200*), containing at least one variable site per 20 Kb on average (*--homozyg-density 20*), and with
1355 possible gaps of up to 50 Kb (*--homozyg-gap 50*). Default values were considered for all other
1356 parameters: *--homozyg-window-het 1 --homozyg-window-snp 50 --homozyg-window-threshold 0.05*.
1357 Finally, the total mean ROH lengths were calculated for each population separately using the “.hom”
1358 output-file for each of four length categories: 0.2 to 0.5 Mb, 0.5 to 1 Mb, 1 to 2 Mb, and 2 to 4 Mb
1359 (**FigureF2x**). Corresponding pipelines are provided in the corresponding GitHub repository
1360 (<https://github.com/Gwennid/africa-wgs-descriptive>).

1361 **Individual pairwise genome-wide genetic differentiation**

1362 We explored genome-wide genetic differentiation between pairs of individuals with Neighbor-Joining
1363 Tree (NJT) (Saitou and Nei, 1987; Gascuel, 1997), and Multi-Dimensional Scaling (MDS) approaches
1364 based on the pairwise matrix of Allele-Sharing Dissimilarities (ASD) (Bowcock *et al.*, 1991) computed
1365 using 14,182,615 genome-wide autosomal SNPs pruned for low LD, with PLINK version 1.90b4.9
1366 (Purcell *et al.*, 2007) *indep-pairwise* function (*--indep-pairwise 50 5 0.1*). We considered at first all 177
1367 unrelated individuals in our dataset to calculate the ASD matrix, and then subset this ASD matrix for
1368 the 73 Central and Southern African unrelated individuals original to this study before computing NJT
1369 and MDS analyses anew (**FigureF3x**).

1370 We computed the ASD matrix considering, for each pair of individuals, only those 14,182,615
1371 SNPs without missing data, using the *asd* software (v1.1.0a; <https://github.com/szpiech/asd>; Szpiech,
1372 2020). We computed the unrooted NJT using the *bionj* function of the R package *ape* and setting the
1373 branch-length option to “true”. We computed the MDS using the *cmdscale* function in R.

1374 The clustering software ADMIXTURE (Alexander, Novembre and Lange, 2009) allows
1375 researchers to further explore inter-individual genome-wide levels of dissimilarity and resemblance.
1376 Indeed, while it is tedious and often cognitively difficult to explore multiple combinations of
1377 dimensions of genetic variation using an MDS or a NJT approach, ADMIXTURE instead allows
1378 exploring K higher such dimensions at once (Pritchard, Stephens and Donnelly, 2000; Rosenberg, 2002;

1379 Falush, Stephens and Pritchard, 2003; Alexander, Novembre and Lange, 2009; Lawson, van Dorp and
1380 Falush, 2018). However, note that the visual representation of relative “distances” across pairs of
1381 individuals is lost in the classical bar-plot representation of ADMIXTURE, hence showing the
1382 complementarity of this descriptive analysis to the above MDS and NJT.

1383 We considered here only 840,031 genome-wide autosomal SNPs pruned for LD (r^2 threshold 0.1)
1384 and with Minimum Allele Frequency above 0.1, for the 177 worldwide unrelated individuals, using
1385 PLINK *indep-pairwise* and *maf* functions. We computed unsupervised ADMIXTURE version 1.3.0
1386 clustering with values of K ranging from 2 to 10, considering 20 independent runs for each value of K
1387 separately. We then calculated ADMIXTURE results symmetric-similarity-coefficient SSC for each
1388 value of K separately in order to find the groups of runs providing highly similar results (SSC>99.8%).
1389 Individual’s genotype membership proportions to each K cluster were then averaged, per individual,
1390 across such highly resembling runs, and then plotted (**FigureF4x**). SSC calculations, averaging results
1391 across similar runs, and producing barplots were conducted using the software PONG with the “greedy”
1392 algorithm (Behr *et al.*, 2016). Corresponding pipelines are provided in the corresponding GitHub
1393 repository (<https://github.com/Gwennid/africa-wgs-descriptive>).

1394 **Machine-Learning Approximate Bayesian Computation scenario-** 1395 **choice and posterior parameter estimation**

1396 We reconstructed the complex demographic history of Central and Southern African populations using
1397 machine-learning Approximate Bayesian Computation (Tavaré *et al.*, 1997; Beaumont, Zhang and
1398 Balding, 2002; Blum and François, 2010; Csilléry, François and Blum, 2012; Pudlo *et al.*, 2016). In
1399 principle, in ABC, researchers first simulate numerous genetic datasets under competing scenarios by
1400 drawing randomly a vector of parameter values for each simulation in distributions set a priori by the
1401 user. For each simulation separately, we then calculate a vector of summary statistics, thus
1402 corresponding to a vector of parameter values used for the simulation. The same set of summary
1403 statistics is then computed on the observed data. ABC scenario-choice then allows the researcher to
1404 identify which one of the competing scenarios produces the simulations for which summary-statistics
1405 are closest to the observed ones. Under this winning scenario, ABC posterior-parameter inference
1406 procedures allow the researcher to estimate the posterior distribution of parameter values most likely
1407 underlying the observed genetic patterns.

1408 **48 competing scenarios for the demographic history of Central and Southern** 1409 **African populations**

1410 We designed 48 competing demographic-history scenarios possibly underlying the genetic patterns
1411 observed in five Central and Southern African populations of five individuals each (**TableT1x**,
1412 **FigureF5x**); one Eastern Rainforest Hunter-Gatherer population (eRHG: Nsua, Ba.Twa, or Mbuti), one
1413 Western Rainforest Hunter-Gatherer population (wRHG: Baka, Ba.Kola, or Aka Mbat), one Eastern
1414 or Western RHG neighboring population (eRHGn: Ba.Kiga or Ba.Konjo; wRHGn: Nzimé or Ngumba),
1415 one Northern Khoe-San population (nKS: Ju|’hoansi or !Xun), and one Southern Khoe-San population
1416 (sKS: Karretjie People or Nama). Note that we did not include the Khutse-San population in ABC
1417 inferences for simplicity, as this population is located at intermediate geographic distances between the
1418 nKS and sKS groups of populations.

1419 The 48 competing scenarios differed in their ancient tree-topologies and relative timing of the more
1420 recent divergence events, combined with different (duration and intensities) possibly asymmetric gene-
1421 flow processes among ancient and recent lineages (**FigureF5x**). This design, while complex, allowed

1422 us, for the first time to our knowledge, to consider explicitly the expected confounding effects of gene-
1423 flow processes on otherwise different tree-topologies, which may also possibly induce certain
1424 reticulations among ancient lineages, while jointly estimating divergence times, gene-flow events'
1425 timing and intensities, and effective population size changes over time.

1426 Importantly, we relied on extensive previous findings having formally demonstrated the common
1427 origin of Congo Basin RHG populations (Patin *et al.*, 2009; Verdu *et al.*, 2009), and that of KS
1428 populations (Schlebusch *et al.*, 2012, 2020), respectively, and thus did not consider all possible tree-
1429 topologies for five tree-leaves. Moreover, we considered only a single RHGn sample in each
1430 combination, in turn from the East or the West of Central Africa, in order to simplify already highly
1431 complex scenarios, as we were not interested here in the detailed demographic history of RHGn which
1432 has been extensively studied previously (e.g. (Patin *et al.*, 2017; Fortes-Lima *et al.*, 2024)). This
1433 simplification was deemed reasonable as RHGn populations throughout the Congo Basin and into
1434 Southern Africa have previously been shown to be strongly more genetically resembling one another,
1435 compared to RHG and KS neighboring populations, or compared to much more genetically dissimilar
1436 RHG and KS populations, respectively (e.g. (Verdu *et al.*, 2009, 2013; Patin *et al.*, 2014, 2017; Lopez
1437 *et al.*, 2018, 2019; Fortes-Lima *et al.*, 2024)); a result that we also obtained here (**FigureF5x**).

1438

1439 *Eight competing topologies*

1440 We investigated eight competing topologies starting with a single ancestral population and resulting in
1441 five different sampled populations in the present (**FigureF5x**). The eight topologies differ in which
1442 lineage ancestral to modern KS, RHG, or RHGn diverged first from the two others, and in the relative
1443 order of divergence events internal to KS (nKS-sKS divergence) and internal to RHG (wRHG-eRHG
1444 divergence), in order to consider the known variable demographic histories of Sub-Saharan populations
1445 at a regional scale.

1446 Note that divergence and introgression times (all t and tad , see **FigureF5x**, **TableT2x**), were each
1447 randomly drawn in uniform prior distributions between 10 and 15,000 generations, thus effectively
1448 setting an upper limit for the most ancient divergence times among lineages at roughly 450,000 years
1449 ago considering an upper-bound of 30 years for human generation duration (Fenner, 2005). This was
1450 vastly anterior to the current estimates for the genetic or morphological emergence of *Homo sapiens*
1451 (Hublin *et al.*, 2017; Richter *et al.*, 2017), which allowed us to estimate *a posteriori* the most ancient
1452 divergences among our populations, without constraining our assumptions based on previous results
1453 obtained with different methods, data, and models. Furthermore, these prior-distribution boundaries
1454 allowed for gene-flow events among ancient lineages (see next section), to influence, potentially, the
1455 timing of the earliest divergence events among human lineages previously estimated (Schlebusch and
1456 Jakobsson, 2018; Fan *et al.*, 2023; Ragsdale *et al.*, 2023). Note that we retained for simulations only
1457 those vectors of randomly-drawn parameter-values that satisfied the chronological order of lineage
1458 divergences set for each eight topologies, respectively (**Figure5x**, **TableT2x**)

1459 In each eight topologies, we incorporated the possibility for changes in effective population sizes,
1460 N (**FigureF5x**), during history along each lineage separately by defining constant diploid effective
1461 population sizes parameters separately for each tree-branch, each drawn randomly in $U[10-100,000]$
1462 (**Figure5x**, **TableT2x**).

1463

1464 *Asymmetric instantaneous or recurring gene-flows and their intensities*

1465 Migration of individuals between populations is ubiquitous in human history (for an overview in Africa,
1466 see e.g. (Schlebusch and Jakobsson, 2018; Pfennig *et al.*, 2023)). Here we aimed at disentangling the
1467 nature of migration processes that may have occurred across pairs of lineages throughout the history of
1468 Central and Southern African populations.

1469 In particular, we first aimed at determining whether gene-flow events during history occurred
1470 relatively instantaneously, leading to reticulations in tree-topologies, or, conversely, occurred
1471 recurrently over longer periods of time leading to un- or weakly-differentiated lineages throughout
1472 history (Henn, Steele and Weaver, 2018; Hollfelder *et al.*, 2021; Ragsdale *et al.*, 2023). To do so, for
1473 each eight topologies described above (**FigureF5x**), we simulated gene-flow either as single-generation
1474 gene-flow pulses across pairs of lineages, or as constant recurring gene-flow across pairs of lineages in-
1475 between each lineage-divergence event. This design resulted in 16 competing scenarios encompassing
1476 eight different possible topologies, contrasting instantaneous gene-flow events with recurring gene-
1477 flows. Note that, for simplicity, we did not consider possible gene-flows between KS and RHG recent
1478 lineages, as genetic signatures of such recent migrations have never been identified in previous genetic
1479 studies nor in historical records to our knowledge.

1480 Importantly, we aimed at determining whether gene-flow occurred symmetrically, or not, across
1481 pairs of lineages, in particular in the past. Indeed, previous studies already identified recent asymmetric
1482 admixture processes between RHG and RHGn (Patin *et al.*, 2009, 2014; Verdu *et al.*, 2009, 2013), and
1483 such possible asymmetries have not been explored among ancient lineages in previous tree studies (Fan *et*
1484 *al.*, 2023; Ragsdale *et al.*, 2023), although they may be influencing ancient tree-topologies or
1485 reticulations across lineages and the subsequent past evolution of extant populations.

1486 To do so, for each gene-flow event in each 16 competing scenarios (**FigureF5x**), we parameterized
1487 separately the introgression of lineage A into lineage B, from that of lineage B into lineage A, by
1488 drawing randomly the corresponding parameter values independently in the same prior distribution
1489 (**TableT2x**). ABC posterior-parameter estimations would thus reveal possible asymmetries in gene-
1490 flows across pairs of lineages for each event separately, if identifiable from genomic data.

1491 Finally, for each one of the 16 competing combinations of topologies and gene-flow processes, we
1492 considered three classes of gene-flow intensities by setting different boundaries of gene-flow
1493 parameters' uniform prior-distributions (**TableT2x**): no to low possible gene-flow ($U[0, 0.001]$ for
1494 recurring processes, and $U[0, 0.01]$ for instantaneous processes), no to moderate gene-flow
1495 ($U[0,0.0125]$ for recurring processes, and $U[0, 0.25]$ for instantaneous processes), or no to intense gene-
1496 flow events ($U[0,0.05]$ for recurring processes and $U[0,1]$ for instantaneous processes). Note that, for
1497 each three classes of gene-flow intensities, each gene-flow parameter from one lineage to another at
1498 each point or period in time was drawn independently in the above prior-distributions' boundaries, and
1499 thus may differ across events within a scenario, and across scenarios (**FigureF5x**, **TableT2x**).

1500 Altogether, this design resulted in $16 \times 3 = 48$ competing scenarios for the complex demographic
1501 history of Central and Southern African populations. Importantly, note that for any given class of gene-
1502 flow process or intensity, the eight topologies are highly nested in certain parts of the space of parameter
1503 values (Robert, Mengersen and Chen, 2010). In particular, despite the fact that we did not consider
1504 competing scenarios with explicitly trifurcating ancient tree-topologies, scenarios in which the
1505 parameter values for the oldest and second oldest divergence times are similar are expected to provide
1506 results highly resembling those obtained with ancient trifurcation scenarios.

1507 All scenario-parameters are represented schematically in **FigureF5x**, and their prior distributions
1508 and constraints are indicated in **TableT2x**.

1509

1510 **Simulating genomic datasets**

1511 We performed simulations under the coalescent using *fastsimcoal2* (Excoffier and Foll, 2011; Excoffier
1512 *et al.*, 2013), for each one of the 48 scenarios described above, separately. For each simulation, we
1513 generated a vector of parameter values randomly drawn in prior-distributions and satisfying the
1514 topological constraints as described above and in **TableT2x**, with custom-made Python and Bash scripts
1515 available in the GitHub repository for this article (<https://github.com/Gwennid/africa-wgs-abc>). The

1516 genetic mutation model is based on the model deployed in (Jay, Boitard and Austerlitz, 2019). We
1517 simulated 100 independent loci (or “chromosomes”) with the same structure. Each such “chromosome”
1518 corresponds to a linkage block, with the following properties: the type of marker was “DNA”, the length
1519 of the loci 1 Mb, the recombination rate was 1×10^{-8} per base pair, the mutation rate 1.25×10^{-8} per
1520 base pair, and there was no transition bias (transition rate of 0.33).

1521 We performed 5,000 such simulations for each 48 competing scenarios separately in order to
1522 conduct Random Forest ABC scenario-choice, hence producing 240,000 separate simulated dataset
1523 each corresponding to a single vector of parameter values randomly drawn in prior distributions
1524 (**TableT2x**). Then we performed an additional 95,000 separate simulations under the winning scenario
1525 obtained with RF-ABC for each 54 separate combinations of observed population samples (see below).
1526 We thus reached 100,000 simulations under each winning scenario identified for each 54 observed
1527 datasets respectively, to be used for Neural Network ABC posterior parameter inference.

1528

1529 **Building observed genome-wide data-sets for 54 sets of five Central and Southern** 1530 **African populations**

1531 We considered 54 separate combinations of four eRHG, wRHG, nKS, sKS, and one eRHGn or wRHGn
1532 populations, each with five unrelated individuals (**TableT1x**). We prepared a callset of high-quality
1533 regions by applying the 1000 Genomes phase 3 accessibility mask and filtering out indels. We pieced
1534 high quality windows together to create 1 Mb-long windows. We selected 100 of these regions for
1535 which the total length from start to end is less than 1.2 Mb. We then extracted these 100 independent
1536 autosomal loci of 1 Mb each from the 25 individuals, hence mimetic of the simulated data.

1537

1538 **Calculating 337 summary statistics**

1539 For each simulated genetic dataset, we computed 337 summary statistics. 41 summary statistics were
1540 computed within each five populations separately (hence 205 statistics for five populations), 132
1541 summary statistics were computed across the five populations, for each one of the 54 combinations of
1542 five populations separately. In brief, all summary statistics were calculated with plink v1.90b4.9
1543 (Purcell *et al.*, 2007), R v3.6.1 (R Core Team, 2015), Python v2.7.15, bash, awk[1], and scripts
1544 developed by (Jay, Boitard and Austerlitz, 2019), available at
1545 https://gitlab.inria.fr/ml_genetics/public/demoseq/-/tree/master, and using the software *asd*
1546 (<https://github.com/szpiech/asd>). The list of calculated summary statistics is provided in **TableT3x**. The
1547 same statistics were computed for each one of the 54 observed data-sets separately, using the same
1548 computational tools and pipeline. Custom-made Python and Bash scripts for computations of all
1549 summary-statistics for each simulation and for the observed data are available in the GitHub repository
1550 for this article in the “fsc-simulations/code” folder (<https://github.com/Gwennid/africa-wgs-abc>).

1551 We used the complete set of 337 statistics to perform Random-Forest ABC scenario-choice, as this
1552 method is relatively fast and is unaffected by correlations among statistics (Pudlo *et al.*, 2016). We then
1553 considered a subset of 202 statistics among the 337 for Neural-Network ABC posterior-parameter
1554 estimation, in order to diminish computation time (**TableT3x**).

1555

1556 **Prior-checking simulations’ fit to the observed data**

1557 Before conducting ABC scenario-choice and posterior parameter estimation inferences, we checked
1558 that summary-statistics calculated on the observed data fell within the range of values of summary-
1559 statistics obtained from our simulations. First, we visually verified that each observed vector of
1560 summary-statistics computed separately from 54 combinations of five populations each, clustered with
1561 the 240,000 vectors of summary-statistics obtained from simulations under the 48 competing scenarios,

1562 using two-dimensional PCA calculated with the *prcomp* function R (**SupplementaryFigureSF1x**).
1563 Second, we used the *gfit* function in R to perform goodness-of-fit 100-permutations tests between the
1564 240,000 vectors of summary-statistics obtained from simulations under the 48 competing scenarios and,
1565 in turn, each vector of summary-statistics obtained from the 54 observed datasets respectively.

1566 Results both showed that our simulation design could successfully mimic summary-statistics
1567 observed in our five-population sample sets, for each 54 combinations of five populations separately.
1568 Hence, ABC inferences could be confidently conducted a priori based on the simulations and observed
1569 data here considered.

1570

1571 **Random Forest ABC grouped scenario-choice**

1572 We conducted series of Random Forest ABC scenario-choice procedures for different groups of
1573 scenarios (Pudlo *et al.*, 2016; Estoup *et al.*, 2018), elaborated specifically to address our different
1574 questions of interest regarding gene-flow processes, their intensities, and topological features of the
1575 history of Central and Southern African populations. Random Forest-ABC scenario-choice has proven
1576 to be performing efficiently and satisfactorily with a significantly lower number of simulations
1577 compared to any other ABC scenario-choice procedure. Moreover, RF-ABC scenario choice is
1578 insensitive to correlations among summary-statistics (Pudlo *et al.*, 2016; Estoup *et al.*, 2018).

1579 Each RF-ABC scenario-choice procedure presented below was conducted using the *predict.abcrf*
1580 and the *abcrf* functions with the “group” option of the *abcrf* package in R (Pudlo *et al.*, 2016), with
1581 1000 decision trees to train the algorithm after checking that error rates were appropriately low with the
1582 *err.abcrf* function, separately for each 54 combinations of five sampled populations (**FigureF6x**).

1583 We conducted cross-validation procedures considering in turn each one of the simulations as
1584 pseudo-observed data and all remaining simulations to train the algorithm, for each RF-ABC analysis
1585 separately (**SupplementaryFigureSF6xPart0**). While not necessarily predicting the outcome of the
1586 scenario-choice for the observed data, these cross-validation procedures provide us with a sense of the
1587 discriminatory power, *a priori*, of RF-ABC for our set of competing-scenarios, as well as empirical
1588 levels of nestedness among scenarios or groups of scenarios (Robert, Mengersen and Chen, 2010;
1589 Fortes-Lima *et al.*, 2021).

1590

1591 ***Instantaneous or recurring gene-flows in the history of Central and Southern African populations?***

1592 The 5,000 simulations performed under each one of the 48 competing scenarios were first gathered into
1593 two groups of 24 scenarios each (**FigureF5x**), corresponding to either instantaneous asymmetric gene-
1594 flow processes or to recurring constant asymmetric gene-flow between each pair of lineages,
1595 respectively. Both groups thus contained all simulations from the eight competing topologies and all
1596 three sets of possible gene-flow intensities (low, moderate, or high, see above), and only differed in the
1597 process of gene-flow itself. For each 54 combinations of five sampled populations separately, we
1598 performed such RF-ABC scenario-choice to determine the winning group of scenarios (**FigureF6x-**
1599 **panelA**).

1600

1601 ***Low, moderate, or high gene-flow intensities?***

1602 For each 54 combinations of five sampled populations separately, we performed RF-ABC scenario-
1603 choice in order to determine which class of intensities of gene-flows best explained the data, everything
1604 else (topology and gene-flow process) being-equal. We thus considered three groups of scenarios in our
1605 RF-ABC scenario-choice, each corresponding to “low”, “moderate”, or possibly “high” intensities, and
1606 each grouping the 16 scenarios corresponding to eight topologies and instantaneous or recurring gene-
1607 flow processes indiscriminately (**FigureF6x-panelB**).

1608 Then, we conducted, for each 54 population combinations separately, the same RF-ABC scenario-
1609 choice procedure to disentangle groups of gene-flow intensities “all topologies being equal”,
1610 considering here only the 24 scenarios from the winning group among instantaneous or recurring gene-
1611 flow processes obtained from the above scenario-choice procedure
1612 (**SupplementaryFigureSF6xPart1**).

1613 Finally, we performed, for each 54 combinations of five populations separately, RF-ABC scenario
1614 choice for the 48 competing scenarios grouped in six different groups (encompassing eight scenarios
1615 each), combining instantaneous or recurring processes with the three classes of intensities respectively
1616 (**FigureF6x-panelC**), thus “intersecting” the two corresponding group analyses (see above).

1617

1618 *Which ancestral Central or Southern African lineage diverged first?*

1619 The eight topologies considered in the 48 competing scenarios can also be grouped according to their
1620 ancient topologies, by considering which ancestral lineage diverged first from all others at the oldest
1621 divergence event in the tree (**FigureF5x**). The ancient lineage which separated first from the two others
1622 can either be the lineage ancestral to Northern and Southern Khoe-San populations (AKS, scenario
1623 topologies 1a, 1b, 1c in **FigureF5x**), the lineage ancestral to eastern and western Rainforest Hunter-
1624 Gatherer populations (ARHG, scenario topologies 2a, 2b, 2c in **FigureF5x**), or the lineage leading to
1625 Rainforest Hunter-Gatherer Neighboring populations (RHGn, scenario topologies 3a, 3b in
1626 **FigureF5x**).

1627 For each 54 combinations of five populations separately, we conducted RF-ABC scenario-choice
1628 across these three groups of scenarios, randomly drawing 2/3 of the 5,000 simulations per scenario for
1629 the scenarios 1a, 1b, and 1c, and for the scenarios 2a, 2b, and 2c, respectively, and kept all simulations
1630 for the scenarios 3a and 3b, in order to even the total number of simulations in competition among the
1631 three groups of scenarios. We thus performed RF-ABC scenario-choice across the three groups of
1632 topologies “all gene-flow processes and intensities being equal” (**FigureF6x-panelD**).

1633 Then we restricted the RF-ABC scenario-choice of the three-competing groups of topologies only
1634 for the 24 scenarios from the winning instantaneous or recurring gene-flow processes obtained above,
1635 all gene-flow intensities being equal (**FigureF6x-panelE**).

1636

1637 *Northern and Southern Khoe-San populations diverged before or after the divergence between* 1638 *Eastern and Western RHG?*

1639 The eight topologies considered in the 48 competing scenarios can alternatively be grouped according
1640 to the relative order of divergence-time between the divergence event among Northern and Southern
1641 KS populations, and that among Eastern and Western RHG populations; an important question to
1642 understand the relative duration of separate evolution of each groups of populations.

1643 To address this specific question we thus conducted, for the 54 combinations of five populations
1644 separately, RF-ABC scenario-choice procedures by grouping the 48 competing scenarios in two
1645 separate groups of 24 scenarios each. One group corresponded to KS populations diverging from one-
1646 another before the RHG divergence (scenario topologies 1b, 2b, 3b, and 1c in **FigureF5x**), and the other
1647 group where RHG diverged from one another before the KS divergence (scenario topologies 1a, 2a, 3a,
1648 and 2c in **FigureF5x**), instantaneous or recurring gene-flow and all “low”, “moderate”, or “high” gene-
1649 flow intensities being equal among the two groups (**FigureF6x-panelF**). Last, we conducted the same
1650 RF-ABC scenario-choice procedure between two groups of four topologies each, but restricted to the
1651 24 scenarios winning among the instantaneous or recurring gene-flow processes as determined from
1652 above procedures (**FigureF6x-panelG**).

1653 Finally, we conducted RF-ABC scenario choice procedures considering, respectively, all 48
1654 competing scenarios separately (**SupplementaryFigureSF6xPart2**), and all 24 scenarios winning

1655 among the instantaneous or recurring gene-flow processes as determined above
1656 (**SupplementaryFigureSF6xPart3**).

1657

1658 **Neural Network ABC posterior-parameter estimation**

1659 We intersected all the RF-ABC scenario-choice results considering the different groups of scenarios or
1660 all scenarios independently, as described above, in order to determine, which single scenario was
1661 winning in the majority of the 54 combinations of five Central and Southern African populations. We
1662 found that the Scenario i1-1b conservatively produced simulations whose observed summary-statistics
1663 most resembled those obtained in 25 out of the 54 possible population combinations here tested (see
1664 **Results**, the detailed list of 25 population combinations is provided in **SupplementaryTableST3x**).
1665 The second most often winning scenario, Scenario i1-3b, only succeeded for four population
1666 combinations among the 54 tested.

1667 Based on this *a posteriori* winning Scenario i1-1b, we performed an additional 95,000 simulations
1668 considering the same prior distributions and constraints among parameters, in order to obtain 100,000
1669 simulations for further ABC posterior parameter inferences, separately for each one of the 25 population
1670 combinations for which it was identified confidently as the winner, separately. As it remains difficult
1671 to estimate jointly the posterior distribution of all parameters of the complex scenarios here explored
1672 with RF-ABC (Raynal *et al.*, 2019), we instead conducted Neural Network ABC joint posterior
1673 parameter inferences using the *abc* package in R (Blum and François, 2010; Csilléry, François and
1674 Blum, 2012), following best-practices previously bench-marked for highly complex demographic
1675 scenarios (Jay, Boitard and Austerlitz, 2019; Fortes-Lima *et al.*, 2021). There are no rules of thumb in
1676 order to determine a priori the best number of neurons and tolerance rate to be set for training the neural
1677 network in NN-ABC (Jay, Boitard and Austerlitz, 2019).

1678 There are no evident criteria to choose *a priori* the best tolerance level and numbers of neurons in
1679 the neural network's hidden layer for parameterizing the NN-ABC posterior-parameter estimation
1680 procedure (Jay, Boitard and Austerlitz, 2019; Huang *et al.*, 2024). As the total number of parameters in
1681 the winning Scenario i1-1b was large (43, **FigureF5x** and **TableT2x**), and as the number of summary-
1682 statistics considered was also large (202, **TableT3x**), we chose to conduct the 25 NN-ABC posterior-
1683 parameter inferences for the 25 combinations of sampled populations considering, in-turn, 7, 14, 21,
1684 28, 35, 42, or 43 neurons in the hidden layer (43 being the assumed number of dimensions equal to the
1685 number of parameters for this scenario (Blum and François, 2010; Jay, Boitard and Austerlitz, 2019)),
1686 and a tolerance level of 0.01, thus considering the 1,000 simulations closest to each observed dataset
1687 out of the 100,000 performed, respectively. We found *a posteriori* that, considering 42 neurons in the
1688 hidden layer provided overall posterior parameter distributions departing from their priors, and in
1689 particular for divergence-times parameters which we were highly interested in, and therefore decided
1690 to provide all posterior parameter distributions' results using this parameterization (tolerance = 0.01,
1691 number of hidden neurons = 42) for the training of the neural network.

1692 We thus performed 25 separate NN-ABC joint parameter posterior inferences, using the
1693 “neuralnet” method option in the function *abc* of the *abc* package in R (Csilléry, François and Blum,
1694 2012), with logit-transformed (“logit” transformation option), an “epanechnikov” kernel, parameter
1695 values within parameter-priors' boundaries, a tolerance level of 0.01, and 42 neurons in the hidden
1696 layer. Adjusted posterior parameter distributions obtained with this method were then plotted for each
1697 parameter separately using a gaussian kernel truncated at the boundaries of the parameter prior-
1698 distributions. We then estimated the mode, median, mean, and 50% and 95% Credibility Intervals of
1699 these distributions *a posteriori*. We provide all posterior-parameter distributions together with their
1700 priors in figure-format (**FigureF7x**, **FigureF8x**, **SupplementaryFigureSF7xPart1-2**), and in table-
1701 format (**TableT4x**, **SupplementaryTableST4x**).

1702

1703 **GitHub repositories for this work**

1704

1705 <https://github.com/Gwennid/africa-wgs-descriptive>

1706

1707 <https://github.com/Gwennid/africa-wgs-abc>

1708

1709

1710 **Data availability statement**

1711 All raw whole-genome sequence FASTQ data original to this study are made available on the European
1712 Genome-Phenome Archive (<https://ega-archive.org/>), provided that requests conform to ethical
1713 requirements and informed consents provided by donors, as detailed in the corresponding Data Access
1714 Policies associated with each data repository. Data access requests are evaluated by the associated Data
1715 Access Committees. **EGA ACCESSION NUMBERS PENDING ACCEPTANCE.**

1716

1717

1718 **Supporting Information**

1719

1720 **8 Supplementary Figures and 4 Supplementary Tables**

1721

1722

1723

1724
 1725
 1726
 1727
 1728
 1729

Supplementary Table ST1x: Whole genome total variant counts among 177 unrelated worldwide individuals

Variant counts compared to the reference sequence for the human genome GRCh38 and previously reported variants in dbSNP 156. “SD” stands for “standard deviation”. See **Material and Methods** for details about the quality control, relatedness filtering, and variant count procedures.

	<u>Total variant count</u>		<u>Per individual variant count among 177 worldwide unrelated individuals</u>			
	Original dataset: 73 Central and Southern African unrelated individuals	Entire dataset: 177 worldwide unrelated individuals	Total	SD	Min	Max
Number of biallelic SNPs	26,780,319	36,272,545	4,245,000	375,379	3,093,785	4,642,544
Number of biallelic SNPs not in dbSNP 156	854,114	1,055,245	8,052	4,563	1,152	27,406
Number of multiallelic SNPs	241,428	257,906	56,330	5,162	38,057	61,716
Number of simple indels	2,454,965	3,159,306	385,391	35,522	278,237	421,976
Number of novel simple indels not in dbSNP 156	114,362	189,124	1,651	514	771	627
Number of complex indels	969,025	977,542	414,037	71,029	289,993	566,951
Number of variants appearing in a single sample	8,404,499	12,638,316	71,403	15,094	22,444	104,211
Number of filtered SNPs	2,940,629	3,296,775	247,172	48,579	66,597	350,377

1730
 1731

1732
1733
1734
1735
1736
1737
1738
1739

Supplementary Table ST2x: Whole genome biallelic SNPs counts in 46 worldwide populations

Mean number of bi-allelic SNPs across unrelated individuals from 46 worldwide populations, and associated standard deviations (when applicable), compared to the reference sequence for the human genome GRCh38 and compared to previously reported variants in dbSNP 156. See **Material and Methods** for details about the quality control, relatedness filtering, and variant count procedures. Populations are ordered in increasing mean number of SNPs per pop. The information original to this study is indicated in bold. “SD” stands for standard deviation. “na” stands for not applicable. Additional population information and geographical location of samples can be found in **Figure F1x** and **Table T1x**.

Population Name ¹	Sampling location	Dataset ²	N ³	Mean coverage X (SD)	Mean number of biallelic SNPs (SD)	Mean number of biallelic SNPs not in dbSNP 156 (SD)
Karitiana	Brazil	SGDP, HGDP	5	31.1 (7.8)	3,214,086 (69,241)	1,030 (547)
Papuan	Papua New Guinea	SGDP, HGDP	6	35.8 (7.8)	3,414,518 (74,893)	841 (290)
French	France	SGDP, HGDP	4	33.8 (8.7)	3,455,358 (77,954)	652 (225)
Dai	China	SGDP, KGP	6	37.7 (10.5)	3,480,744 (67,644)	695 (407)
CEU	United States of America	KGP	2	67.3 (0.2)	3,551,792 (7,819)	1,539 (57)
Saharawi	Western Sahara	SGDP	2	43.7 (2.2)	3,726,018 (27,793)	590 (21)
Mozabite	Algeria	SGDP	2	34.9 (1.6)	3,736,308 (69,393)	563 (11)
Somali	Kenya	SGDP	1	37.7 (na)	3,940,793 (na)	540 (na)
Coloured	South Africa	SAHGP	8	44.9 (2.2)	4,000,396 (97,447)	10,169 (2,319)
Dinka	Sudan	SGDP, HGDP	4	31.3 (5.8)	4,142,138 (97,001)	3,146 (4,811)
Maasai	Kenya	SGDP	2	38.6 (5.5)	4,169,732 (5,281)	654 (27)
Mandinka	Gambia	KGP	1	27.2 (na)	4,193,651 (na)	1,602 (na)
Mandenka	Senegal	SGDP, HGDP	4	30 (7.6)	4,217,398 (107,847)	914 (464)
Yoruba	Nigeria	SGDP, HGDP	4	31.1 (5)	4,233,500 (11,5371)	755 (195)
Ba.Kiga	Mukono (Uganda)	this study	5	39.4 (4.6)	4,274,466 (25,201)	5,391 (191)
Gambian	Gambia	SGDP	2	35.5 (0.8)	4,277,874 (12,976)	816 (107)
Esan	Nigeria	SGDP, KGP	3	39.7 (9.5)	4,284,041 (18,359)	704 (67)
Bantu_Kenya	Kenya	SGDP	2	38.7 (6.8)	4,288,533 (10,954)	611 (88)
Luo	Kenya	SGDP	2	34 (0.2)	4,293,632 (1,003)	722 (4)
Luhya	Kenya	SGDP, KGP	3	40.6 (12.1)	4,294,198 (28,096)	704 (37)
Igbo	Nigeria	SGDP	2	36.3 (8)	4,294,870 (1,867)	1,166 (9)
Lemande	Cameroon	SGDP	2	36.3 (0.5)	4,306,880 (5,284)	1,530 (37)
Bantu_Herero	Namibia	SGDP	2	38.4 (2.2)	4,318,604 (33,525)	690 (161)
Ba.Konjo	Mulimassenge (Uganda)	this study	5	43 (6.4)	4,322,177 (6,336)	5,349 (773)
Nzime	Messea (Cameroon)	this study	5	35.5 (2.6)	4,322,734 (18,671)	5,305 (208)
Ngumba	Dispersed between Lolodorf and Kribi (Cameroon)	this study	5	40.2 (2.2)	4,332,324 (9,273)	4,810 (552)
Kongo	Cameroon	SGDP	1	42.8 (na)	4,336,731 (na)	1233 (na)
Mende	Sierra Leone	SGDP, KGP	3	39.6 (8.3)	4,341,512 (13,486)	866 (180)
Zulu	South Africa	SAHGP	1	42.8 (na)	4,376,691 (na)	8,146 (na)
Xhosa	South Africa	SAHGP	8	42.4 (1.1)	4,385,032 (56,902)	10,085 (1,300)
Ba.Twa	Kebiremu, Byumba, Kitariro, Mgunu, Nteko (Uganda)	this study	6	53.8 (17.7)	4,426,223 (32,552)	20,251 (1,022)
Bantu_Tswana	South Africa	SGDP	2	36.7 (3.1)	4,430,870 (55,537)	720 (37)
Mbuti	Democratic Republic of Congo	SGDP, HGDP	5	32.1 (7.4)	4,434,425 (110,658)	1,227 (488)
Sotho	South Africa	SAHGP	7	42.9 (6.4)	4,446,796 (23,172)	12,040 (1,150)
Ba.Kola	Dispersed between Lolodorf and Kribi (Cameroon)	this study	5	39.2 (1.9)	4,475,922 (50,261)	6,995 (1,948)
Nsua	Bundimassoli (Uganda)	this study	5	37.6 (4.9)	4,486,217 (22,774)	16,443 (984)
Baka	Bosquet (Cameroon)	this study	7	40.7 (3.5)	4,503,977 (23,946)	3,927 (144)
Bi.Aka_Mbati	Bombeketi section of Bagandou (Central African Republic)	this study	5	43.4 (13.6)	4,514,318 (17,376)	4,587 (198)
Biaka	Central African Republic	SGDP	2	39.1 (0.3)	4,521,140 (8,155)	726 (52)
Nama	Windhoek (Namibia)	this study	5	48.9 (4.3)	4,530,168 (18,807)	27,902 (462)
Ju 'hoansi_comp	Namibia	SGDP, HGDP	4	33.9 (7)	4,554,593 (103,099)	1,382 (273)
#Khomani	South Africa	SGDP	2	44.7 (2.8)	4,559,846 (79,851)	1,364 (48)
Khutse_San	Kutse Game reserve (Botswana)	this study	5	46.3 (6.7)	4,592,852 (27,794)	27,989 (1,353)
Karretjie	Colesberg (South Africa)	this study	5	46.6 (5.4)	4,601,724 (26,951)	27,720 (1,968)
!Xun	Omega camp (Namibia) and Schmidtsdrift (Sout Africa)⁴	this study	5	45.6 (1.9)	4,614,835 (10,047)	28,923 (1,371)
Ju 'hoansi	Tsumkwe (Namibia)	this study	5	49 (4.6)	4,628,957 (8,139)	26,665 (6,491)

1740
1741
1742
1743
1744
1745
1746

¹Population Names are self reported for the original dataset presented in this study

²SGDP: Mallick et al 2016; KGP: Auton et al. 2015; HGDP: Meyer et al. 2012, Rasmussen et al. 2014; SAHGP: Choudhury et al. 2017

³Number of unrelated individuals considered in all analyses in this study (see Material and Methods)

⁴The place of origin of the !Xun is around Menongue in Angola.

1747
 1748
 1749
 1750
 1751
 1752
 1753

SupplementaryTableST3x: 25 combinations of five population samples out of 54 combinations tested for which Scenario i1-1b is winning with Random-Forest ABC among the 48 competing scenarios.

Values in the table correspond to posterior densities plotted in the left panels of **FigureF7x**, **FigureF8x**, and **SupplementaryFigureSF7xPart1-Part2**. Parameter definitions and priors are provided in **TableT2x** and represented graphically in the corresponding scenario panel of **FigureF5x**. Values in italic are not satisfactorily departing from the priors.

1754

	nKS	sKS	RHGn	wRHG	eRHG
<i>Combination 1</i>	Ju hoansi	Karretjie	Nzime	Baka	Nsua
<i>Combination 2</i>	Ju hoansi	Karretjie	Ngumba	Ba.Kola	Nsua
<i>Combination 3</i>	Ju hoansi	Karretjie	Ngumba	Ba.Kola	Ba.Twa
<i>Combination 4</i>	Ju hoansi	Karretjie	Nzime	Bi.Aka_Mbati	Nsua
<i>Combination 5</i>	Ju hoansi	Karretjie	Nzime	Bi.Aka_Mbati	Ba.Twa
<i>Combination 6</i>	Ju hoansi	Karretjie	Ngumba	Ba.Kola	Mbuti
<i>Combination 7</i>	Ju hoansi	Karretjie	Nzime	Bi.Aka_Mbati	Mbuti
<i>Combination 8</i>	Ju hoansi	Karretjie	Ba.Konjo	Baka	Mbuti
<i>Combination 9</i>	Ju hoansi	Karretjie	Ba.Konjo	Ba.Kola	Mbuti
<i>Combination 10</i>	Ju hoansi	Karretjie	Ba.Konjo	Bi.Aka_Mbati	Mbuti
<i>Combination 11</i>	Ju hoansi	Nama	Ba.Konjo	Baka	Nsua
<i>Combination 12</i>	!Xun	Nama	Ba.Konjo	Baka	Nsua
<i>Combination 13</i>	Ju hoansi	Nama	Nzime	Baka	Nsua
<i>Combination 14</i>	Ju hoansi	Nama	Ngumba	Ba.Kola	Nsua
<i>Combination 15</i>	Ju hoansi	Nama	Nzime	Baka	Ba.Twa
<i>Combination 16</i>	Ju hoansi	Nama	Ngumba	Ba.Kola	Ba.Twa
<i>Combination 17</i>	Ju hoansi	Nama	Nzime	Bi.Aka_Mbati	Ba.Twa
<i>Combination 18</i>	Ju hoansi	Nama	Ba.Konjo	Ba.Kola	Nsua
<i>Combination 19</i>	Ju hoansi	Nama	Ba.Kiga	Ba.Kola	Ba.Twa
<i>Combination 20</i>	Ju hoansi	Nama	Ba.Kiga	Bi.Aka_Mbati	Ba.Twa
<i>Combination 21</i>	!Xun	Nama	Nzime	Bi.Aka_Mbati	Nsua
<i>Combination 22</i>	!Xun	Nama	Nzime	Bi.Aka_Mbati	Ba.Twa
<i>Combination 23</i>	!Xun	Karretjie	Nzime	Bi.Aka_Mbati	Nsua
<i>Combination 24</i>	!Xun	Karretjie	Nzime	Bi.Aka_Mbati	Ba.Twa
<i>Combination 25</i>	!Xun	Karretjie	Ba.Konjo	Bi.Aka_Mbati	Nsua

1755 **SupplementaryTableST4x: NN-ABC posterior parameter estimation of all parameters in Scenario i1-1b for results**
1756 **from each 25 sets of five populations each, separately.**

1757 Each sampled-population combination is provided as a single line. Table-headers in bold indicate the categorization of sampled
1758 populations into Northern Khoisan (nKS), Southern Khoisan (sKS), Rainforest Hunter-Gatherer neighbors (RHGn), Western
1759 Rainforest Hunter-Gatherer (wRHG), or Eastern Rainsforest Hunter-Gatherer (eRHG) groups, as per scenario topology and
1760 simulation design explicated in **Material and Methods** and in **FigureF5x**. Note that each sampled populations considered in
1761 the RF-ABC scenario choice analyses are represented twice or more in the 25 combinations. See “.xlsx” file as the table is to
1762 large to reasonably fit in a A4 page format.

1763

1764

1765

1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777

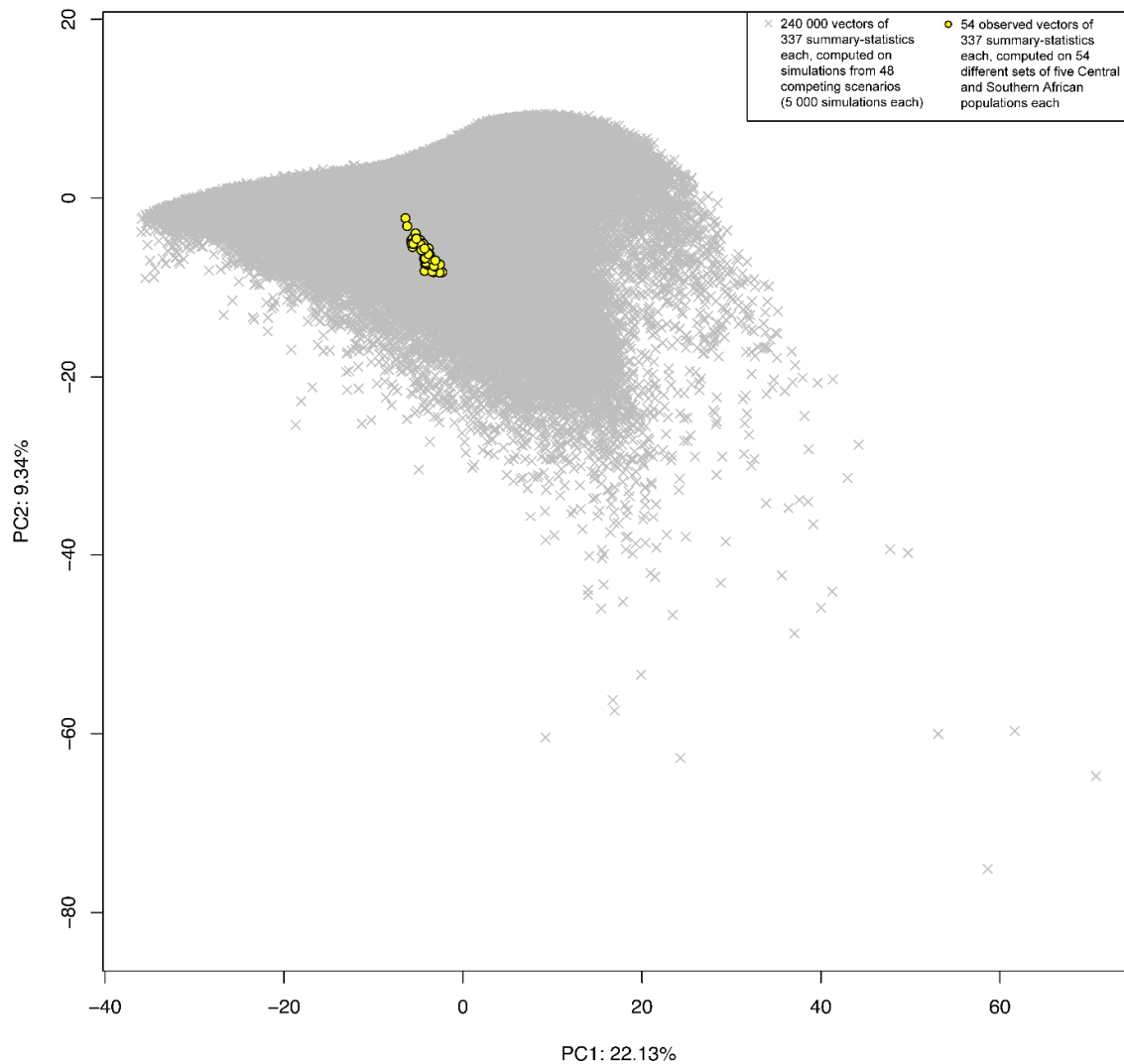
SupplementaryFigureSF1x: PCA for ABC prior-checking of simulations fit to the observed data.

We calculated 337 summary-statistics for each one of the 240 000 simulated data sets, 5000 under each one of the 48 competing scenarios (**FigureF5x** and **Material and Methods**), which we projected on the first two axes of a principal component analysis. Each vector (corresponding thus to a single simulation) is represented by a gray cross. We separately computed the same 337 summary-statistics, separately for each one of the 54 sets of five observed Central and Southern African populations included in our analyses. Each observed vector is then projected, in turn, on the PCA obtained from simulations only. Each such observed vector is represented as a single yellow dot. We can see that all the 54 observed sets of populations fall well within the space of simulated dataset, thus allowing us a priori to conduct machine-learning ABC scenario choice and posterior parameter estimation procedures.

1778
1779
1780

SupplementaryFigureSF1x

PCA prior checking of simulation fit to the observed data



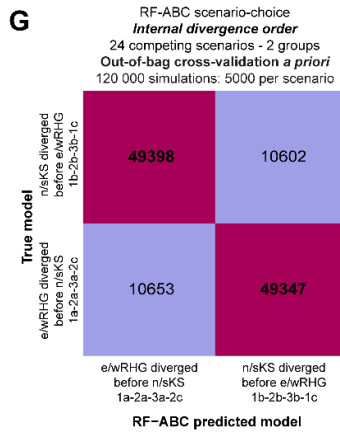
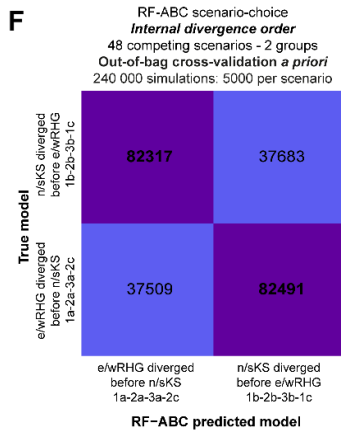
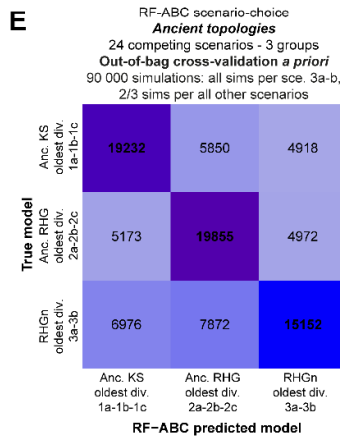
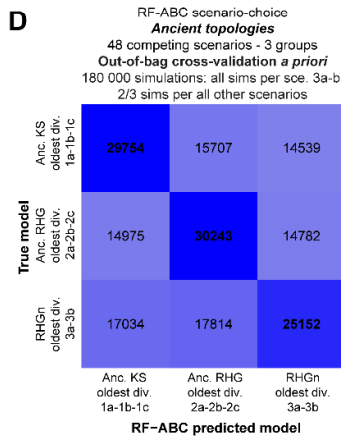
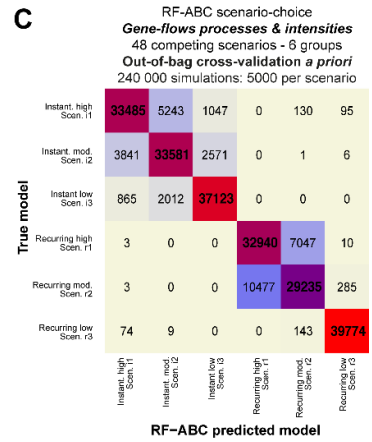
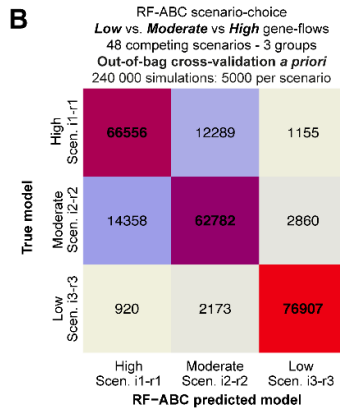
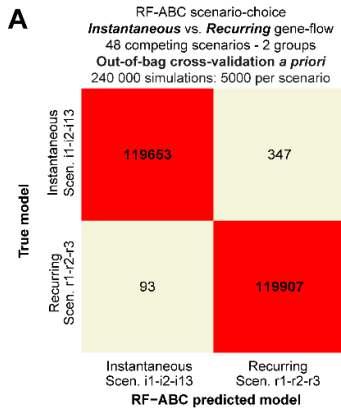
1781 **SupplementaryFigureSF6xPart0: Out-of-bag cross-validation errors for Random-Forest ABC scenario-choices among**
1782 **48 competing scenarios.**

1783 All RF-ABC cross-validation analyses are conducted *a priori*, without using any observed data, by considering, in-turn, each
1784 simulation as observed data while all remaining simulations are used to train the Random Forest. Each panel provides the *a*
1785 *priori* cross-validation errors for the corresponding RF-ABC scenario-choice analysis detailed in each panel in **FigureF6x**.
1786 While these cross-validation errors do not predict the specific outcome of RF-ABC prediction for the observed data, they
1787 provide levels of scenarios and groups of scenarios nestedness *a priori* for the entire space of parameters used for simulations.
1788 (A) shows that intermediate and recurring gene-flow processes are, all topologies and gene-flow intensities “being equal”,
1789 clearly differentiable by RF-ABC and little nested. (B) and (C) show that different classes of gene-flow intensities are also
1790 relatively clearly identifiable, albeit more nested, as expected by design since the “high gene-flow” class comprises parameter
1791 values of the moderate and the low intensity classes of gene-flows (**Material and Methods, TableT2x**). (D) and (E) show
1792 that ancient tree topologies are relatively well distinguishable *a priori* all gene-flow processes “being equal”, despite some
1793 amounts of errors due to the expected nestedness in the space of parameter values where ancient divergence times are
1794 resembling and thus undistinguishable. (F) and (G) show that the chronological order between the divergence-time of Northern
1795 and Southern KS lineages, and that of Western and Eastern RHG lineages, are also relatively well distinguishable a priori with
1796 RF-ABC scenario-choice, albeit with an increased cross-validation error expected due to high levels of nestedness in the spaces
1797 of parameter-values where both separate events may occur at similar time.

1798
1799
1800

1801
1802
1803
1804

SupplementaryFigureSF6xPart0



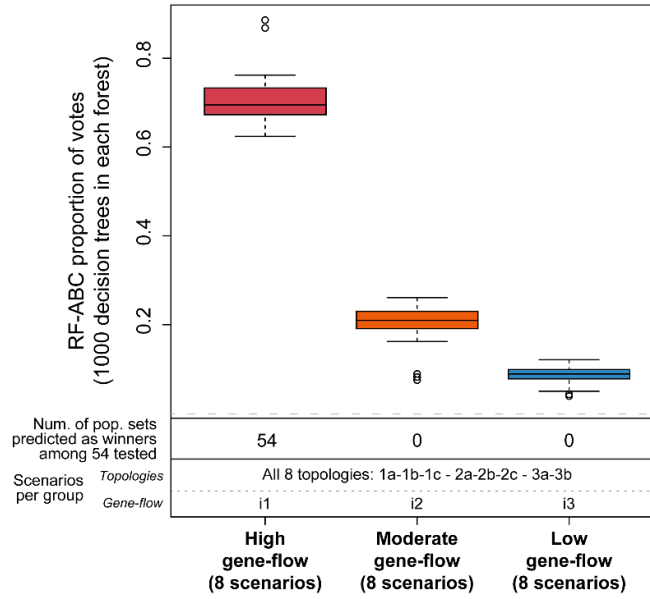
1805 **SupplementaryFigureSF6xPart1: RF-ABC scenario choice among 24 competing scenarios with instantaneous**
1806 **asymmetric gene-flow processes, for three groups of gene-flow intensities, all topologies being equal.**
1807 (A) corresponds to the scenario-choice results for the same test as in **FigureF6x-PanelC**, restricted to the 24 competing
1808 scenarios considering instantaneous gene-flow processes only, all scenarios topologies “being equal”. Description of the x and
1809 y axis legends are provided in **FigureF6x**. (B) corresponds to the RF-ABC cross-validation prior error for the scenario-choice
1810 procedure conducted in (A). Errors were obtained without considering observed data and using, instead, all simulations in turn
1811 as pseudo-observed data and the remaining simulations to train the RF.
1812
1813
1814

1815
 1816
 1817
 1818

SupplementaryFigureSF6xPart1

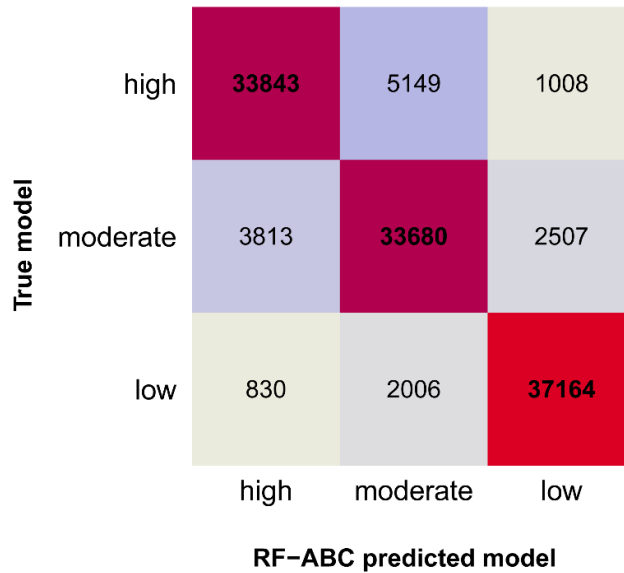
A

RF-ABC scenario-choice
Low vs. Moderate vs High
Instantaneous gene-flows
 24 competing scenarios - 3 groups
 54 population sets



B

RF-ABC scenario-choice
Out-of-bag cross-validation a priori
 120 000 simulations (no real data)



1819 **SupplementaryFigureSF6xPart2: RF-ABC scenario choice among 48 competing scenarios without groups.**
1820 (A) corresponds to the RF-ABC scenario-choice results obtained when considering all 48 scenarios (5,000 simulations per
1821 scenario), as separate competitors without grouping them, for each 54 combinations of five sampled-populations separately.
1822 Description of the x and y axis legends are provided in **FigureF6x. (B)** corresponds to the RF-ABC cross-validation prior error
1823 for the scenario-choice procedure conducted in (A). Errors were obtained without considering observed data and using, instead,
1824 all simulations in turn as pseudo-observed data and the remaining simulations to train the RF. Note that, in (B), whichever the
1825 intensity of the gene-flow, when considering only recurring gene-flow processes, topologies are harder to distinguish from one
1826 another a priori. Conversely, such difficulty to distinguish a priori among topologies remains overall limited for instantaneous
1827 gene-flow processes, and errors increase only when considering the possibility of intense instantaneous gene-flows. However,
1828 as mentioned throughout the article, the capacity to distinguish a priori among scenarios in the entire space of summary-
1829 statistics values does not predict the power to predict a winning scenario in the specific space occupied by observed data (Pudlo
1830 et al. 2016).
1831
1832
1833

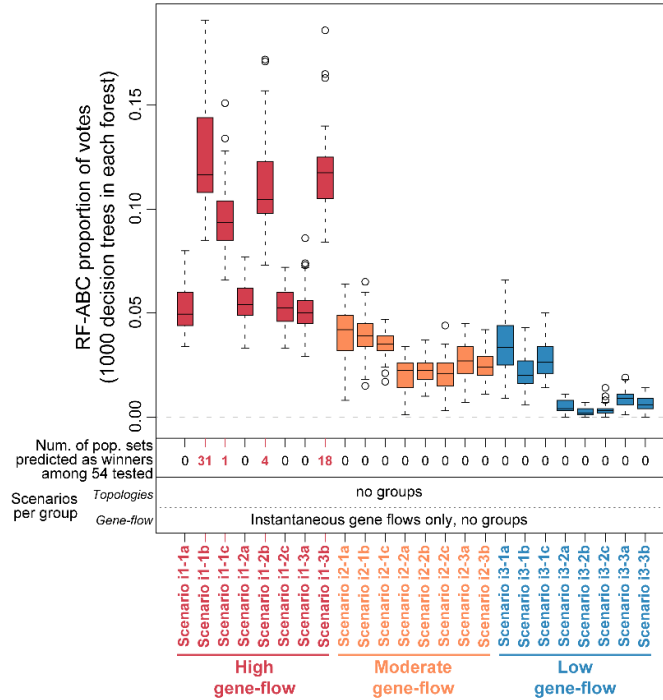
1838 **SupplementaryFigureSF6xPart3: RF-ABC scenario choice among 24 competing scenarios without groups.**
1839 (A) corresponds to the RF-ABC scenario-choice results obtained when considering all 24 instantaneous gene-flow scenarios
1840 (5,000 simulations per scenario), as separate competitors without grouping them, for each 54 combinations of five sampled-
1841 populations separately. Description of the x and y axis legends are provided in **FigureF6x**. (B) corresponds to the RF-ABC
1842 cross-validation prior error for the scenario-choice procedure conducted in (A). Errors were obtained without considering
1843 observed data and using, instead, all simulations in turn as pseudo-observed data and the remaining simulations to train the
1844 RF. Note that prediction and cross-validation error results resemble those obtained for these 24 specific scenarios obtained
1845 when considering instead all the 48 scenarios in competition (**SupplementaryFigureSF6xPart2**).
1846
1847
1848

1849
1850
1851
1852

SupplementaryFigureSF6xPart3

A

RF-ABC scenario-choice
Topologies and Low vs. Moderate vs High
Instantaneous gene-flows
24 competing scenarios - no groups
54 population sets



B

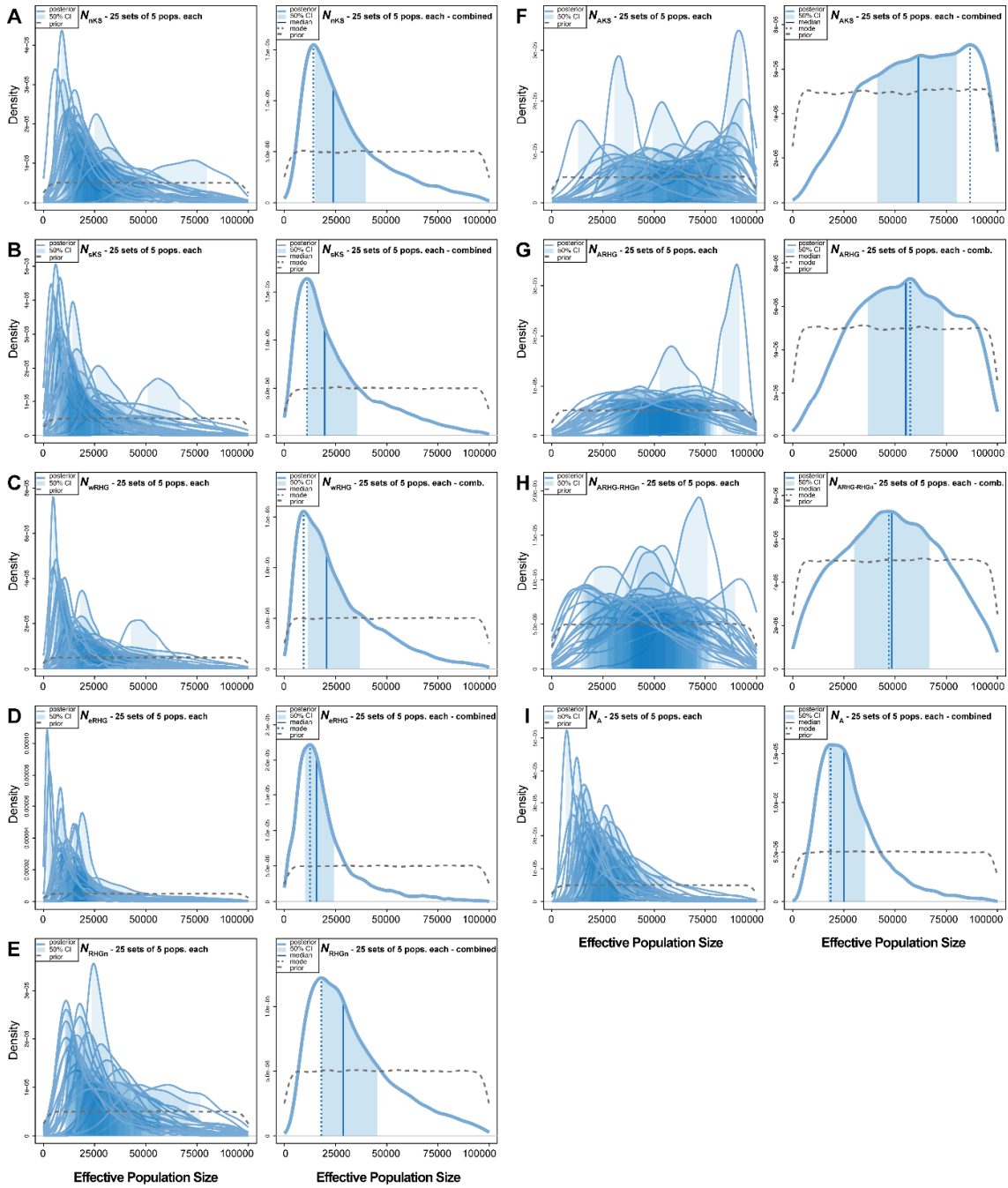
RF-ABC scenario-choice
24 competing scenarios - no groups
Out-of-bag cross-validation a priori
120 000 simulations: 5000 per scenario
no real data

True model	Scenario i1-1a	Scenario i1-1b	Scenario i1-1c	Scenario i1-2a	Scenario i1-2b	Scenario i1-2c	Scenario i1-3a	Scenario i1-3b	Scenario i2-1a	Scenario i2-1b	Scenario i2-1c	Scenario i2-2a	Scenario i2-2b	Scenario i2-2c	Scenario i2-3a	Scenario i2-3b	Scenario i3-1a	Scenario i3-1b	Scenario i3-1c	Scenario i3-2a	Scenario i3-2b	Scenario i3-2c	Scenario i3-3a	Scenario i3-3b	
Scenario i1-1a	830	202	226	647	188	819	615	193	122	32	55	93	92	350	112	58	40	5	27	35	60	174	13	12	
Scenario i1-1b	209	624	989	215	605	220	210	599	125	95	351	33	123	96	86	113	57	23	126	11	33	25	11	21	
Scenario i1-1c	145	471	954	133	430	105	124	444	80	78	382	25	97	57	101	40	22	146	5	47	31	12	19		
Scenario i1-2a	738	189	248	626	209	678	659	185	108	27	58	104	98	354	112	57	33	2	19	29	78	161	13	15	
Scenario i1-2b	191	599	940	214	695	232	199	624	78	97	316	48	128	74	95	118	69	20	132	19	52	38	7	15	
Scenario i1-2c	471	188	154	494	186	707	499	179	60	15	36	85	57	411	108	40	27	4	17	25	44	183	17	13	
Scenario i1-3a	726	182	216	622	226	950	657	197	110	28	64	108	82	323	117	54	30	7	16	22	72	169	11	11	
Scenario i1-3b	199	619	985	210	660	210	211	591	101	100	323	55	129	78	76	138	58	19	133	13	37	37	5	12	
Scenario i2-1a	20	37	23	7	25	9	14	30	148	318	282	638	243	472	560	198	181	52	50	110	57	153	44	21	
Scenario i2-1b	12	57	39	3	35	3	7	57	372	1099	108	154	604	74	153	602	93	108	230	37	84	22	16	34	
Scenario i2-1c	5	91	118	4	81	3	5	81	133	531	259	36	293	64	57	290	40	59	434	8	25	17	7	37	
Scenario i2-2a	47	6	3	57	14	46	35	14	675	127	69	1078	354	1103	547	159	82	25	23	159	94	245	22	14	
Scenario i2-2b	19	27	14	40	35	19	39	23	248	743	536	315	1230	282	190	613	56	86	162	51	156	51	13	42	
Scenario i2-2c	130	9	17	79	12	165	98	8	298	42	53	527	151	100	255	69	26	7	13	73	34	411	33	14	
Scenario i2-3a	30	22	16	17	23	31	27	19	704	117	79	630	206	520	1411	455	97	27	35	129	67	189	93	56	
Scenario i2-3b	12	24	33	9	39	13	12	43	218	624	577	109	581	76	459	1499	61	105	182	36	122	27	36	103	
Scenario i3-1a	0	8	3	0	9	0	1	4	64	16	17	6	5	42	28	238	552	399	334	102	45	382	92		
Scenario i3-1b	0	4	4	0	4	0	0	3	20	39	22	3	18	2	16	80	726	2050	923	112	439	24	122	389	
Scenario i3-1c	0	18	25	1	18	1	1	17	11	61	150	3	14	2	5	51	248	751	423	19	168	16	40	165	
Scenario i3-2a	5	0	0	15	1	5	7	1	20	3	1	50	17	31	77	30	458	89	27	20	19	712	913	396	123
Scenario i3-2b	6	1	0	7	1	5	14	1	4	2	7	15	55	8	27	57	111	294	49	580	283	411	104	416	
Scenario i3-2c	43	1	0	31	3	29	28	1	22	1	3	50	10	131	43	7	143	14	17	813	270	155	154	33	
Scenario i3-3a	0	0	0	0	1	2	1	0	13	3	0	1	1	1	73	23	614	81	23	400	118	64	291	764	
Scenario i3-3b	0	0	2	1	0	0	0	1	2	5	8	1	3	0	18	89	143	334	88	94	548	28	773	222	

1853 **SupplementaryFigureSF7xPart1: ABC posterior distribution of effective population size parameters.**
1854 Neural Network Approximate Bayesian Computation (Csilléry et al. 2012), posterior parameter joint estimations of Effective
1855 population sizes N_e (in generations before present) for 25 sets of five Central and Southern African populations for which the
1856 winning scenario identified by RF-ABC was Scenario i1-1b (**FigureF5x, SupplementaryTableST3x**). NN ABC posterior
1857 parameter estimation procedures were conducted using 100,000 simulations under Scenario i1-1b, each simulation
1858 corresponding to a single vector of parameter values drawn randomly from prior distributions provided in **TableT2x**. We
1859 considered 42 neurons in the hidden layer of the NN and a tolerance level of 0.01, corresponding to the 1,000 simulations
1860 providing summary-statistics closest to the observed ones, for each 25 separate analyses. NN posterior estimates are based on
1861 the logit transformation of parameter values using an Epanechnikov kernel between the corresponding parameter's prior
1862 bounds (see **Material and Methods** and **TableT2x**). Posterior parameter densities are represented with solid blue lines. 50%
1863 Credibility Intervals are represented as the light blue area under the density. The median and mode values are represented as
1864 a solid and dotted blue vertical line, respectively. Parameter prior distributions are represented as dotted grey lines. For all
1865 panels, the left plots represent the NN-ABC posterior parameter distributions for each 25 sets of five Central and Southern
1866 African populations winning under Scenario i1-1b, separately (**SupplementaryTableST3x** and **SupplementaryTableST4x**).
1867 See **FigureF5x** and **TableT2x** for all parameters' descriptions in each panel. Results are also provided in **TableT4x**.
1868
1869
1870

1871
1872
1873
1874

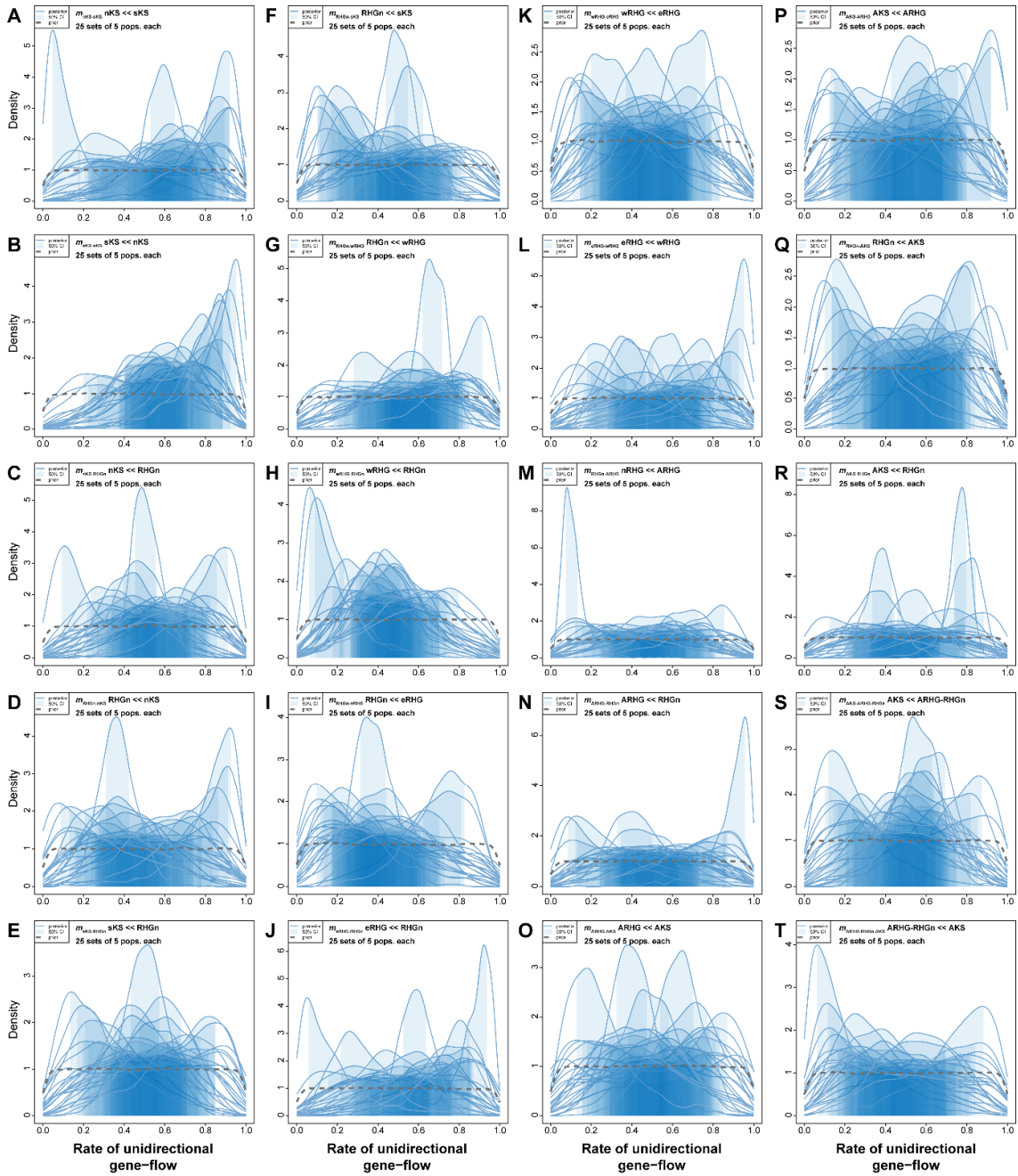
SupplementaryFigureSF7xPart1



1875 **SupplementaryFigureSF7xPart2: ABC posterior distribution of potentially asymmetric instantaneous gene-flow**
1876 **intensity parameters.**
1877 Neural Network Approximate Bayesian Computation (Csilléry et al. 2012), posterior parameter joint estimations of gene-flow
1878 intensity parameters m (in generations before present) for 25 sets of five Central and Southern African populations for which
1879 the winning scenario identified by RF-ABC was Scenario i1-1b (**FigureF5x, SupplementaryTableST3x**). NN ABC posterior
1880 parameter estimation procedures were conducted using 100,000 simulations under Scenario i1-1b, each simulation
1881 corresponding to a single vector of parameter values drawn randomly from prior distributions provided in **TableT2x**. We
1882 considered 42 neurons in the hidden layer of the NN and a tolerance level of 0.01, corresponding to the 1,000 simulations
1883 providing summary-statistics closest to the observed ones, for each 25 separate analyses. NN posterior estimates are based on
1884 the logit transformation of parameter values using an Epanechnikov kernel between the corresponding parameter's prior
1885 bounds (see **Material and Methods** and **TableT2x**). Posterior parameter densities are represented with solid blue lines. 50%
1886 Credibility Intervals are represented as the light blue area under the density. Parameter prior distributions are represented as
1887 dotted grey lines. For all panels, the plots represent the NN-ABC posterior parameter distributions for each 25 sets of five
1888 Central and Southern African populations winning under Scenario i1-1b, separately (**SupplementaryTableST3x** and
1889 **SupplementaryTableST4x**). See **FigureF5x** and **TableT2x** for all parameters' descriptions in each panel. Note that, overall,
1890 parameters are relatively little departing from their priors for numerous sets of population combinations among the 25. Also,
1891 posterior parameter distributions that are substantially departing from their priors are often highly differing from one another
1892 for each parameter. Therefore, we considered these parameters as unsatisfactorily estimated in our analyses and discuss this
1893 limitation in the **Discussion** section of the main text.
1894
1895
1896

1897
1898
1899
1900

SupplementaryFigureSF7xPart2



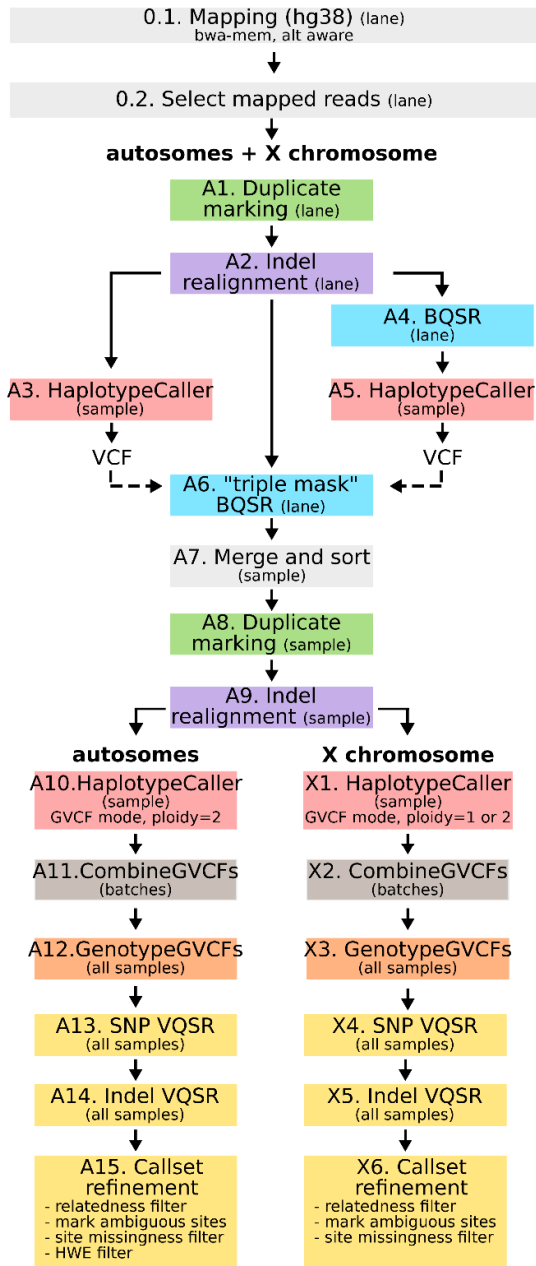
1901 **SupplementaryFigureSF8x**: Schematic overview of the sequencing data processing pipeline.
1902 All template codes with accompanying detailed explanations for all steps are provided in **Material and Methods** and in the
1903 corresponding GitHub repository (<https://github.com/Gwennid/africa-wgs-descriptive>).
1904
1905

1906 **SupplementaryFigureSF8x**

1907

1908

1909



1910 **Acknowledgments**

1911 We are grateful to all subjects who participated in this research. The computations were performed at
1912 the Swedish National Infrastructure for Computing (SNIC-UPPMAX). Sequencing was performed by
1913 the SNP&SEQ Technology Platform in Uppsala. We thank the platform Paléogénomique et Génétique
1914 Moléculaire (P2GM) of the MNHN at the Musée de l'Homme in Paris for helping to prepare, in part,
1915 the Central African samples for sequencing. We thank the Working Group of Indigenous Minorities in
1916 Southern Africa (WIMSA) and the South African San Council for their support and facilitating
1917 fieldwork. The project was approved by the South African San Council.
1918

1919 **Ethical statement**

1920 DNA samples from individuals were collected with the subjects' informed consent following the
1921 Declaration of Helsinki guidelines and approved by the following Ethics committees: Institut de
1922 Recherche pour le Développement (IRD, France); Cameroon Ministry of Scientific Research and
1923 Technology; Central African Republic Ministère de l'Enseignement Supérieur et de la Recherche
1924 Scientifique; France Ministère de l'Enseignement Supérieur, de la Recherche et de l'Innovation
1925 (Comité de Protection des Personnes DC-2016-2740); French Commission Nationale Informatique et
1926 Liberté CNIL (déclaration n°1972648); Uganda National Council for Science and Technology;
1927 Institutional Review Board of the University of Chicago Biological Sciences Division and the Pritzker
1928 School of Medicine University of Chicago Hospitals (protocol number 16986A); Committee for the
1929 Protection of Human Subjects of the Trustees of Dartmouth College and the Dartmouth-Hitchcock
1930 Medical Center (protocol number 22410); Human Research Ethics Committee (Medical) at the
1931 University of the Witwatersrand, Johannesburg (Protocol Numbers: M980553, M050902, M090576,
1932 M10270, M180654); Working Group of Indigenous Minorities in Southern Africa (WIMSA); South
1933 African San Council (SASC); and the Swedish Ethical Review Authority (Dnr 2019-05174).
1934

1935 **Funding**

1936 CS is funded by the European Research Council (ERC) under the European Union's Horizon 2020
1937 research and innovation program (grant agreement No. 759933) and the Knut and Alice Wallenberg
1938 foundation. MJ is funded by the Knut and Alice Wallenberg foundation, and the
1939 Vetenskapsrådet/Swedish Research Council (no: 642-2013-8019 and 2018-05537). PV was funded in
1940 part by the UMR7206 Eco-anthropology and the French Agence Nationale de la Recherche ANR-
1941 METHIS (15-CE32-0009-1).
1942
1943

1944 **REFERENCES**

- 1945 Agranat-Tamir, L., Mooney, J.A. and Rosenberg, N.A. (2024) 'Counting the genetic
1946 ancestors from source populations in members of an admixed population', *GENETICS*.
1947 Edited by J. Novembre, 226(4), p. iyae011. Available at:
1948 <https://doi.org/10.1093/genetics/iyae011>.
- 1949 Alexander, D.H., Novembre, J. and Lange, K. (2009) 'Fast model-based estimation of
1950 ancestry in unrelated individuals', *Genome Research*, 19(9), pp. 1655–1664. Available at:
1951 <https://doi.org/10.1101/gr.094052.109>.
- 1952 Bader, G.D. *et al.* (2022) 'Rethinking the Middle to Later Stone Age transition in southern
1953 Africa - A perspective from the highveld of Eswatini', *Quaternary Science Reviews*, 286, p.
1954 107540. Available at: <https://doi.org/10.1016/j.quascirev.2022.107540>.
- 1955 Bartlein, P.J. *et al.* (2011) 'Pollen-based continental climate reconstructions at 6 and 21 ka: a
1956 global synthesis', *Climate Dynamics*, 37(3–4), pp. 775–802. Available at:
1957 <https://doi.org/10.1007/s00382-010-0904-1>.
- 1958 Beaumont, M.A., Zhang, W. and Balding, D.J. (2002) 'Approximate Bayesian computation in
1959 population genetics', *Genetics*, 162(4), pp. 2025–2035.
- 1960 Behr, A.A. *et al.* (2016) 'pong: fast analysis and visualization of latent clusters in population
1961 genetic data', *Bioinformatics*, 32(18), pp. 2817–2823. Available at:
1962 <https://doi.org/10.1093/bioinformatics/btw327>.
- 1963 Bergström, A. *et al.* (2021) 'Origins of modern human ancestry', *Nature*, 590(7845), pp. 229–
1964 237. Available at: <https://doi.org/10.1038/s41586-021-03244-5>.
- 1965 Beyer, R.M. *et al.* (2021) 'Climatic windows for human migration out of Africa in the past
1966 300,000 years', *Nature Communications*, 12(1), p. 4889. Available at:
1967 <https://doi.org/10.1038/s41467-021-24779-1>.
- 1968 Blum, M.G.B. and François, O. (2010) 'Non-linear regression models for Approximate
1969 Bayesian Computation', *Statistics and Computing*, 20(1), pp. 63–73. Available at:
1970 <https://doi.org/10.1007/s11222-009-9116-0>.
- 1971 Boitard, S. *et al.* (2016) 'Inferring Population Size History from Large Samples of Genome-
1972 Wide Molecular Data - An Approximate Bayesian Computation Approach', *PLOS Genetics*.
1973 Edited by M.A. Beaumont, 12(3), p. e1005877. Available at:
1974 <https://doi.org/10.1371/journal.pgen.1005877>.
- 1975 Bostoen, K. *et al.* (2015) 'Middle to Late Holocene Paleoclimatic Change and the Early
1976 Bantu Expansion in the Rain Forests of Western Central Africa', *Current Anthropology*,
1977 56(3), pp. 354–384. Available at: <https://doi.org/10.1086/681436>.
- 1978 Bowcock, A.M. *et al.* (1991) 'Drift, admixture, and selection in human evolution: a study with
1979 DNA polymorphisms.', *Proceedings of the National Academy of Sciences*, 88(3), pp. 839–
1980 843. Available at: <https://doi.org/10.1073/pnas.88.3.839>.
- 1981 Breton, G. *et al.* (2014) 'Lactase Persistence Alleles Reveal Partial East African Ancestry of
1982 Southern African Khoe Pastoralists', *Current Biology*, 24(8), pp. 852–858. Available at:
1983 <https://doi.org/10.1016/j.cub.2014.02.041>.

- 1984 Breton, G. *et al.* (2021) 'Comparison of sequencing data processing pipelines and
1985 application to underrepresented African human populations', *BMC Bioinformatics*, 22(1), p.
1986 488. Available at: <https://doi.org/10.1186/s12859-021-04407-x>.
- 1987 Breton, G., Fortes-Lima, C. and Schlebusch, C.M. (2021) 'Revisiting the demographic history
1988 of Central African populations from a genetic perspective', *Human Population Genetics and*
1989 *Genomics*, pp. 1–29. Available at: <https://doi.org/10.47248/hpgg2101010004>.
- 1990 Busby, G.B. *et al.* (2016) 'Admixture into and within sub-Saharan Africa', *eLife*, 5, p. e15266.
1991 Available at: <https://doi.org/10.7554/eLife.15266>.
- 1992 Chen, J. *et al.* (2020) 'Improved XGBoost model based on genetic algorithm', *International*
1993 *Journal of Computer Applications in Technology*, 62(3), p. 240. Available at:
1994 <https://doi.org/10.1504/IJCAT.2020.106571>.
- 1995 Choudhury, A. *et al.* (2017) 'Whole-genome sequencing for an enhanced understanding of
1996 genetic variation among South Africans', *Nature Communications*, 8(1), p. 2062. Available
1997 at: <https://doi.org/10.1038/s41467-017-00663-9>.
- 1998 Cornelissen, E. (2002) '[No title found]', *Journal of World Prehistory*, 16(3), pp. 197–235.
1999 Available at: <https://doi.org/10.1023/A:1020949501304>.
- 2000 Csilléry, K., François, O. and Blum, M.G.B. (2012) 'abc: an R package for approximate
2001 Bayesian computation (ABC): *R package: abc*', *Methods in Ecology and Evolution*, 3(3), pp.
2002 475–479. Available at: <https://doi.org/10.1111/j.2041-210X.2011.00179.x>.
- 2003 Cuadros-Espinoza, S. *et al.* (2022) 'The genomic signatures of natural selection in admixed
2004 human populations', *The American Journal of Human Genetics*, 109(4), pp. 710–726.
2005 Available at: <https://doi.org/10.1016/j.ajhg.2022.02.011>.
- 2006 Danecek, P. *et al.* (2011) 'The variant call format and VCFtools', *Bioinformatics*, 27(15), pp.
2007 2156–2158. Available at: <https://doi.org/10.1093/bioinformatics/btr330>.
- 2008 DePristo, M.A. *et al.* (2011) 'A framework for variation discovery and genotyping using next-
2009 generation DNA sequencing data', *Nature Genetics*, 43(5), pp. 491–498. Available at:
2010 <https://doi.org/10.1038/ng.806>.
- 2011 Estoup, A. *et al.* (2018) 'Model choice using Approximate Bayesian Computation and
2012 Random Forests: analyses based on model grouping to make inferences about the genetic
2013 history of Pygmy human populations.', *Journal de la Société Française de Statistique*. Tome
2014 159, pp. 167–190. Available at: http://www.numdam.org/item/JSFS_2018__159_3_167_0/.
- 2015 Excoffier, L. *et al.* (2013) 'Robust Demographic Inference from Genomic and SNP Data',
2016 *PLoS Genetics*. Edited by J.M. Akey, 9(10), p. e1003905. Available at:
2017 <https://doi.org/10.1371/journal.pgen.1003905>.
- 2018 Excoffier, L. and Foll, M. (2011) 'fastsimcoal: a continuous-time coalescent simulator of
2019 genomic diversity under arbitrarily complex evolutionary scenarios', *Bioinformatics*, 27(9),
2020 pp. 1332–1334. Available at: <https://doi.org/10.1093/bioinformatics/btr124>.
- 2021 Falush, D., Stephens, M. and Pritchard, J.K. (2003) 'Inference of population structure using
2022 multilocus genotype data: linked loci and correlated allele frequencies', *Genetics*, 164(4), pp.
2023 1567–1587.

- 2024 Fan, S. *et al.* (2023) 'Whole-genome sequencing reveals a complex African population
2025 demographic history and signatures of local adaptation', *Cell*, 186(5), pp. 923-939.e14.
2026 Available at: <https://doi.org/10.1016/j.cell.2023.01.042>.
- 2027 Fenner, J.N. (2005) 'Cross-cultural estimation of the human generation interval for use in
2028 genetics-based population divergence studies', *American Journal of Physical Anthropology*,
2029 128(2), pp. 415–423. Available at: <https://doi.org/10.1002/ajpa.20188>.
- 2030 Fortes-Lima, C. *et al.* (2022) 'Demographic and Selection Histories of Populations Across
2031 the Sahel/Savannah Belt', *Molecular Biology and Evolution*. Edited by E. Heyer, 39(10), p.
2032 msac209. Available at: <https://doi.org/10.1093/molbev/msac209>.
- 2033 Fortes-Lima, C.A. *et al.* (2021) 'Complex genetic admixture histories reconstructed with
2034 Approximate Bayesian Computation', *Molecular Ecology Resources*, 21(4), pp. 1098–1117.
2035 Available at: <https://doi.org/10.1111/1755-0998.13325>.
- 2036 Fortes-Lima, C.A. *et al.* (2024) 'The genetic legacy of the expansion of Bantu-speaking
2037 peoples in Africa', *Nature*, 625(7995), pp. 540–547. Available at:
2038 <https://doi.org/10.1038/s41586-023-06770-6>.
- 2039 Gascuel, O. (1997) 'BIONJ: an improved version of the NJ algorithm based on a simple
2040 model of sequence data', *Molecular Biology and Evolution*, 14(7), pp. 685–695. Available at:
2041 <https://doi.org/10.1093/oxfordjournals.molbev.a025808>.
- 2042 Gosling, W.D., Scerri, E.M.L. and Kaboth-Bahr, S. (2022) 'The climate and vegetation
2043 backdrop to hominin evolution in Africa', *Philosophical Transactions of the Royal Society B:*
2044 *Biological Sciences*, 377(1849), p. 20200483. Available at:
2045 <https://doi.org/10.1098/rstb.2020.0483>.
- 2046 Gravel, S. (2012) 'Population Genetics Models of Local Ancestry', *Genetics*, 191(2), pp.
2047 607–619. Available at: <https://doi.org/10.1534/genetics.112.139808>.
- 2048 Hamid, I. *et al.* (2021) 'Rapid adaptation to malaria facilitated by admixture in the human
2049 population of Cabo Verde', *eLife*, 10, p. e63177. Available at:
2050 <https://doi.org/10.7554/eLife.63177>.
- 2051 Harris, K. and Nielsen, R. (2016) 'The Genetic Cost of Neanderthal Introgression', *Genetics*,
2052 203(2), pp. 881–891. Available at: <https://doi.org/10.1534/genetics.116.186890>.
- 2053 Henn, B.M., Steele, T.E. and Weaver, T.D. (2018) 'Clarifying distinct models of modern
2054 human origins in Africa', *Current Opinion in Genetics & Development*, 53, pp. 148–156.
2055 Available at: <https://doi.org/10.1016/j.gde.2018.10.003>.
- 2056 Hewlett, B.S. (ed.) (2017) *Hunter-gatherers of the Congo Basin: cultures, histories and*
2057 *biology of African Pygmies*. First published by Transaction Publishers 2014. London New
2058 York: Routledge, Taylor and Francis Group.
- 2059 Hollfelder, N. *et al.* (2021) 'The deep population history in Africa', *Human Molecular*
2060 *Genetics*, 30(R1), pp. R2–R10. Available at: <https://doi.org/10.1093/hmg/ddab005>.
- 2061 Huang, X. *et al.* (2024) 'Harnessing deep learning for population genetic inference', *Nature*
2062 *Reviews Genetics*, 25(1), pp. 61–78. Available at: [https://doi.org/10.1038/s41576-023-](https://doi.org/10.1038/s41576-023-00636-3)
2063 00636-3.

- 2064 Hublin, J.-J. *et al.* (2017) 'New fossils from Jebel Irhoud, Morocco and the pan-African origin
2065 of *Homo sapiens*', *Nature*, 546(7657), pp. 289–292. Available at:
2066 <https://doi.org/10.1038/nature22336>.
- 2067 Jay, F., Boitard, S. and Austerlitz, F. (2019) 'An ABC Method for Whole-Genome Sequence
2068 Data: Inferring Paleolithic and Neolithic Human Expansions', *Molecular Biology and*
2069 *Evolution*. Edited by R. Hernandez, 36(7), pp. 1565–1579. Available at:
2070 <https://doi.org/10.1093/molbev/msz038>.
- 2071 Kamm, J. *et al.* (2020) 'Efficiently Inferring the Demographic History of Many Populations
2072 With Allele Count Data', *Journal of the American Statistical Association*, 115(531), pp. 1472–
2073 1487. Available at: <https://doi.org/10.1080/01621459.2019.1635482>.
- 2074 Lachance, J. *et al.* (2012) 'Evolutionary History and Adaptation from High-Coverage Whole-
2075 Genome Sequences of Diverse African Hunter-Gatherers', *Cell*, 150(3), pp. 457–469.
2076 Available at: <https://doi.org/10.1016/j.cell.2012.07.009>.
- 2077 Laurent, R. *et al.* (2023) 'A genetic and linguistic analysis of the admixture histories of the
2078 islands of Cabo Verde', *eLife*, 12, p. e79827. Available at:
2079 <https://doi.org/10.7554/eLife.79827>.
- 2080 Lawson, D.J., van Dorp, L. and Falush, D. (2018) 'A tutorial on how not to over-interpret
2081 STRUCTURE and ADMIXTURE bar plots', *Nature Communications*, 9(1), p. 3258. Available
2082 at: <https://doi.org/10.1038/s41467-018-05257-7>.
- 2083 Leung, F.H.F. *et al.* (2003) 'Tuning of the structure and parameters of a neural network using
2084 an improved genetic algorithm', *IEEE Transactions on Neural Networks*, 14(1), pp. 79–88.
2085 Available at: <https://doi.org/10.1109/TNN.2002.804317>.
- 2086 Lézine, A.-M. *et al.* (2019) 'A 90,000-year record of Afromontane forest responses to climate
2087 change', *Science*, 363(6423), pp. 177–181. Available at:
2088 <https://doi.org/10.1126/science.aav6821>.
- 2089 Li, H. and Durbin, R. (2009) 'Fast and accurate short read alignment with Burrows–Wheeler
2090 transform', *Bioinformatics*, 25(14), pp. 1754–1760. Available at:
2091 <https://doi.org/10.1093/bioinformatics/btp324>.
- 2092 Li, J.Z. *et al.* (2008) 'Worldwide Human Relationships Inferred from Genome-Wide Patterns
2093 of Variation', *Science*, 319(5866), pp. 1100–1104. Available at:
2094 <https://doi.org/10.1126/science.1153717>.
- 2095 Lipson, M. *et al.* (2022) 'Ancient DNA and deep population structure in sub-Saharan African
2096 foragers', *Nature*, 603(7900), pp. 290–296. Available at: <https://doi.org/10.1038/s41586-022-04430-9>.
- 2098 Long, J.C. (1991) 'The genetic structure of admixed populations.', *Genetics*, 127(2), pp.
2099 417–428. Available at: <https://doi.org/10.1093/genetics/127.2.417>.
- 2100 Lopez, M. *et al.* (2018) 'The demographic history and mutational load of African hunter-
2101 gatherers and farmers', *Nature Ecology & Evolution*, 2(4), pp. 721–730. Available at:
2102 <https://doi.org/10.1038/s41559-018-0496-4>.
- 2103 Lopez, M. *et al.* (2019) 'Genomic Evidence for Local Adaptation of Hunter-Gatherers to the
2104 African Rainforest', *Current Biology*, 29(17), pp. 2926–2935.e4. Available at:
2105 <https://doi.org/10.1016/j.cub.2019.07.013>.

- 2106 Lorente-Galdos, B. *et al.* (2019) 'Whole-genome sequence analysis of a Pan African set of
2107 samples reveals archaic gene flow from an extinct basal population of modern humans into
2108 sub-Saharan populations', *Genome Biology*, 20(1), p. 77. Available at:
2109 <https://doi.org/10.1186/s13059-019-1684-5>.
- 2110 Lucas-Sánchez, M., Serradell, J.M. and Comas, D. (2021) 'Population history of North Africa
2111 based on modern and ancient genomes', *Human Molecular Genetics*, 30(R1), pp. R17–R23.
2112 Available at: <https://doi.org/10.1093/hmg/ddaa261>.
- 2113 Mallick, S. *et al.* (2016) 'The Simons Genome Diversity Project: 300 genomes from 142
2114 diverse populations', *Nature*, 538(7624), pp. 201–206. Available at:
2115 <https://doi.org/10.1038/nature18964>.
- 2116 Mazet, O. *et al.* (2016) 'On the importance of being structured: instantaneous coalescence
2117 rates and human evolution—lessons for ancestral population size inference?', *Heredity*,
2118 116(4), pp. 362–371. Available at: <https://doi.org/10.1038/hdy.2015.104>.
- 2119 McKenna, A. *et al.* (2010) 'The Genome Analysis Toolkit: A MapReduce framework for
2120 analyzing next-generation DNA sequencing data', *Genome Research*, 20(9), pp. 1297–1303.
2121 Available at: <https://doi.org/10.1101/gr.107524.110>.
- 2122 Mesfin, I., Oslisly, R. and Forestier, H. (2021) 'Technological analysis of the quartz industry
2123 of Maboué 5 – Layer 3 (Lopé National Park, Gabon): Implications for the Late Stone Age
2124 emergence in western Central Africa', *Journal of Archaeological Science: Reports*, 39, p.
2125 103130. Available at: <https://doi.org/10.1016/j.jasrep.2021.103130>.
- 2126 Meyer, M. *et al.* (2012) 'A High-Coverage Genome Sequence from an Archaic Denisovan
2127 Individual', *Science*, 338(6104), pp. 222–226. Available at:
2128 <https://doi.org/10.1126/science.1224344>.
- 2129 Miller, S.A., Dykes, D.D. and Polesky, H.F. (1988) 'A simple salting out procedure for
2130 extracting DNA from human nucleated cells', *Nucleic Acids Research*, 16(3), pp. 1215–1215.
2131 Available at: <https://doi.org/10.1093/nar/16.3.1215>.
- 2132 Mills, R.E. *et al.* (2011) 'Natural genetic variation caused by small insertions and deletions in
2133 the human genome', *Genome Research*, 21(6), pp. 830–839. Available at:
2134 <https://doi.org/10.1101/gr.115907.110>.
- 2135 Mooney, J.A. *et al.* (2018) 'Understanding the Hidden Complexity of Latin American
2136 Population Isolates', *The American Journal of Human Genetics*, 103(5), pp. 707–726.
2137 Available at: <https://doi.org/10.1016/j.ajhg.2018.09.013>.
- 2138 Mooney, J.A. *et al.* (2023) 'On the number of genealogical ancestors tracing to the source
2139 groups of an admixed population', *GENETICS*. Edited by J. Novembre, 224(3), p. iyad079.
2140 Available at: <https://doi.org/10.1093/genetics/iyad079>.
- 2141 Murtagh, F. (1991) 'Multilayer perceptrons for classification and regression',
2142 *Neurocomputing*, 2(5–6), pp. 183–197. Available at: [https://doi.org/10.1016/0925-2312\(91\)90023-5](https://doi.org/10.1016/0925-2312(91)90023-5).
- 2144 Nei, M. (1978) 'Estimation of average heterozygosity and genetic distance from a small
2145 number of individuals', *Genetics*, 89(3), pp. 583–590. Available at:
2146 <https://doi.org/10.1093/genetics/89.3.583>.

- 2147 Okonechnikov, K., Conesa, A. and García-Alcalde, F. (2016) 'Qualimap 2: advanced multi-
2148 sample quality control for high-throughput sequencing data', *Bioinformatics*, 32(2), pp. 292–
2149 294. Available at: <https://doi.org/10.1093/bioinformatics/btv566>.
- 2150 Patin, E. *et al.* (2009) 'Inferring the Demographic History of African Farmers and Pygmy
2151 Hunter–Gatherers Using a Multilocus Resequencing Data Set', *PLoS Genetics*. Edited by A.
2152 Di Rienzo, 5(4), p. e1000448. Available at: <https://doi.org/10.1371/journal.pgen.1000448>.
- 2153 Patin, E. *et al.* (2014) 'The impact of agricultural emergence on the genetic history of African
2154 rainforest hunter-gatherers and agriculturalists', *Nature Communications*, 5(1), p. 3163.
2155 Available at: <https://doi.org/10.1038/ncomms4163>.
- 2156 Patin, E. *et al.* (2017) 'Dispersals and genetic adaptation of Bantu-speaking populations in
2157 Africa and North America', *Science*, 356(6337), pp. 543–546. Available at:
2158 <https://doi.org/10.1126/science.aal1988>.
- 2159 Pemberton, T.J. *et al.* (2012) 'Genomic Patterns of Homozygosity in Worldwide Human
2160 Populations', *The American Journal of Human Genetics*, 91(2), pp. 275–292. Available at:
2161 <https://doi.org/10.1016/j.ajhg.2012.06.014>.
- 2162 Perry, G.H. *et al.* (2014) 'Adaptive, convergent origins of the pygmy phenotype in African
2163 rainforest hunter-gatherers', *Proceedings of the National Academy of Sciences*, 111(35).
2164 Available at: <https://doi.org/10.1073/pnas.1402875111>.
- 2165 Perry, G.H. and Verdu, P. (2017) 'Genomic perspectives on the history and evolutionary
2166 ecology of tropical rainforest occupation by humans', *Quaternary International*, 448, pp. 150–
2167 157. Available at: <https://doi.org/10.1016/j.quaint.2016.04.038>.
- 2168 Peter, B.M. (2022) 'A geometric relationship of F_2 , F_3 and F_4 -statistics with principal
2169 component analysis', *Philosophical Transactions of the Royal Society B: Biological
2170 Sciences*, 377(1852), p. 20200413. Available at: <https://doi.org/10.1098/rstb.2020.0413>.
- 2171 Pfennig, A. *et al.* (2023) 'Evolutionary Genetics and Admixture in African Populations',
2172 *Genome Biology and Evolution*. Edited by A. Eyre-Walker, 15(4), p. evad054. Available at:
2173 <https://doi.org/10.1093/gbe/evad054>.
- 2174 Phillipson, D.W. (2005) *African Archaeology*. 3rd edn. Cambridge University Press. Available
2175 at: <https://doi.org/10.1017/CBO9780511800313>.
- 2176 Pierron, D. *et al.* (2017) 'Genomic landscape of human diversity across Madagascar',
2177 *Proceedings of the National Academy of Sciences*, 114(32). Available at:
2178 <https://doi.org/10.1073/pnas.1704906114>.
- 2179 Pritchard, J.K. *et al.* (1999) 'Population growth of human Y chromosomes: a study of Y
2180 chromosome microsatellites', *Molecular Biology and Evolution*, 16(12), pp. 1791–1798.
2181 Available at: <https://doi.org/10.1093/oxfordjournals.molbev.a026091>.
- 2182 Pritchard, J.K., Stephens, M. and Donnelly, P. (2000) 'Inference of population structure using
2183 multilocus genotype data', *Genetics*, 155(2), pp. 945–959.
- 2184 Prüfer, K. *et al.* (2014) 'The complete genome sequence of a Neanderthal from the Altai
2185 Mountains', *Nature*, 505(7481), pp. 43–49. Available at: <https://doi.org/10.1038/nature12886>.
- 2186 Pudlo, P. *et al.* (2016) 'Reliable ABC model choice via random forests', *Bioinformatics*, 32(6),
2187 pp. 859–866. Available at: <https://doi.org/10.1093/bioinformatics/btv684>.

- 2188 Purcell, S. *et al.* (2007) 'PLINK: A Tool Set for Whole-Genome Association and Population-
2189 Based Linkage Analyses', *The American Journal of Human Genetics*, 81(3), pp. 559–575.
2190 Available at: <https://doi.org/10.1086/519795>.
- 2191 R Core Team (2015) 'R: A Language and Environment for Statistical Computing'. Vienna,
2192 Austria: R Foundation for Statistical Computing. Available at: <https://www.R-project.org>.
- 2193 Ragsdale, A.P. *et al.* (2023) 'A weakly structured stem for human origins in Africa', *Nature*,
2194 617(7962), pp. 755–763. Available at: <https://doi.org/10.1038/s41586-023-06055-y>.
- 2195 Ragsdale, A.P. and Gravel, S. (2019) 'Models of archaic admixture and recent history from
2196 two-locus statistics', *PLOS Genetics*. Edited by J.M. Akey, 15(6), p. e1008204. Available at:
2197 <https://doi.org/10.1371/journal.pgen.1008204>.
- 2198 Rasmussen, M. *et al.* (2014) 'The genome of a Late Pleistocene human from a Clovis burial
2199 site in western Montana', *Nature*, 506(7487), pp. 225–229. Available at:
2200 <https://doi.org/10.1038/nature13025>.
- 2201 Raynal, L. *et al.* (2019) 'ABC random forests for Bayesian parameter inference',
2202 *Bioinformatics*. Edited by O. Stegle, 35(10), pp. 1720–1728. Available at:
2203 <https://doi.org/10.1093/bioinformatics/bty867>.
- 2204 Richter, D. *et al.* (2017) 'The age of the hominin fossils from Jebel Irhoud, Morocco, and the
2205 origins of the Middle Stone Age', *Nature*, 546(7657), pp. 293–296. Available at:
2206 <https://doi.org/10.1038/nature22335>.
- 2207 Robert, C.P., Mengersen, K. and Chen, C. (2010) 'Model choice versus model criticism',
2208 *Proceedings of the National Academy of Sciences*, 107(3), pp. E5–E5. Available at:
2209 <https://doi.org/10.1073/pnas.0911260107>.
- 2210 Rosenberg, N.A. (2002) 'Genetic Structure of Human Populations', *Science*, 298(5602), pp.
2211 2381–2385. Available at: <https://doi.org/10.1126/science.1078311>.
- 2212 Saitou, N. and Nei, M. (1987) 'The neighbor-joining method: a new method for reconstructing
2213 phylogenetic trees.', *Molecular Biology and Evolution* [Preprint]. Available at:
2214 <https://doi.org/10.1093/oxfordjournals.molbev.a040454>.
- 2215 Scerri, E.M.L. *et al.* (2018) 'Did Our Species Evolve in Subdivided Populations across Africa,
2216 and Why Does It Matter?', *Trends in Ecology & Evolution*, 33(8), pp. 582–594. Available at:
2217 <https://doi.org/10.1016/j.tree.2018.05.005>.
- 2218 Schlebusch, C. (2010) 'Issues raised by use of ethnic-group names in genome study',
2219 *Nature*, 464(7288), pp. 487–487. Available at: <https://doi.org/10.1038/464487a>.
- 2220 Schlebusch, C.M. *et al.* (2012) 'Genomic Variation in Seven Khoe-San Groups Reveals
2221 Adaptation and Complex African History', *Science*, 338(6105), pp. 374–379. Available at:
2222 <https://doi.org/10.1126/science.1227721>.
- 2223 Schlebusch, C.M. *et al.* (2020) 'Khoe-San Genomes Reveal Unique Variation and Confirm
2224 the Deepest Population Divergence in Homo sapiens', *Molecular Biology and Evolution*.
2225 Edited by C. Mulligan, 37(10), pp. 2944–2954. Available at:
2226 <https://doi.org/10.1093/molbev/msaa140>.

- 2227 Schlebusch, C.M. and Jakobsson, M. (2018) 'Tales of Human Migration, Admixture, and
2228 Selection in Africa', *Annual Review of Genomics and Human Genetics*, 19(1), pp. 405–428.
2229 Available at: <https://doi.org/10.1146/annurev-genom-083117-021759>.
- 2230 Seidensticker, D. *et al.* (2021) 'Population collapse in Congo rainforest from 400 CE urges
2231 reassessment of the Bantu Expansion', *Science Advances*, 7(7), p. eabd8352. Available at:
2232 <https://doi.org/10.1126/sciadv.abd8352>.
- 2233 Semo, A. *et al.* (2020) 'Along the Indian Ocean Coast: Genomic Variation in Mozambique
2234 Provides New Insights into the Bantu Expansion', *Molecular Biology and Evolution*. Edited by
2235 D. Falush, 37(2), pp. 406–416. Available at: <https://doi.org/10.1093/molbev/msz224>.
- 2236 Sengupta, D. *et al.* (2021) 'Genetic substructure and complex demographic history of South
2237 African Bantu speakers', *Nature Communications*, 12(1), p. 2080. Available at:
2238 <https://doi.org/10.1038/s41467-021-22207-y>.
- 2239 Skoglund, P. *et al.* (2017) 'Reconstructing Prehistoric African Population Structure', *Cell*,
2240 171(1), pp. 59–71.e21. Available at: <https://doi.org/10.1016/j.cell.2017.08.049>.
- 2241 Skoglund, P. and Mathieson, I. (2018) 'Ancient Genomics of Modern Humans: The First
2242 Decade', *Annual Review of Genomics and Human Genetics*, 19(1), pp. 381–404. Available
2243 at: <https://doi.org/10.1146/annurev-genom-083117-021749>.
- 2244 Stringer, C. (2002) 'Modern human origins: progress and prospects', *Philosophical
2245 Transactions of the Royal Society of London. Series B: Biological Sciences*, 357(1420), pp.
2246 563–579. Available at: <https://doi.org/10.1098/rstb.2001.1057>.
- 2247 Szpiech, Z.A. *et al.* (2019) 'Ancestry-Dependent Enrichment of Deleterious Homozygotes in
2248 Runs of Homozygosity', *The American Journal of Human Genetics*, 105(4), pp. 747–762.
2249 Available at: <https://doi.org/10.1016/j.ajhg.2019.08.011>.
- 2250 Tavaré, S. *et al.* (1997) 'Inferring coalescence times from DNA sequence data', *Genetics*,
2251 145(2), pp. 505–518.
- 2252 The 1000 Genomes Project Consortium *et al.* (2015) 'A global reference for human genetic
2253 variation', *Nature*, 526(7571), pp. 68–74. Available at: <https://doi.org/10.1038/nature15393>.
- 2254 Thomas, D.S.G. *et al.* (2022) 'Lacustrine geoarchaeology in the central Kalahari:
2255 Implications for Middle Stone Age behaviour and adaptation in dryland conditions',
2256 *Quaternary Science Reviews*, 297, p. 107826. Available at:
2257 <https://doi.org/10.1016/j.quascirev.2022.107826>.
- 2258 Tishkoff, S.A. *et al.* (2009) 'The Genetic Structure and History of Africans and African
2259 Americans', *Science*, 324(5930), pp. 1035–1044. Available at:
2260 <https://doi.org/10.1126/science.1172257>.
- 2261 Van Der Auwera, G.A. *et al.* (2013) 'From FastQ Data to High-Confidence Variant Calls: The
2262 Genome Analysis Toolkit Best Practices Pipeline', *Current Protocols in Bioinformatics*, 43(1).
2263 Available at: <https://doi.org/10.1002/0471250953.bi1110s43>.
- 2264 Verdu, P. *et al.* (2009) 'Origins and Genetic Diversity of Pygmy Hunter-Gatherers from
2265 Western Central Africa', *Current Biology*, 19(4), pp. 312–318. Available at:
2266 <https://doi.org/10.1016/j.cub.2008.12.049>.

- 2267 Verdu, P. *et al.* (2013) 'Sociocultural Behavior, Sex-Biased Admixture, and Effective
2268 Population Sizes in Central African Pygmies and Non-Pygmies', *Molecular Biology and*
2269 *Evolution*, 30(4), pp. 918–937. Available at: <https://doi.org/10.1093/molbev/mss328>.
- 2270 Verdu, P. and Rosenberg, N.A. (2011) 'A General Mechanistic Model for Admixture Histories
2271 of Hybrid Populations', *Genetics*, 189(4), pp. 1413–1426. Available at:
2272 <https://doi.org/10.1534/genetics.111.132787>.
- 2273 Vicente, M. and Schlebusch, C.M. (2020) 'African population history: an ancient DNA
2274 perspective', *Current Opinion in Genetics & Development*, 62, pp. 8–15. Available at:
2275 <https://doi.org/10.1016/j.gde.2020.05.008>.
- 2276 Wang, H., Czerminski, R. and Jamieson, A.C. (2021) 'Neural Networks and Deep Learning',
2277 in M. Einhorn *et al.* (eds) *The Machine Age of Customer Insight*. Emerald Publishing Limited,
2278 pp. 91–101. Available at: <https://doi.org/10.1108/978-1-83909-694-520211010>.
- 2279 Yelmen, B. and Jay, F. (2023) 'An Overview of Deep Generative Models in Functional and
2280 Evolutionary Genomics', *Annual Review of Biomedical Data Science*, 6(1), pp. 173–189.
2281 Available at: <https://doi.org/10.1146/annurev-biodatasci-020722-115651>.
- 2282 Ziegler, M. *et al.* (2013) 'Development of Middle Stone Age innovation linked to rapid climate
2283 change', *Nature Communications*, 4(1), p. 1905. Available at:
2284 <https://doi.org/10.1038/ncomms2897>.
- 2285
2286