

On the limits of fitting complex models of population history to genetic data

Robert Maier^{1,*}, Pavel Flegontov^{1,2,*}, Olga Flegontova², Piya Changmai² and David Reich^{1,3,4,5,§}

¹ *Department of Human Evolutionary Biology, Harvard University, Cambridge, MA, USA*

² *Department of Biology and Ecology, Faculty of Science, University of Ostrava, Ostrava, Czechia*

³ *Howard Hughes Medical Institute, Harvard Medical School, Boston, MA, USA*

⁴ *Broad Institute of Harvard and MIT, Cambridge, MA, USA*

⁵ *Department of Genetics, Harvard Medical School, Boston, MA, USA*

* These authors contributed equally

§ corresponding author's email address: rmaier@broadinstitute.org, pavel.flegontov@osu.cz, reich@genetics.med.harvard.edu

1 **Abstract**

2

3 Our understanding of human population history in deep time has been assisted by fitting “admixture
4 graphs” to data: models that specify the ordering of population splits and mixtures which is the only
5 information needed to capture the patterns of allele frequency correlation among populations. Not
6 needing to specify population size changes, split times, or whether admixture events were sudden or
7 drawn out simplifies the space of models that need to be searched. However, the space of possible
8 admixture graphs relating populations is vast and cannot be sampled fully, and thus most published
9 studies have identified fitting admixture graphs through a manual process driven by prior
10 hypotheses, leaving the vast majority of alternative models unexplored. Here, we develop a method
11 for systematically searching the space of all admixture graphs that can incorporate non-genetic
12 information in the form of topology constraints. We implement this *findGraphs* tool within a
13 software package, *ADMIXTOOLS 2*, which is a reimplement of the *ADMIXTOOLS* software with
14 new features and large performance gains. We apply this methodology to identify alternative
15 models to admixture graphs that played key roles in eight published studies and find that graphs
16 modeling more than six populations and two or three admixture events are often not unique, with
17 many alternative models fitting nominally or significantly better than the published one. Our results
18 suggest that strong claims about population history from admixture graphs should only be made
19 when all well-fitting and temporally plausible models share common topological features. Our re-
20 evaluation of published data also provides insight into the population histories of humans, dogs, and
21 horses, identifying features that are stable across the models we explored, as well as scenarios of
22 populations relationships that differ in important ways from models that have been highlighted in
23 the literature, that fit the allele frequency correlation data, and that are not obviously wrong.

24 Introduction

25

26 Admixture graph models provide a powerful intellectual framework for describing the relationships
27 among populations that allows not only branching of populations from a common ancestor but also
28 mixture events. Admixture graphs can precisely summarize important features of population history
29 without requiring specification of all parameters such as population sizes, split times, mixture times,
30 and distinguishing between sudden splits or drawn-out separations. All these parameters describe
31 important features of demographic history and are considered by several tools for fitting
32 demographic models (Gutenkunst et al. 2009, Gronau et al. 2011, Schiffels et al. 2016, Flegontov et
33 al. 2019, Kamm et al. 2019, Rogers 2019, Hubisz et al. 2020). However, the fact that it is possible to
34 first infer important aspects of the topology (admixture graphs), and then fit these additional
35 parameters, simplifies demographic inference (Patterson et al. 2012, Pickrell and Pritchard 2012,
36 Lipson et al. 2013, Leppälä et al. 2017, Lipson 2020, Molloy et al. 2021, Yan et al. 2021). Admixture
37 graphs thus serve both as conceptual frameworks that allow us to think about the relationships of
38 populations deep in time, and as mathematical models we can fit to genetic data.

39

40 A challenge for fitting admixture graph models is that they are often not uniquely constrained by the
41 data, with many providing equally good fits to the f_2 -, f_3 -, and f_4 -statistics used to constrain them
42 within the limits of statistical resolution. Previously published methods for finding fitting admixture
43 graphs were not well-equipped to handle the large range of equally well-fitting models for three
44 reasons: (1) They did not reliably provide information on whether there is a uniquely fitting
45 parsimonious model or alternatively whether there are many models that fit equally well to within
46 the limits of statistical resolution, (2) they did not provide formal goodness-of-fit tests, and related
47 to this, (3) they not provide tests for whether the difference between the fits of any two models is
48 statistically significant. As a consequence and as we demonstrate in what follows, many published
49 admixture graph models have been interpreted as providing more confidence than is merited about
50 the extent to which genetic data allows us to disentangle ancestral relationships.

51

52 There are two main approaches to studying demographic history with admixture graphs.

53

54 The first approach is to identify admixture graphs automatically, either without human intervention
55 or with guidance. It is possible in theory to exhaustively test all possible graphs for a given set of
56 populations and pre-specified number of admixture events, as implemented, for example, in the
57 *admixturegraph* R package (Leppälä et al. 2017). An exhaustive approach can provide a complete
58 and unbiased view on the kinds of models that are consistent with the data for a specified level of
59 parsimony (total number of admixture events allowed in the graph), but this approach is limited to
60 small graphs (typically up to 6 groups, 2 admixture events) due to the rapid increase in the number
61 of possible admixture graphs as the number of populations and admixture events grows. In addition,
62 as we show in our discussion of case studies, the assumption of parsimony that needs to be made in
63 order to use an exhaustive approach can lead to misleading models of population history because
64 not including additional populations can blind users to additional mixture events that occurred (and
65 whose existence is revealed by examining data from additional populations). Specifically, models
66 with additional admixture events that are qualitatively profoundly different to the best fitting
67 parsimonious graph and that capture the true history, will sometimes be completely missed when
68 applying the parsimony assumption. Alternatively, the programs *TreeMix* (Pickrell and Pritchard
69 2012, Molloy et al. 2021), *MixMapper* (Lipson et al. 2013), and *miqoGraph* (Yan et al. 2021) all
70 address the problem of how to rapidly explore the vast space of admixture graphs relating a set of
71 populations by applying algorithmic ideas or heuristics; all of these methods speed up model search
72 by orders of magnitude. The new method we introduce here, *findGraphs*, belongs to this class of
73 algorithms. Algorithmic innovations in *findGraphs* enable us to get around some of the limitations
74 associated with parsimony assumptions by more efficiently exploring a larger proportion of plausible
75 of admixture graph space, and by increasing the speed of evaluation of individual graphs.

76

77 The second approach to fitting admixture graphs is to manually build them up by grafting additional
78 populations onto simpler smaller graphs that fit the data. This approach involves stepwise addition
79 of populations in an order that is chosen based on the best judgment of the user, and for each newly
80 added population involves adding admixture events or tweaks in the graph until a fit is obtained; the
81 user then moves on to adding the next population (see Reich et al. Nature 2009, Reich et al. AJHG
82 2011, Reich et al. Nature 2012, Lazaridis et al. 2014, Seguin-Orlando et al. 2014, Fu et al. 2016,
83 Skoglund et al. 2016, Yang et al. 2017, McColl et al. 2018, Moreno-Mayar et al. 2018, Tambets et al.
84 2018, van de Loosdrecht et al. 2018, Flegontov et al. 2019, Sikora et al. 2019, Wang et al. 2019,
85 Lipson et al. 2020, Shinde et al. 2019, Yang et al. 2020, Hajdinjak et al. 2021, Wang et al. 2021 for
86 examples). The program *qpGraph* in the *ADMIXTOOLS* package (Patterson et al. 2012) has been the
87 most common computational method used for testing fits of individual admixture graphs. Most
88 published admixture graphs have been constructed manually, often acknowledging the existence of
89 alternative models by presenting plausible models side-by-side, and this approach has been the
90 basis for many claims about population history (Lazaridis et al. 2014, Yang et al. 2017, Posth et al.
91 2018, Sikora et al. 2019, Shinde et al. 2019, Bergström et al. 2020, Lipson et al. 2020, Hajdinjak et al.
92 2021, Wang et al. 2021). A strength of this approach is that it takes advantage of human judgment
93 and outside knowledge about what graphs best fit the history of the human populations being
94 analyzed. This external information is powerful as it can incorporate non-genetic evidence such as
95 geographic plausibility and temporal ordering of populations or linguistic similarity, or other genetic
96 data such as estimates of population split times or shared Y chromosomes or rejection of proposed
97 scenarios based on joint analysis of much larger numbers of populations than can reasonably be
98 analyzed within a single admixture graph. Thus, while manual approaches explore many orders of
99 magnitude fewer topologies than automatic approaches often do, they still may provide inferences
100 about population history that are more useful than those provided by automatic approaches. These
101 methods' strength is also their weakness: by relying on intuition, following a manual approach has
102 the potential to validate the biases users have as to what types of histories are most plausible (these
103 may be the only types of histories that will be carefully explored). This can blind users to surprises:
104 to profoundly different topologies that may correspond more closely to the true history, and we
105 discuss examples of this in the Results section.

106
107 The *findGraphs* method combines the advantages of automated and manual topology exploration by
108 allowing users to encode various sources of information as constraints on the space of admixture
109 graphs, which is then explored automatically.

110
111 Regardless of the approach used to search through the space of possibly fitting admixture graphs, a
112 challenge in the effort to find a uniquely well-fitting admixture graph is that it is difficult to quantify
113 the absolute quality of the fit of a model, as well as the relative quality of the fits of multiple models.
114 Performance gains relative to the original implementation of *qpGraph* allow us to address this
115 problem by obtaining bootstrap confidence intervals and *p*-values for estimated parameters of
116 single models, as well as for the difference in fit quality of two models. Existing methods for
117 comparing admixture graph models (for example based on Akaike information criterion (AIC) or
118 Bayesian information criterion (BIC), see Flegontov et al. 2019, Shinde et al. 2019) do not take the
119 variability across SNPs into account appropriately since they do not rely on dataset resampling
120 approaches such as bootstrap, and thus they tend to overestimate our ability to differentiate
121 competing models. In combination with the previously described approach to automating the search
122 of admixture graphs, this leads to a situation where we are able to find and test a large number of
123 models, many of which fit equally well despite often having very different topological features.

124
125 The methods for automated graph topology inference and model comparison relying on bootstrap
126 resampling described above are implemented in *ADMIXTOOLS 2*, a comprehensive platform for
127 learning about population history from *f*-statistics. It is built to provide a stand-alone workspace for
128 research in this area and is implemented as an R package. For all computations, *ADMIXTOOLS 2*
129 exhibits large speedups relative to previously published platforms for *f*-statistic analysis (e.g.
130 *popstats* and *ADMIXTOOLS* version 6.0 which we call "Classic *ADMIXTOOLS*" in what follows to

131 distinguish it from updated *ADMIXTOOLS* version 7.0.2 which implements some of the speed-up
132 ideas also implemented in *ADMIXTOOLS 2*). This is achieved by deploying a series of algorithmic
133 improvements, most notably storage of pre-computed f -statistics in random access memory, which
134 avoids having to rely on reading in extremely large genotype matrices to perform most
135 computations. In addition to the new algorithmic ideas allowing efficient searching through the
136 space of admixture graphs and comparing the fits of two admixture graphs, *ADMIXTOOLS 2* also
137 provides a solution to the question of which parameters of an admixture graph are identifiable in
138 the limit of infinite data. The most important algorithmic improvements presented in *ADMIXTOOLS*
139 *2* and its philosophical differences relative to classic *ADMIXTOOLS* are described in the next section,
140 while the full methodological details are presented in the Methods and Supplement.

141

142

143 **Results**

144

145 ***New ideas implemented in ADMIXTOOLS 2***

146

147 We present a new implementation of the popular *ADMIXTOOLS* software (called “Classic
148 *ADMIXTOOLS*”) (Patterson et al. 2012, Haak et al. 2015). Our implementation (*ADMIXTOOLS 2*, see
149 documentation at <https://uqrmaie1.github.io/admixtools>) enhances performance by greatly
150 reducing runtime and memory requirements across a wide range of different methods, relative to
151 Classic *ADMIXTOOLS* (**Figure 1a**). We note that some of these improvements have now been
152 implemented in version 7.0.2 of *ADMIXTOOLS*. The present study focuses not on the performance
153 differences between Classic *ADMIXTOOLS* and *ADMIXTOOLS 2*, but on the description of new ideas
154 implemented in one or both of these tools.

155

156

157 *Computation and use of f-statistics*

158

159 A key idea that facilitates the performance increases shared by *ADMIXTOOLS 2* and *ADMIXTOOLS v.*
160 7.0.2 is that any f -statistic (which form the basis of almost all *ADMIXTOOLS* programs as well as
161 other toolkits for studying population history such as *popstats*) can be computed from a small
162 number of f_2 -statistics. For most f -statistic-based analyses (for example *qpWave*, *qpAdm*, and
163 *qpGraph*; **Figure 1c**), the time required to process f -statistics is trivial compared to the time required
164 to compute f -statistics from genotype data. These f_2 -statistics can be stored and re-used to
165 compute f_3 - and f_4 -statistics, thus reducing the size of the input data, runtime, and memory
166 requirements by orders of magnitude (**Figure 1a, 1d**).

167

168 Using precomputed f_2 -statistics is not always the best solution. In data sets with large amounts of
169 missing data, computing f_3 - and f_4 -statistics from pre-computed f_2 -statistics may introduce bias. In
170 this case, it is necessary to compute f_3 - and f_4 -statistics directly, using different SNPs for each f -
171 statistic (all available SNPs in each population triplet or quadruplet). However, even without the use
172 of pre-computed f_2 -statistics, *ADMIXTOOLS 2* often achieves large performance gains (**Figure 1a**).

173

174 The program *qpfstats* in Classic *ADMIXTOOLS* implements an idea which strikes a balance between
175 these two extremes. It increases the accuracy of estimation of f -statistics by using a regression
176 approach to jointly estimate the values of all f_2 -, f_3 - and f_4 -statistics relating a set of populations,
177 taking advantage of the algebraic relationships of the expected values of these statistics. This
178 approach makes it possible to obtain more precise estimates of the values of these statistics than
179 can be obtained by inferring them only using SNPs that are covered in each of the populations being
180 compared. This feature is available in *ADMIXTOOLS 2* through the *qpfstats* option in the *extract_f2*
181 function.

182

183 Another improvement introduced in *ADMIXTOOLS 2* relates to accurate evaluation of the match
184 between observed and expected f_3 -statistics when fitting admixture graphs where at least one

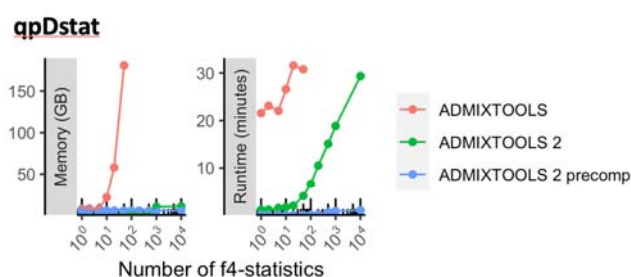
185 population is represented by a single individual with genotypes derived from randomly selected
186 sequencing reads (“pseudo-diploid” or “pseudo-haploid” data). f_3 -statistic computations need to be
187 modified when analyzing pseudo-diploid data, because heterozygosity cannot be computed using
188 comparisons of sequences within the same individual; however, computation of heterozygosity is
189 essential to calculate “admixture” f_3 -statistics where negative values provide proof of the mixed
190 nature of the target population. When a population is represented by multiple individuals, unbiased
191 estimation of admixture f_3 -statistics can be carried out even for pseudo-diploid data by analyzing
192 positions covered by sequences from at least two individuals and only computing variation rates
193 across individuals. This approach is implemented in Classic *ADMIXTOOLS* with the “inbreed: YES”
194 option. However, no admixture f_3 -statistic can be computed when the “inbreed: YES” option is
195 turned on and the target population is represented by a single individual (as no variation across
196 individuals within a population can be detected in this case). Classic *ADMIXTOOLS* deals with this
197 case by failing to run if any population in an analysis is represented by a single individual and the
198 “inbreed: YES” option is turned on. Because the datasets from all the admixture graphs revisited
199 here included at least one population represented by a single individual (**Table S1**), the “inbreed:
200 YES” option could not be used in the original studies (the program failed with this option, by design).
201 Thus, admixture graph fitting in those studies relied on the incorrect algorithm for calculating f_3 -
202 statistics (except for Librado *et al.* 2021, who used *TreeMix* instead of *qpGraph*) and, as a result,
203 some f_3 -statistics that are negative and could provide important constraints for admixture graph
204 fitting were evaluated as positive. These concerns are relevant for the Shinde *et al.* (2019), Lipson *et*
205 *al.* (2020), and Wang *et al.* (2021) studies we revisit below (see Table S1 for datasets where negative
206 f_3 -statistics were encountered). To be able to detect negative f_3 -statistics and thus take advantage of
207 their power for constraining the space of possibly fitting historical models, in *ADMIXTOOLS 2* we
208 introduced an option which makes it possible to compute negative f_3 -statistics on pseudo-diploid
209 data, at a cost of removing sites with only one chromosome genotyped in any population that is
210 represented by at least two individuals (so that it is possible in theory to compute heterozygosity in
211 these populations). Admixture f_3 -statistics continue to be incorrectly computed using *ADMIXTOOLS 2*
212 for targets that are singleton populations represented by pseudo-diploid data, as there is no
213 avoiding this particular problem. See Methods for a description of the new algorithm for calculating
214 f_3 -statistics.

215

216

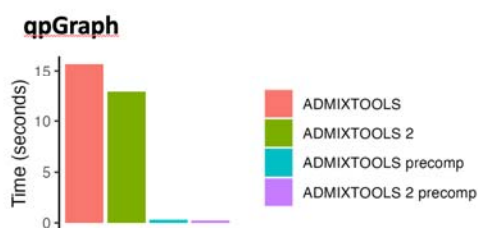
217 **Figure 1**

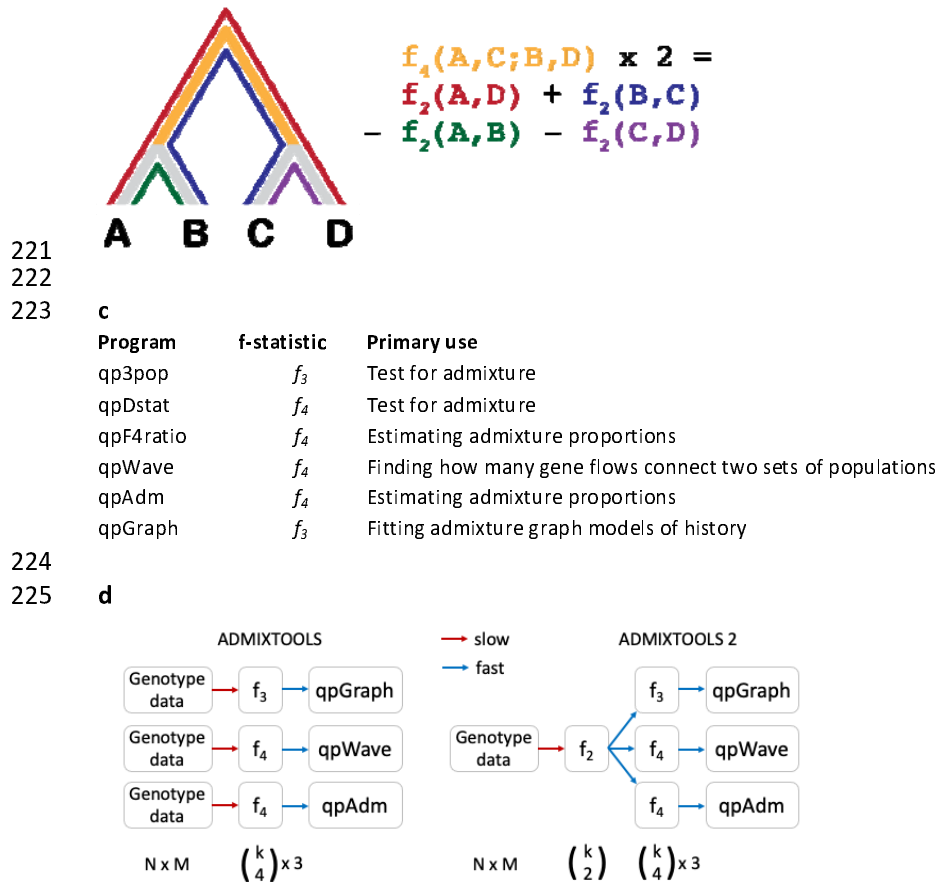
218 **a**



219

220 **b**





226

227 **Figure 1:**

228 **a** Performance comparison of f -statistic computation and admixture graph fitting. Top: Memory usage and
229 runtime for computing f -statistics using (1) the *qpDstat* program in *ADMIXTOOLS* v7.0.2 released in 06/2021
230 (2) the f function in *ADMIXTOOLS* 2 without precomputing f -statistics, and (3) the f function in
231 *ADMIXTOOLS* 2 with precomputed f -statistics. (1) and (2) give identical results, whereas (3) only gives
232 identical results in the absence of missing data, which limits usefulness beyond a moderate number of
233 populations. Bottom: Runtime comparison of *qpGraph* with and without precomputed f -statistics.

234 **b** Illustration of f_3 - and f_4 -statistics. f_3 measures the amount of drift separating any two populations, while
235 f_4 measures the amount of drift shared between two population pairs. Every f_4 -statistic is a linear combination of
236 four f_2 -statistics.

237 **c** Overview of the major *ADMIXTOOLS* programs, their primary use cases, and their associated f -statistics.

238 **d** Schematic representation of the computations behind the *ADMIXTOOLS* programs *qpGraph*, *qpWave*, and
239 *qpAdm*. *ADMIXTOOLS* 2 separates the computation of f_2 -statistics from the later steps in the pipeline. Shown
240 below are the number of data points for individuals, SNPs, and populations. The exact number of all
241 possible non-redundant f_3 , f_4 , and f_2 -statistics for k populations are $\binom{k}{2}$, $\binom{k}{4}$, and $\binom{k}{2}$. A small number of
242 f_2 -statistics can be used to obtain a much larger number of f_3 - and f_4 -statistics and requires much less space
243 than the raw genotype data.

244

245

246 *Admixture graph fitting, model comparison, and interpretation*

247

248 There are several challenges that arise when modeling the ancestral relationships among
249 populations with admixture graphs, and *ADMIXTOOLS* 2 implements solutions to several key
250 problems that were not adequately addressed with previous approaches: (1) Automated admixture
251 graph inference; (2) Estimating confidence intervals for admixture graph parameters; (3) Comparing
252 fits of different admixture graphs; (4) Determining identifiability of admixture graph parameters; and
253 (5) Drawing conclusions from a large number of fitting graphs.

254

255 Here, we describe these challenges, how we address them, and how our approaches compare to
256 other approaches, while the Methods section gives detailed descriptions.

257

258 (1) Automated admixture graph inference

259 Constructing admixture graphs manually runs the risk of overlooking models that challenge
260 conventional hypotheses. On the other hand, current methods for inferring admixture graphs
261 automatically (Leppälä et al. 2017, Molloy et al. 2021, Pickrell and Pritchard 2012, Yan et al. 2021) do
262 not allow external information to be integrated into the analysis, and often result in models that
263 may fit the genetic data but can be rejected on other grounds. In addition, *TreeMix* (Pickrell and
264 Pritchard 2012), as well as *OrientAGraph* (Molloy et al. 2021), an improved version of *TreeMix*, can
265 miss admixture graph topologies that exist on parts of the non-convex likelihood surface that are
266 bypassed by these algorithms for exploring admixture graphs (for example, topology M7 in Figure 4
267 of (Molloy et al. 2021)). *MixMapper* (Lipson et al. 2013) and *miqoGraph* (Yan et al. 2021) have a
268 different limitation: exploring topologies with more than one admixture event in the history of any
269 group is not possible. Due to these limitations, many published findings are based on manual
270 proposal of topologies and evaluation of fit, and the great majority of studies using this manual
271 approach (see, for example, Reich et al. 2011, Reich et al. 2012, Lazaridis et al. 2014, Fu et al. 2016,
272 Skoglund et al. 2016, Yang et al. 2017, McColl et al. 2018, Moreno-Mayar et al. 2018, Tambets et al.
273 2018, van de Loosdrecht et al. 2018, Flegontov et al. 2019, Sikora et al. 2019, Wang et al. 2019,
274 Lipson et al. 2020, Shinde et al. 2019, Yang et al. 2020, Hajdinjak et al. 2021, Wang et al. 2021) rely
275 on the software *qpGraph*. We introduce an approach for finding well-fitting admixture graphs
276 automatically that can integrate external information, and that recovers graph topologies more
277 accurately than *TreeMix* (**Figure 2**). External information can be integrated by specifying a set of
278 constraints that admixture graphs must satisfy. This not only ensures that resulting models are
279 temporally plausible, but also cleanly separates prior assumptions from the independent constraints
280 provided by genetic data. Our strategy implemented in the function “*findGraphs*”, differs from
281 *TreeMix/OrientAGraph* in several deep ways, most notably in that it optimizes graphs directly, rather
282 than optimizing trees first and adding admixture events later. This makes it less prone to getting
283 stuck in local optima: our simulation results show that *findGraphs* is more accurate for random
284 graphs (**Figure 2**), and that it can recover specific topologies that pose problems for *TreeMix* and
285 *OrientAGraph*.

286

287 (2) Estimating confidence intervals for admixture graph parameters

288 Since our new implementation of *qpGraph* can evaluate models much more rapidly, it becomes
289 feasible to evaluate the same model multiple times on different SNP sets. This allows us to derive
290 bootstrap confidence intervals (Boos 2003) for all parameters estimated by *qpGraph*, including drift
291 lengths, admixture weights, log-likelihood (LL) scores, and f_4 -statistic residuals. Parameters with
292 extremely wide confidence intervals can thus be immediately shown to be poorly determined. It
293 should also be noted that the estimated confidence intervals do not take into account uncertainty
294 about the graph topology.

295

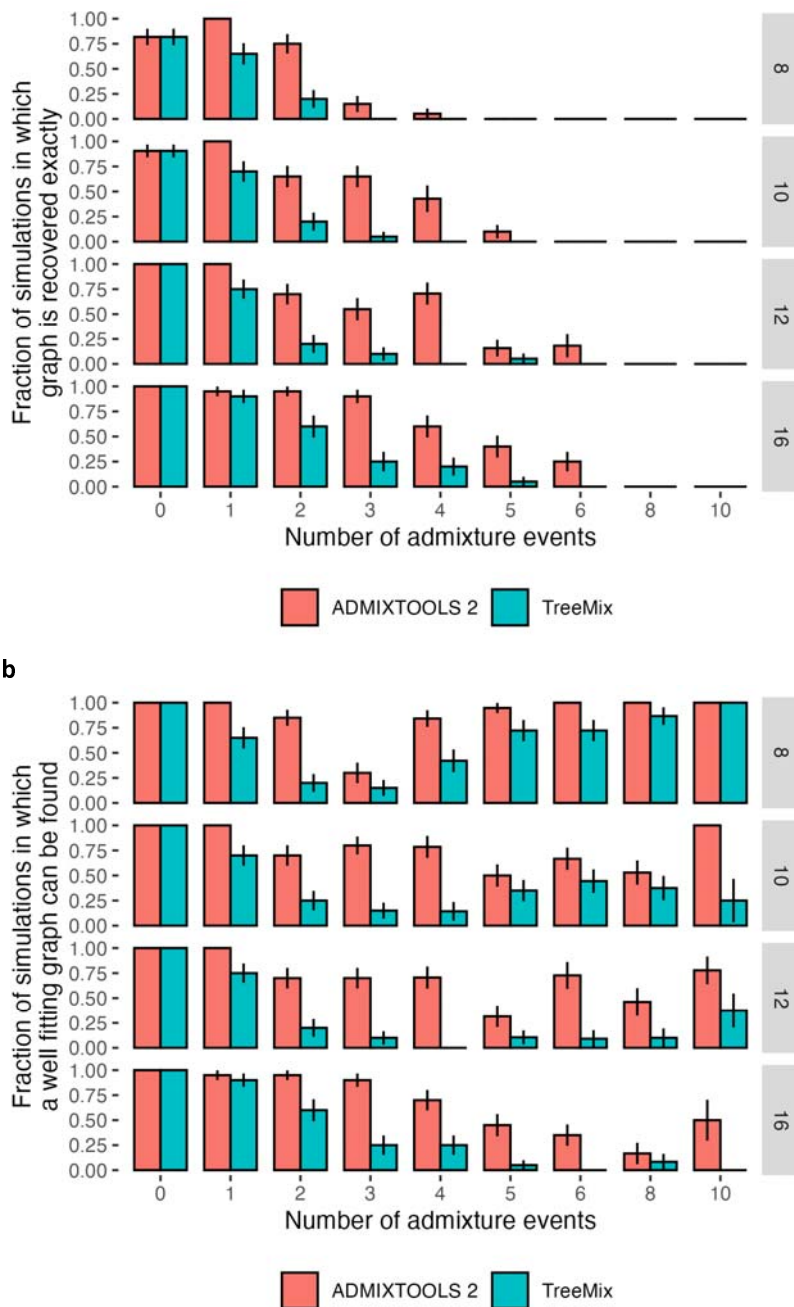
296 (3) Comparing the fits of different admixture graphs

297 Using the bootstrap method for evaluating a graph multiple times on different SNP sets not only
298 allows us to obtain confidence intervals for single graphs, but also allows us to test whether the fit of
299 one graph is significantly better than the fit of another graph, by obtaining confidence intervals for
300 the log-likelihood score difference or worst-residual (WR) difference of two graphs. When we apply
301 this approach to a range of data sets, we find that models with modest log-likelihood differences are
302 often not distinguishable after accounting for the variability across SNPs, even if one might expect
303 them to be distinguishable based on the magnitude of the likelihood difference (**Figure 3a,b**). Thus,
304 previous methods relying on AIC or BIC (such as Shinde et al. 2019, Flegontov et al. 2019) that used
305 specified likelihood difference thresholds to reject some models over others, were over-aggressive.

306

307

308 **a**



309
310

b

311
312

Figure 2:

313 Comparison of accuracy of automated search for optimal topology in the *findGraphs* function of *ADMIXTOOLS*
314 2 and *TreeMix* using simulated graphs with 8, 10, 12, and 16 populations, and 0 to 10 admixture events. Error
315 bars show standard errors calculated as $SE^2 = p(1 - p)/n$ where p is the fraction on the y-axis and n is the
316 number of simulations in each group (typically 20). In the case of *ADMIXTOOLS 2*, we applied *findGraphs* three
317 times on each simulated data set and picked a result with the best fit score. More details are provided in the
318 Methods section. **a** Fraction of simulations where the simulated graph is recovered exactly. **b** Fraction of
319 simulations where the simulated graph is either recovered exactly, or the score is at least as good as the score
320 of the simulated graph, when both graphs are evaluated by *ADMIXTOOLS 2*. More admixture edges greatly
321 increase the search space and make it more difficult to recover the simulated graph, but they do make it easier
322 to find alternative graphs with good fits.

323

324 A second challenge in comparing different admixture graph models arises when comparing models
325 of different complexity (i.e., with a different number of admixture events). Established methods
326 such as AIC and BIC are applicable and also can account for different model complexity if the

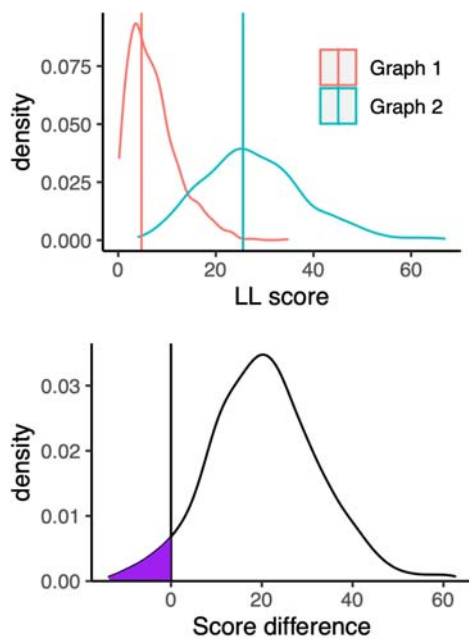
327 number of independent parameters in a model is known. However, the number of independent
328 parameters estimated in an admixture graph is not simply determined by the number of groups,
329 drift edges, or admixture events, as it also depends on the graph topology in a complex way. We
330 implement a method to compare admixture graph models of different complexity by using a new
331 scoring function, which uses different blocks of SNPs for deriving fitted and estimated f -statistics.
332 This ensures that our model comparison test does not favor more complex models by allowing them
333 to overfit the data. This cross-validation approach can also be used to rank alternative models of the
334 same complexity and deal with overfitting. We note that the calculation of cross-validated log-
335 likelihood scores is not turned on in *findGraphs* by default, and to make our results more
336 comparable to those of the published studies we revisited, we relied on standard log-likelihood
337 scores below.

338

339 To test if our method is well calibrated, we simulated 100,000 SNPs under the same graph in 1000
340 replicates. We then created two new topologies by removing one out of two symmetric edges from
341 the first graph (**Figure 3c**). These new incorrect models are symmetrically related to the first graph
342 and can be used to test the null hypothesis that the true difference in log-likelihood scores of these
343 two graphs is zero. The uniform distribution of p -values confirms that our method is well calibrated
344 (**Figure 3d**). A caveat is that only one symmetric topology was explored in this way.

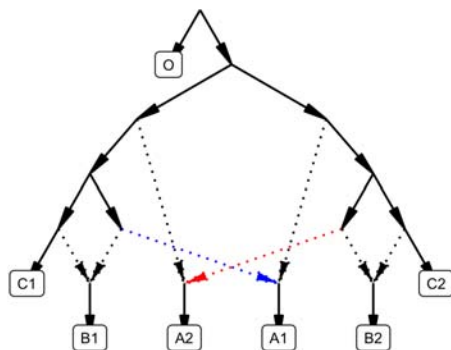
345

346 **a, b**



347

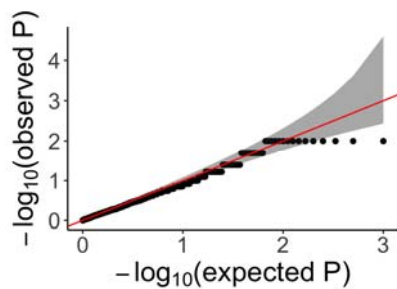
348 **c**



349

350

351 **d**



352

353

Figure 3: Calibrating the bootstrap model comparison approach

354

a Bootstrap sampling distributions of the log-likelihood scores for two admixture graphs (shown in **Figure S1**) of the same populations fitted using real data. Vertical lines show the log-likelihood scores computed on all SNP blocks.

356

357

b Distribution of differences of the bootstrap log-likelihood scores for both graphs (same data as in **a**). The purple area shows the proportion of resamplings in which the first graph has a higher score than the second graph. The two-sided p -value for the null hypothesis of no difference is equivalent to twice that area (or one over the number of bootstrap iterations if all values fall on one side of zero). In this case it is 0.078.

358

359

360

361

c The admixture graph which was used to evaluate our method for testing the significance of the difference of two graph fits on simulated data. We simulated under the full graph and fitted two graphs that result from deleting either the red admixture edge or the blue admixture edge. These two graphs have the same expected $qpGraph$ fit score but can have different scores in any one simulation iteration.

362

363

364

365

d QQ plot of p -values testing for a score difference between the two graphs (on simulated data) under the null hypothesis of no difference, confirming that the method is well calibrated.

366

367

368

369

(4) Determining identifiability of admixture graph parameters

370

Fitting admixture graphs results in an estimate of the overall model fit, as well as in estimates of branch lengths and admixture weights. However, even with infinite data some of these parameters cannot be estimated, as they are not identifiable from the system of equations that corresponds to the admixture graph. Issues like this have been well described for simple topological features of a graph. For example, the lengths of the two branches connected to the root node cannot be estimated without either fixing one of them or forcing the lengths to be evenly distributed.

374

375

376

Furthermore, in a graph with populations and admixture events, at least one parameter will not be identifiable unless the inequality – is satisfied (Lipson 2020). However, even in

377

378

graphs that meet this criterion, some parameters are not identifiable, and until *findGraphs*, there was no method for testing whether any given parameter in an admixture graph is identifiable. We introduce a method for testing which parameters in an admixture graph are identifiable, and which are not, based on the Jacobi matrix of the graph's system of f_2 equations. Like our method for deriving confidence intervals for admixture graph parameters, this can improve the interpretability of admixture graph analyses.

379

380

381

382

383

384

385

Our methods for automated topology inference, for bootstrapping log-likelihood scores or worst residuals for comparing model fits, for cross-validation of admixture graphs of different complexities, for estimating confidence intervals and determining identifiability of admixture graph parameters can greatly improve the interpretability of admixture graph analyses. We implement all these methods in *findGraphs* to assist the user in building a series of models that can explain admixture processes of ancient populations in a manner that surpasses all other programs in accuracy.

386

387

388

389

390

391

392

393

394

(5) Drawing conclusions from a large number of fitted models

395

As discussed in the Results section, when we apply our methods for finding optimal graphs and comparing admixture graphs to a number of previously published models, we find that there often exists a much larger number of fitting models than has previously been appreciated. In these cases, we are unable to prioritize a single model, or even a small number of models, based on the evidence

396

397

398

399 we have. However, we are still able to reject the vast majority of all tested models. This suggests
 400 that insight can be gained by identifying common features among the well-fitting models. We
 401 therefore introduce methods for summarizing collections of well-fitting admixture graphs to
 402 determine which features they share. In practice, we find that these methods can aid the manual
 403 inspection of *findGraphs* results, but the high diversity of well-fitting topologies we see in most case
 404 studies and the importance of fitted parameters (especially admixture proportions) for historical
 405 interpretation of topologies makes it difficult to reliably automatize the process of interpreting fitted
 406 admixture graph models.

407

408 **Revisiting published admixture graphs**

409

410 We studied admixture graphs from eight publications (Lazaridis et al. 2014, Shinde et al. 2019, Sikora
 411 et al. 2019, Bergström et al. 2020, Lipson et al. 2020, Hajdinjak et al. 2021, Librado et al. 2021, Wang
 412 et al. 2021) with the goal of comparing published models to models identified by our algorithm for
 413 automatically inferring optimal (best-fitting) admixture graphs (**Table 1**). In all but one study
 414 *qpGraph* or its automated reimplementation (*admixturegraph*, Leppälä et al. 2017) was used for
 415 fitting topologies to genetic data, while Librado et al. (2021) relied on the automated *OrientAGraph*
 416 method. The main question we were interested in is whether we can find alternative models which
 417 (1) fit as well as, or better than the published graph, (2) differ in important ways from the published
 418 graph, and (3) cannot immediately be rejected based on other evidence such as temporal
 419 plausibility. The studies were selected according to the criterion that an admixture graph model
 420 inferred in the study is used as primary evidence for at least one statement about population history
 421 in the main text of the study. In other words, the admixture graph method was used in the original
 422 studies to generate new insights into population history, and not simply to show that there is a
 423 model that exists (even if it is not unique) that does not contradict results of other genetic analyses,
 424 an approach that is a valid use of admixture graphs and has been taken in some studies (e.g., Seguin-
 425 Orlando et al. 2014, Narasimhan et al. 2019, Wang et al. 2019). There are many published studies
 426 that could have been included in our re-evaluation exercise as they meet our key criterion (e.g.,
 427 Yang et al. 2017, McColl et al. 2018, Posth et al. 2018, Flegontov et al. 2019, Calhoff et al. 2021,
 428 Lipson et al. 2022, Vallini et al. 2022). However, critical re-evaluation of each published graph was an
 429 intensive process, and the sample of studies we revisited is diverse enough to identify some general
 430 patterns and to support strong conclusions.

Publication	Figure in the original publication	Groups (populations)	Admixture events	SNPs used	Publ. model: worst residual, SE	Distinct alternative topologies found	Significantly better fitting topologies, %	Non-significantly better fitting topologies, %	Significantly worse fitting topologies, %	Non-significantly worse fitting topologies, %
Bergström et al. 2020	1e	7	3	312,282	2.1	221	0.5	2.3	16.7	80.5
Lazaridis et al. 2014	3	7	4	642,247	2.2	306	1.0	12.1	80.7	6.2
Shinde et al. 2019	3	8*	3	249,009	2.6	143	0.0	2.8	3.5	93.7
Librado et al. 2021	3b		3		23.9	324	6.8	15.7	24.1	53.4
	Ext 5d	10*	4	1,767,419	14.1	535	0.0	0.0	4.5	95.5
	Ext 5e		5		6.9	784	0.0	0.3	28.4	71.3
Hajdinjak et al. 2021	2d	12	8	263,698	4.8	1,988	15.7	55.7	6.6	22.0
Lipson et al. 2020	Ext 4	12	11**	211,738	2.3	2,000	0.0	11.9	77.1	10.4
Wang et al. 2021	Ext 6	12	8	203,753	3.8	1,778	12.6	84.3	3.1	0.0
Sikora et al. 2019	3f (left)	13	6**	344,903	3.8	894	0.3	17.1	34.6	48.0
	3f (right)	14	6**	613,509	4.2	2,785	0.1	0.9	9.8	89.2

* the population composition was modified, see Table S1 and the text

** certain gene flows were removed from the published model for simplicity, see Table S1 and the text

431

432

433

Table 1: Published graphs in the context of automatically found graphs.

434 We compared graphs from 8 publications to alternative graphs inferred on the same or very similar data (see
435 **Table S1** for details).
436 **Publication:** Last name of the first author and year of the relevant publication.
437 **Figure in the original publication:** Figure number in the original paper where the admixture graph is
438 presented.
439 **Groups (populations):** The number of populations in each graph.
440 **Admixture events:** The number of admixture events in each graph.
441 **SNPs used:** The number of SNPs (with no missing data at the group level) used for fitting the admixture graphs.
442 For all case studies, we tested the original data (SNPs, population composition, and the published graph
443 topology) and obtained model fits very similar to the published ones. However, for the purpose of efficient
444 topology search we adjusted settings for f_3 -statistic calculation, population composition, or graph complexity
445 as noted in the footnotes, in **Table S1**, and discussed in the text.
446 **Publ. model: worst residual, SE:** The worst f -statistic residual of the published graph fitted to the SNP set
447 shown in the “SNPs used” column, measured in standard errors (SE).
448 **Distinct alternative topologies found:** The number of distinct newly found topologies differing from the
449 published one.
450 **Significantly better fitting topologies, %:** The percentage of distinct alternative topologies that fit significantly
451 better than the published graph according to the bootstrap model comparison test (two-tailed empirical p -
452 value <0.05). If the number of distinct topologies was very large, a representative sample of models (1/20 to
453 1/3 of models evenly distributed along the log-likelihood spectrum) was compared to the published one
454 instead, and the percentages in this and following columns were calculated on this sample.
455 **Non-significantly better fitting topologies, %:** The percentage of distinct topologies that fit non-significantly
456 (nominally) better than the published graph according to the bootstrap model comparison test (two-tailed
457 empirical p -value ≥ 0.05).
458 **Non-significantly worse fitting topologies, %:** The percentage of distinct topologies that fit non-significantly
459 (nominally) worse than the published graph according to the bootstrap model comparison test (two-tailed
460 empirical p -value ≥ 0.05).
461 **Significantly worse fitting topologies, %:** The percentage of distinct topologies that fit significantly worse than
462 the published graph according to the bootstrap model comparison test (two-tailed empirical p -value <0.05).
463
464

465 Here we present a high-level summary of these analyses. Discussions of individual graphs, as well an
466 overview of the methodology, can be found in the next section.
467

468 For 19 out of 22 published graphs we examined we were able to find at least one, but usually many,
469 graphs of the same complexity (number of groups and admixture events) and with a log-likelihood
470 score that was nominally better than that of the published graph (see results for 11 selected graphs
471 in **Table 1** and full results for all 22 in **Table S1**). The 22 graphs were drawn from the 8 publications
472 as there were multiple final graphs presented in some of the publications (Shinde et al. 2019, Sikora
473 et al. 2019, Librado et al. 2021), or we examined selected intermediates in the model construction
474 process (Bergström et al. 2020, Lazaridis et al. 2014, Lipson et al. 2020, Wang et al. 2021), or we
475 introduced an outgroup not used in the original study (Hajdinjak et al. 2021, Sikora et al. 2019), or
476 we tested additional graph complexity classes dropping “unnecessary” admixture events (Lipson et
477 al. 2020, Sikora et al. 2019).
478

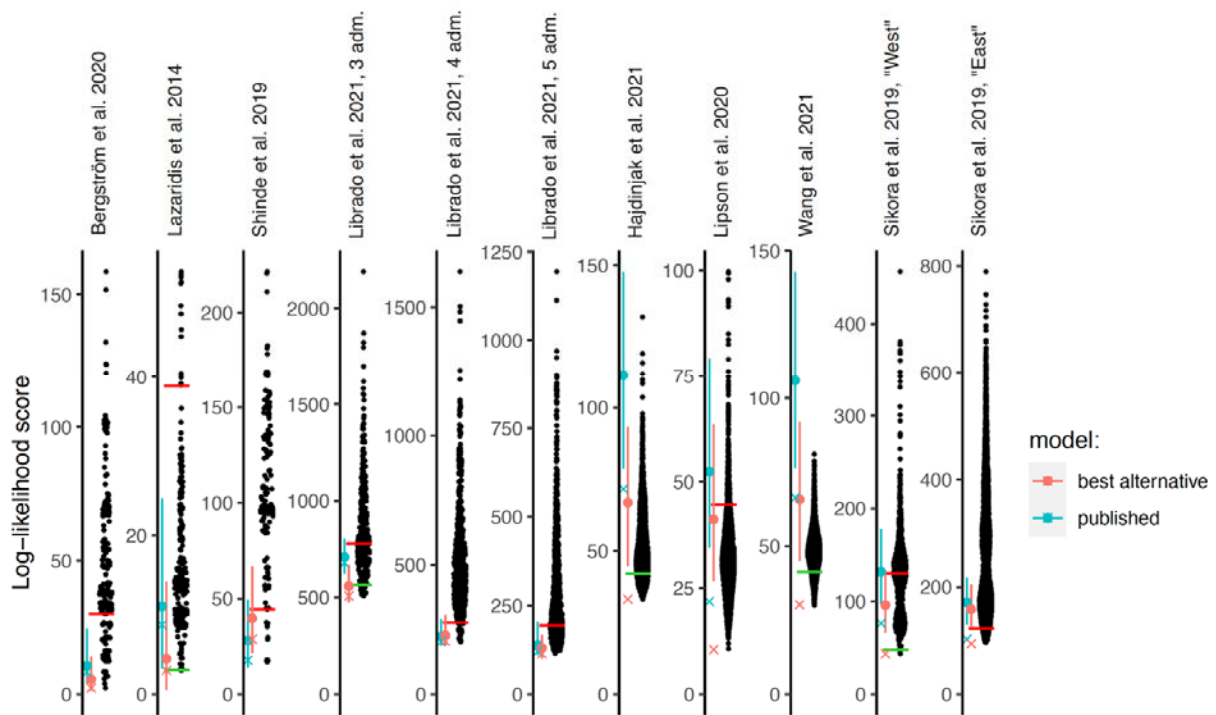
479 Oftentimes these alternative graphs are not significantly better than the published one after taking
480 into account variability across SNPs via bootstrapping. In the following cases, at least one model that
481 fits significantly better than the published one according to our bootstrap model comparison
482 method was found: the Bergström *et al.* and Lazaridis *et al.* seven-population graphs; the Librado *et al.*
483 graph with 3 admixture events; the Hajdinjak *et al.* graphs with or without adding a chimpanzee
484 outgroup; the Lipson *et al.* intermediate graphs with 7 groups and 4 admixture events and with 10
485 groups and 8 admixture events; the Wang *et al.* 12-population graph; and the Sikora *et al.* graphs for
486 West Eurasians and for East Eurasians with 10 or 6 admixture events (**Table S1**). In nearly all cases
487 (except for the Lazaridis *et al.* six-population graph, Shinde *et al.* graph with 8 populations and 3
488 admixture events, and the Librado *et al.* graph with 4 admixture events), we also identified a large
489 number of graphs that fit the data not significantly worse than the published ones. In every example,

490 some of these graphs have topologies that are qualitatively different in important ways from those
 491 of the published graphs. Features such as which populations are admixed or unadmixed, direction of
 492 gene flow, or the order of split events, if not constrained *a priori*, are generally not the same
 493 between alternative fitting models for the same populations. While some of these graphs can be
 494 rejected since their topologies appear highly unlikely because of non-genetic or unrelated genetic
 495 evidence, for all of the publications except one (Shinde et al. 2019), there are alternative equally-
 496 well-or-better-fitting graphs we identified and examined manually that differ in qualitatively
 497 important ways with regard to the implications about history, are temporally plausible (for instance,
 498 very ancient populations do not receive gene flows from sources closely related to much less ancient
 499 groups), and not obviously wrong based on other lines of evidence. These findings suggest the
 500 possibility that complex admixture graph models, even with a very good fit to the data, may differ in
 501 important ways from true population histories.

502

503 The previous statements are valid if the original parsimony constraints are applied, i.e., if the graph
 504 complexity (the number of admixture events) is not altered. Below in selected case studies (Shinde
 505 *et al.*, Librado *et al.*) we explore the effect of relaxing the parsimony constraint.

506



507

508

Figure 4:

509 Log-likelihood scores of published graphs (those shown in **Table 1**) and automatically inferred graphs. Each dot
 510 represents the log likelihood score of a best-fitting graph from one *findGraphs* iteration (low values of the
 511 score indicate a better fit); only topologically distinct graphs are shown. Log-likelihood scores for the published
 512 models and best-fitting alternative models found are shown by blue and pink x's, respectively. Bootstrap
 513 distributions of log-likelihood scores for these models (vertical lines, 90% CI) and their medians (solid dots) are
 514 also shown. Lower scores of the fits obtained using all SNPs, relative to the bootstrap distribution, indicate
 515 overfitting to the full data set. Green and red horizontal lines show the approximate locations where newly
 516 found models consistently have fits significantly better or worse, respectively, than those of the published
 517 model. In the case of the Bergström *et al.*, Lazaridis *et al.* and Hajdinjak *et al.* studies, one or more worst-fitting
 518 models were removed for improving the visualization. The setups shown here (population composition,
 519 number of groups and admixture events, topology search constraints) match those shown in **Table 1**.

520

521

522 **Detailed reconsideration of eight published admixture graphs**

523

524 We investigated admixture graphs from eight publications. We usually focused on one final graph
525 from each publication, and in some cases we also discuss simpler intermediates in the published
526 model-building process, apply various topological constraints to the model inference process, or
527 decrease/increase the number of admixture events to explore the influence of parsimony
528 constraints. **Table 1** and **Figure 4** summarize these results for one or a few graphs from each
529 publication, while **Table S1** contains the full results for all studied graphs and setups. **Table 2**
530 summarizes our assessment of inferences in the original publications that were supported by the
531 published graphs.

532

533 To identify alternative models, we ran many iterations of *findGraphs* for each set of input
534 populations, constraints, and the number of admixture events we investigated, and we selected the
535 best-fitting graph in each iteration, that is, the graph with the lowest log-likelihood (LL) score. Each
536 algorithm iteration was initiated from a random graph, and the algorithm is non-deterministic so
537 that in each iteration it takes a different trajectory through graph space, possibly terminating in a
538 different final best graph. The number of admixture events in the initial random graphs and in the
539 output graphs was always kept equal to that of the published graph. For each example, we counted
540 how many distinct topologies were found with significantly or non-significantly better or worse LL
541 scores than that of the published graph (**Table 1**, **Table S1**). To obtain a formally correct comparison
542 of model fit, the published graph and each alternative model were fitted to resampled replicates of
543 the dataset and the resulting LL score distributions were compared (see Methods). As shown in
544 **Figure 4**, for 4 of the 8 publications we re-analyzed, the LL score of the published graph run on the
545 full data is better than almost all the bootstrap replicates on the same data (it falls below the 5th
546 percentile), which is a sign of overfitting, and underscores the importance of applying bootstrap to
547 assess the robustness of fitted models and conclusions drawn from them.

548

549 The fraction of graphs with scores better than the score of the published graph should not be
550 overinterpreted, as it is influenced by the *findGraphs* algorithm which does not guarantee ergodic
551 sampling from the space of well-fitting admixture graphs. In particular, it is possible that despite
552 *findGraphs*'s strategies for efficiently identifying classes of well-fitting admixture graphs (see
553 Methods), it has a bias toward missing certain classes of graphs. However, even a single alternative
554 graph which is no worse-fitting than the published graph suggests that we are not able to identify a
555 single best-fitting model. Many of these alternatives, despite providing a good fit to the data, appear
556 unlikely, for example, because they suggest that Paleolithic era humans are mixed between different
557 lineages of present-day humans. We were mainly interested in alternative models which are also
558 plausible, and so we constrained the space of allowed topologies in *findGraphs* to those we
559 considered plausible *a priori*, in cases where this was necessary for reducing the search space size.
560 Constraints were either integrated into the topology search itself, or were applied to outcomes of
561 unconstrained searches, as detailed below.

562

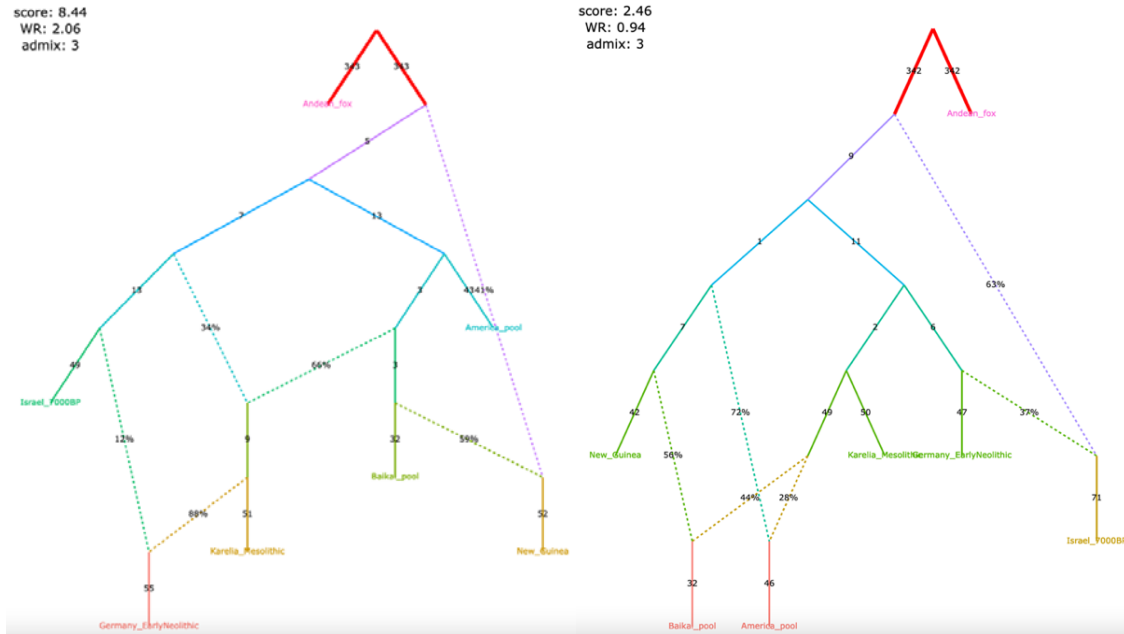
563 **Figure 5:**

564 Published graphs from six studies for which we explored alternative admixture graph fits. In all cases, we
565 selected a temporally plausible alternative model that fits nominally or significantly better than the published
566 model and has important qualitative differences compared to the published model with respect to the
567 interpretation about population relationships. In all but one case, the model has the same complexity as the
568 published model shown on the left with respect to the number of admixture events; the exception is the
569 reanalysis of Librado *et al.* 2021 horse dataset since the published model with 3 admixture events is a poor fit
570 (worst Z-score comparing the observed and expected *f*-statistics has an absolute value of 23.9 even when
571 changing the composition of the population groups to increase their homogeneity and improve the fit relative
572 to the composition used in the published study). For this case, we show an alternative model with 8 admixture
573 events that fits well and has important qualitative differences from the point of view of population history
574 interpretation relative to the published model. The existence of well-fitting admixture graph models does not
575 mean that the alternative models are the correct models; however, their identification is important because
576 they prove that alternative reasonable scenarios exist to published models.

577

578 a, Bergström *et al.* 2020 published
579
580
581
582
583
584

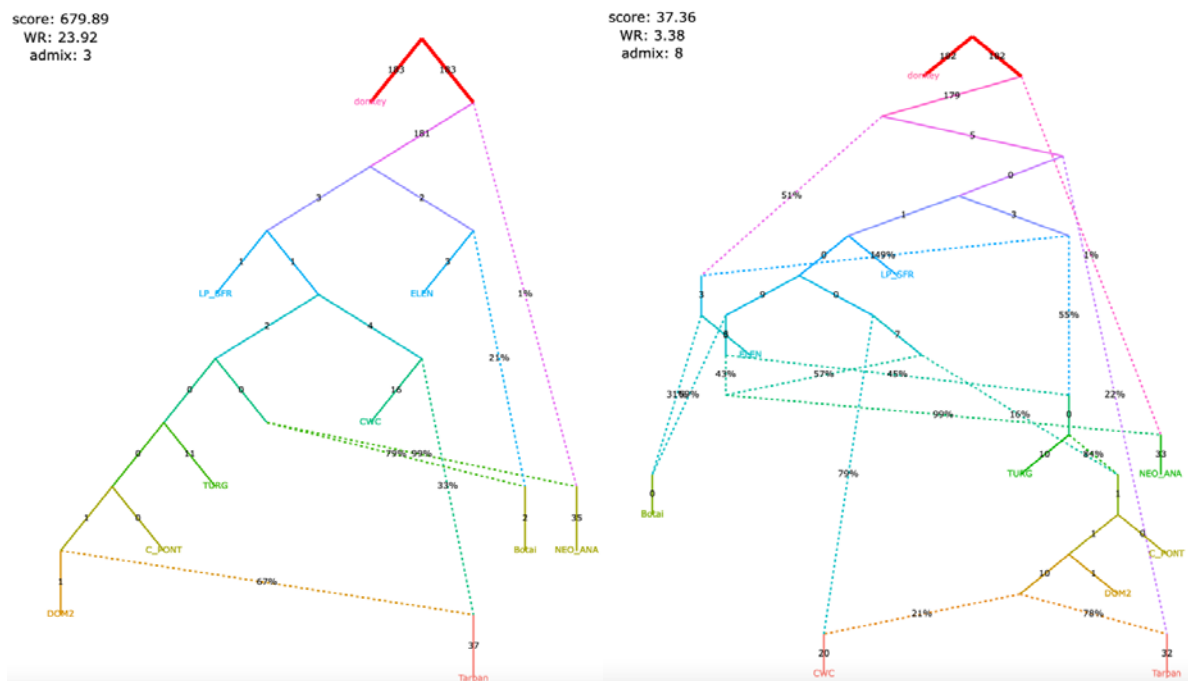
Nominally better fitting graph for dogs that is more congruent to human history. For both species, Baikal and Native American groups are mixed between European- and East Asian-related lineages, and a “Basal Eurasian” lineage contributes to Near Eastern groups; these features are all characteristic of human history but absent in the published dog graph.



585

586 **b**, Librado *et al.* 2021 (modified
587 population composition) published
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602

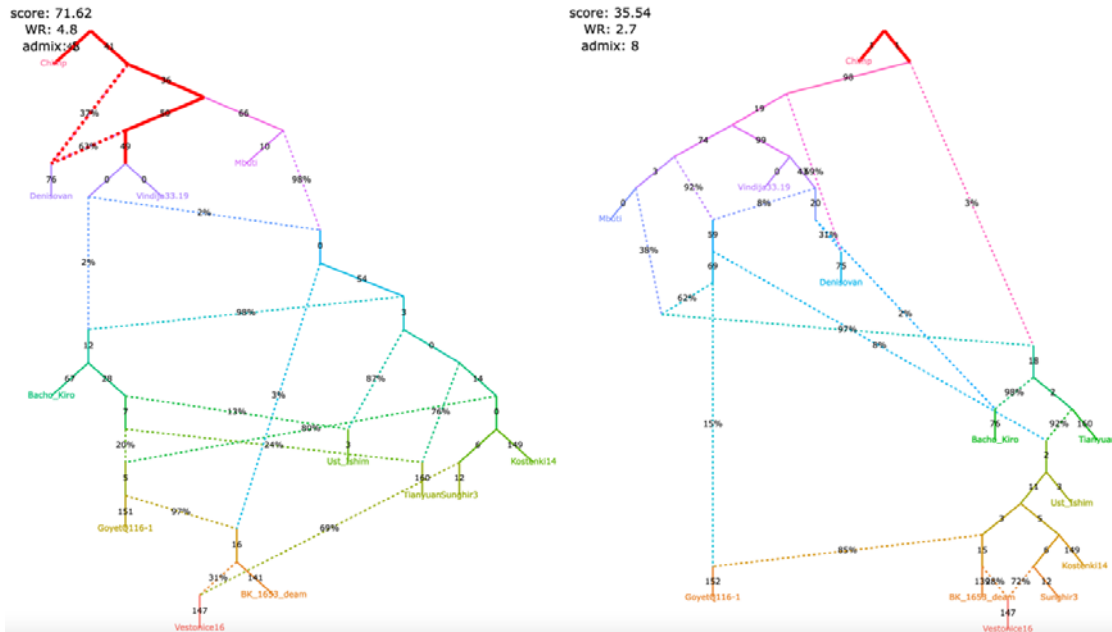
Significantly better fitting, temporally/geographically plausible.
In contrast to the published graph, in this graph with 8 mixture events (the minimum necessary to obtain an acceptable statistical fit to the data), a lineage maximized in horses associated with Yamnaya steppe pastoralists or their Sintashta descendants (C-PONT, TURG, or DOM2) contributes a substantial amount of ancestry to the horses from the Corded Ware archaeological context (CWC). Thus, in this model both CWC humans and horses are mixtures of Yamnaya and European farmer-associated lineages. This is qualitatively different from the suggestion that there was no Yamnaya-associated contribution to CWC horses which was a possibility raised in the paper. The eight-admixture graph also is different from the published model in that it shows a fitting model where the Tarpan horse does not have the history claimed in the study (as an admixture of the CWC and DOM2 horses).



604

605 c, Hajdinjak *et al.* 2021 published
606
607
608
609
610
611
612
613
614
615
616
617

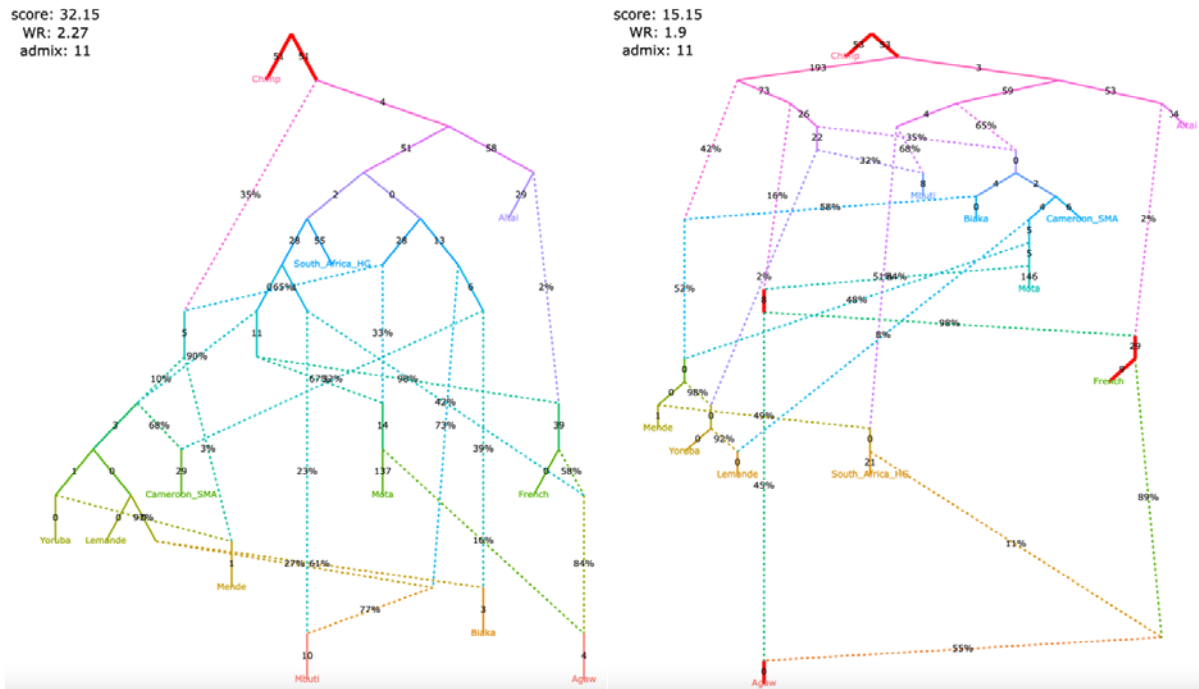
Significantly better fitting, but without a specific lineage shared between Bacho Kiro Initial Upper Paleolithic and East Asians. In this model, all the lineages shared between Bacho Kiro IUP and East Asians contributed a large fraction of the ancestry of later European hunter-gatherers as well, and thus this graph does not imply distinctive shared ancestry between the earliest modern humans in Europe and later people in East Asia, and instead could be explained by a quite different and also archaeologically plausible scenario of a primary modern human expansion out of the Near East contributing serially to the major lineages leading to Bacho Kiro, then later East Asians, then Ust'-Ishim, then the primary ancestry in later European hunter-gatherers.



618

619 d, Lipson *et al.* 2020 published
620
621
622
623
624
625
626
627
628
629
630
631
632

Nominally better fitting. In contrast to the published graph, there is no single lineage specific to modern rainforest hunter-gatherers (Biaka and Mbuti) and Shum Laka (Cameroon_SMA). Rather, the primary ancestries in each group are separate deep-branching lineages (the deeper lineage they all share is also the source of the majority of ancestry in all anatomically modern humans modeled here). In contrast to the graph in the published paper, there is no West African-maximized ancestry present in mixed form in Biaka, Mbuti, and Shum Laka; archaic admixture is not limited to a subset of Africans, but is present in all anatomically modern humans in various proportions; and there is no ghost modern human ancestry in Agaw, Biaka, Lemande, Mbuti, Mende, Mota, Shum Laka, and Yoruba.

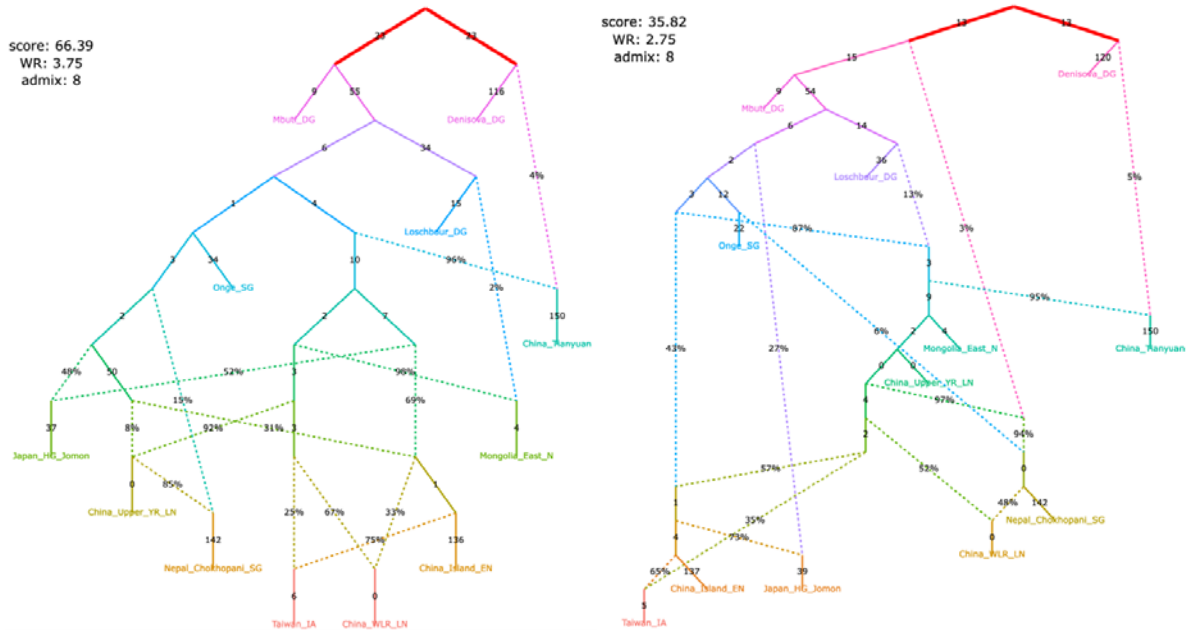


633
634

635
636
637
638
639
640

e, Wang *et al.* 2021 published

Significantly better fitting while meeting the constraints used to inform model-building in the published paper. The finding of Onge-related admixture that is widespread in East Asia suggesting an early peopling via a coastal route is not a feature of this model.

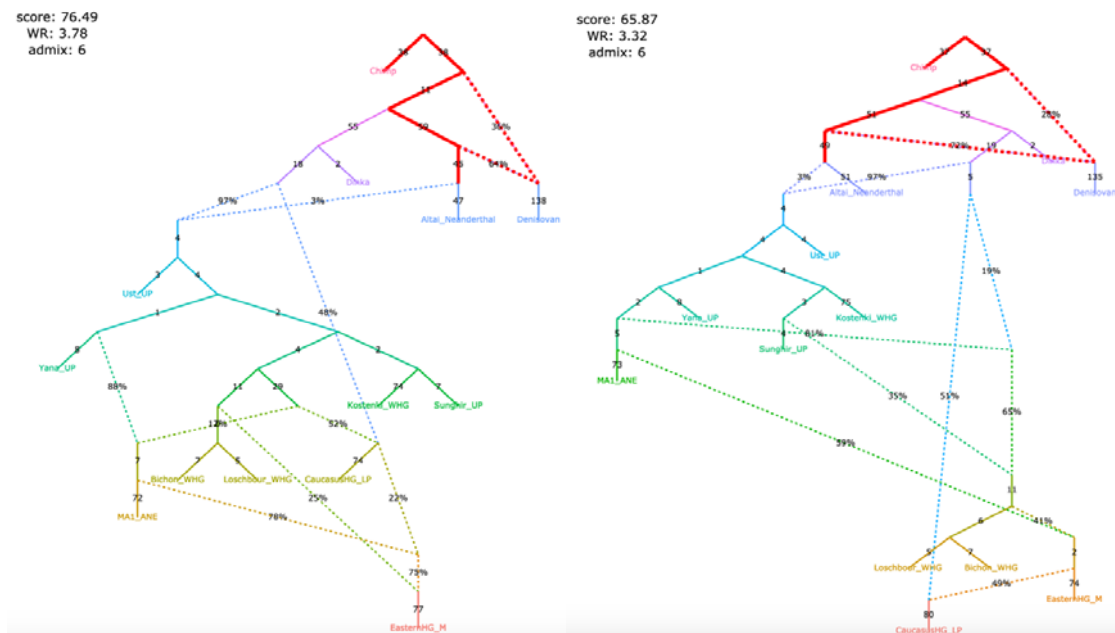


642
643
644
645
646
647

f, Sikora *et al.* 2019 (simplified Western graph)

641

Nominally better fitting. The striking feature of the admixture graph suggested in the paper whereby Mal'ta (MA1_ANE) derives some ancestry from the CHG-associated lineage is not a feature of this alternative model.



648

649
650

651 **Bergström et al. 2020.** The admixture graph for dogs in Figure 1e of Bergström et al. 2020 was
652 inferred by exhaustively evaluating all graphs with two admixture events and outgroup ‘Andean fox’
653 for the six populations that remain in the graph after excluding an Early Neolithic dog from
654 Germany. The only six-population graph with a worst residual (WR) below 3 standard errors (SE) was
655 then chosen as a scaffold onto which the Early Neolithic dog genome from Germany was mapped,
656 allowing for one more admixture event, and a seven-population graph with the lowest LL score was
657 shown as the final model in the paper. Alternative six-population scaffolds were not explored in the
658 original publication, although two six-population graphs with fits very similar to the best one were
659 found. LL was not used as a ranking metric for alternative models in the original study; instead, the
660 number of f_4 -statistics having residuals above 3 SE was considered. Since no f_3 -statistics were
661 negative when all sites available for each population triplet were used (the “useallsnps: YES” option),
662 we did not use the upgraded algorithm for calculating f_3 -statistics on pseudo-diploid data.

663

664 Our *findGraphs* results confirm the published six-population graph in that no graph with lower LL
665 score is identified, but 3 of 14 unique alternative graphs found fit not significantly worse than the
666 published graph (the published graph was also recovered by *findGraphs*) (**Table 1, Table S1**). When
667 we used *findGraphs* to infer seven-population graphs with three admixture events (again fixing
668 Andean fox as the outgroup), we identified 5 graphs with log-likelihood score nominally better than
669 that of the published graph and one with a score that is slightly lower than that of the published
670 graph but actually significantly better according to our model comparison methodology (this model
671 is very similar to the published graph, **Figure S2**). In the newly found seven-population graph with
672 the best log-likelihood score (**Figure S2**), the Siberian (Baikal), American, and Levantine dogs are
673 admixed, and the West European, East European (Karelia), and dogs of Southeast Asian origin (New
674 Guinea singing dog) are unadmixed, while the opposite pattern is found in the published graph
675 (**Table 2**). The best-fitting graph does not fit the data significantly better than the published graph
676 (two-tailed empirical p -value = 0.332), but it bears a closer resemblance to the human population
677 history (see the third-best graph found by *findGraphs* on human data from Bergström *et al.* 2020 in
678 **Figure S3**) than the published seven-population graph (**Figure 5a, Figure S2**).

679

680 In this new seven-population model (**Figure 5a**), both American and Siberian dog lineages represent
681 a mixture between groups related to the Asian and East European dog lineages, and robust genetic
682 results suggest that in the time horizon investigated in the original publication (after ca. 10,900 years
683 ago) nearly all Siberian (Jeong et al. 2019, Sikora et al. 2019) and all American (Raghavan et al. 2014,
684 Raghavan et al. 2015, Moreno-Mayar et al. 2018) human populations were admixed between groups
685 most closely related to Europeans and Asians. According to this model, Levantine dogs are modeled
686 as a mixture of a basal branch (splitting deeper than the divergence of the Asian and European dogs)
687 and West European dogs, again in agreement with current models of genetic history of Middle
688 Eastern human populations who are modeled as a mixture of “basal Eurasians” and West European
689 hunter-gatherers (Lazaridis et al. 2016, Lipson et al. 2017). Although greater congruence with human
690 history increases the plausibility of *findGraphs*’s newly identified model relative to the published
691 model, to make unbiased comparisons between the history of the two species, model selection
692 should be done strictly independently for each species, and so the genetic data alone does not favor
693 one model more than another. Our results provide a specific alternative hypothesis that differs in
694 qualitatively important ways from the published model and can be tested against new genetic data
695 as it becomes available as well as other lines of genetic analysis of existing data.

696

697 To explain why the original paper on the population history of dogs missed the model that
698 *findGraphs* identified that is plausibly a closer match to the true history, we observe that the
699 Bergström *et al.* 2020 admixture graph search was exhaustive under the parsimony constraint (no
700 more than 2 admixture events for 6 populations), and thus missed the potentially true topology
701 including 3 admixture events for these 6 populations. This case study also illustrates that even in a
702 relatively low complexity context (7 groups and 3 admixture events) applying manual approaches for
703 finding optimal models is risky. When any new group such as an Early Neolithic dog from Germany is
704 added to the model, it may introduce crucial information into the system, and re-exploring the

705 whole graph space in an automated way is advisable. In contrast, mapping a newly added group on a
 706 simple skeleton graph (even when that skeleton is a uniquely best-fitting model) may yield a
 707 topology that is at odds with the true history. As the original Bergström *et al.* paper noted (Fig. 3C of
 708 that study), no congruent six-population graph models were found for humans and dogs under the
 709 parsimony assumption: the three most congruent graphs for dogs resulted in poorly fitting models
 710 for the corresponding human populations (WR above 10 SE), and the three most congruent graphs
 711 for humans resulted in poorly fitting models for the corresponding dog populations (WR between 5
 712 and 10 SE). We added a West European hunter-gatherer group to the set of human groups from the
 713 original publication and using *findGraphs* on the original set of 77K transversion SNPs we found that
 714 the third best-fitting model for humans (**Figure S3**) (which is not significantly different in fit from the
 715 first one) is topologically identical to the newly found dog graph.

716

717 Even though *findGraphs* identified an admixture graph topology that fits the data as well as the
 718 seven-population graph in Bergström *et al.* and is qualitatively quite different with respect to which
 719 populations were admixed, the new topology continues to support another of the key inferences of
 720 that study: that many of the early divergences among domesticated dog lineages occurred prior to
 721 the date of the Karelian dog (~10,900 ya). Thus, both graphs concur in providing strong evidence
 722 that the radiation of domesticated dog lineages occurred by the early Holocene, prior to the
 723 domestication of other animals. We further emphasize that the Bergstrom *et al.* 2020 graph is the
 724 best-case scenario (along with Lazaridis *et al.* 2014 discussed below) for published admixture graphs.
 725 Most published graphs are far less stable even than this.

726

727 **Table 2: Features of the published admixture graphs that support inferences in the original studies**
 728 **and the level of their support in our re-analysis.**

729 The table lists key features whose support we assessed in sets of alternative well-fitting and temporally
 730 plausible models generated by *findGraphs*. Since this assessment had to be performed manually, only in two
 731 cases (marked by asterisks) all models fitting better and non-significantly worse than the published one were
 732 scrutinized; in other cases only a subset of best-fitting models was examined (see the sections for details).

733

Study	Groups / admixture events	Features of the published model supported by temporally plausible alternative models generated by <i>findGraphs</i>	Features of the published model <u>not</u> supported by temporally plausible alternative models generated by <i>findGraphs</i>
Bergström et al. 2020	7 / 3 *	Early divergence of domesticated dog lineages (prior to the date of the Karelian dog, 10,900 ya).	Siberian (Baikal), American, and Levantine dog lineages are unadmixed, and the West European (Germany Early Neolithic), East European (Karelia), and dogs of Southeast Asian origin (New Guinea singing dog) are admixed.
Lazaridis et al. 2014	7 / 4	Present-day Europeans represent a mixture of three ancestral sources related to the following groups: Mal'ta (MA1), West European hunter-gatherers, and early European farmers.	N/A
Shinde et al. 2019	8 / 3 *	(1) Iranian farmer-related ancestry in the Indus Periphery group is not derived from the Hajji Firuz Neolithic or Tepe Hissar Chalcolithic groups. (2) There is Asian-related ancestry in the Indus Periphery group.	N/A

	8 / 4	(2) There is Asian ancestry in the Indus Periphery group.	(1) Iranian farmer-related ancestry in the Indus Periphery group is not derived from the Hajji Firuz Neolithic or Tepe Hissar Chalcolithic groups.
Librado et al. 2021	10 / 8 or 9	(2) DOM2 and C-PONT are sister groups (they form a clade); (4) there was gene flow from a deep-branching ghost group to the NEO-ANA group.	(1) NEO-ANA-related admixture is absent in the DOM2 group; (3) there is no gene flow connecting the CWC group and the cluster associated with Yamnaya horses and horses of the later Sintashta culture whose ancestry is maximized in the Western Steppe (DOM2, C-PONT, TURG); (5) Tarpan is a mixture of a CWC-related and a DOM2-related lineage.
Hajdinjak et al. 2021	12 / 8	(3) the Vestonice16 lineage is a mixture of a Sunghir-related and a BK1653-related lineage.	(1) there are gene flows from the lineage found in the ~45,000-43,000-years-old Bacho Kiro Initial Upper Paleolithic (IUP) associated lineage to the Ust'-Ishim, Tianyuan, and GoyetQ116-1 lineages; (2) the ~35,000-years-old Bacho Kiro Cave individual BK1653 belonged to a population that was related, but not identical, to that of the GoyetQ116-1 individual.
Lipson et al. 2020	12 / 11	N/A	(1) A lineage maximized in present-day West African groups (Lemane, Mende, and Yoruba) also contributed some ancestry to the ancient Shum Laka individual and to present-day Biaka and Mbuti; (2) another ancestry component in Shum Laka is a deep-branching lineage maximized in the rainforest hunter-gatherers Biaka and Mbuti; (3) "super-archaic" ancestry (i.e., diverging at the modern human/Neanderthal split point or deeper) contributed to Biaka, Mbuti, Shum Laka, Lemane, Mende, and Yoruba; (4) a ghost modern human lineage (or lineages) contributed to Agaw, Mota, Biaka, Mbuti, Shum Laka, Lemane, Mende,

			and Yoruba.
Wang et al. 2021	12 / 8	N/A	Admixture from a source related to Andamanese hunter-gatherers is almost universal in East Asians, occurring in the Jomon, Tibetan, Upper Yellow River Late Neolithic, West Liao River Late Neolithic, Taiwan Iron Age, and China Island Early Neolithic (Liangdao) groups.
Sikora et al. 2019 “West”	13 / 6	N/A	The Mal’ta (MA1_ANE) lineage received a gene flow from the Caucasus hunter-gatherer (CaucasusHG_LP or CHG) lineage.
Sikora et al. 2019 “East”	14 / 6	(2) European-related ancestry in the Kolyma, USR1, and Clovis lineages is closer to Mal’ta than to Yana.	(1) the Mal’ta (MA1_ANE) and Yana (Yana_UP) lineages received gene flow from a common East Asian-associated source diverging before the ones contributing to the Devil’s Cave (DevilsCave_N), Kolyma (Kolyma_M), USR1 (Alaska_LP), and Clovis (Clovis_LP) lineages; (3) the Devil’s Cave lineage received no European-related gene flows, and Kolyma has less European-related ancestry than ancient Americans (USR1 and Clovis).

734

735

736

737 **Lazaridis et al. 2014.** The graph in Figure 3 (Lazaridis et al. 2014) was inferred in the following
738 manner. First, a phylogenetic tree without admixture was constructed which was the best fit for all
739 f_4 -statistics among the populations “Mbuti”, “WHG Loschbour” (Lazaridis et al. 2014), “LBK
740 Stuttgart” (Lazaridis et al. 2014), “Onge”, and “Karitiana”, with “Mbuti” fixed as an outgroup. Next,
741 all possible admixture graphs were considered that result from adding a single admixture edge to
742 this tree. After it was found that each of them had a $WR > 3 SE$, several graphs with two admixture
743 events were considered, and some of them had $WR < 3 SE$. The “MA1” genome (Raghavan et al.
744 2014) was added to these graphs in several different ways, and only one of these configurations was
745 found to have $WR < 3 SE$. This was then used as a skeleton graph onto which a European population
746 (represented by different present-day groups) was added. No fitting graph was found in which
747 present-day Europeans could be modeled as a two-way mixture (adding one admixture event to the
748 graph). After inspecting the non-fitting f -statistics of one of these graphs, it was found that modeling
749 modern Europeans as a three-way mixture (adding two admixture events to the graph) is consistent
750 with all f -statistics. Six f_3 -statistics were negative when all sites available for each population triplet
751 were used (the “useallsnps: YES” option, **Table S1**), but the upgraded algorithm for calculating f_3 -
752 statistics on pseudo-diploid data had no effect since the only pseudo-diploid group in the dataset
753 (MA1) was a singleton population, and the algorithm removes sites with only one chromosome
754 genotyped in any non-singleton population. Thus, below we show results generated using the
755 standard algorithm for calculating f -statistics.

755

756

757

758

First, we considered the published skeleton graph onto which a European population was later added (Lazaridis et al. 2014). As in the (Bergström et al. 2020) example, the best six-population graph with two admixture events found by *findGraphs* is identical to the published six-population

759 graph, which has a LL score of 3.0 (**Table S1**). The second-best graph found has a LL score of 31.8.
760 When computing the bootstrap p -value for the difference between these two graphs, we find that in
761 1.6% of all SNP resamplings the second-best graph has a better score than the published graph,
762 resulting in a two-tailed empirical p -value of 0.032 for a difference in fits between these two graphs.
763 All 14 alternative graphs found by our algorithm fit significantly worse than the published graph
764 (**Table S1**).

765
766 When we add the European population (French) and consider seven-population graphs with four
767 admixture events, we find 40 out of 306 distinct graphs with a score better than that of the
768 published graph (10 of those graphs are shown in **Figure S4**). The best-fitting newly found model and
769 two other models fit the data significantly better than the published model (**Table S1**), but their
770 topology is qualitatively very similar to that of the published graph (**Figure S4**). In the best-fitting
771 newly found model, French and Karitiana share some drift to the exclusion of MA1, while in the
772 published model the source of MA1-related ancestry in French is closer to MA1 than to Karitiana.

773
774 It is important to point out that not all of the 40 alternative graphs that fit nominally or significantly
775 better than the published one are consistent with the conclusion that modern European populations
776 are admixed between three different ancestral populations (**Figure S4**). For example, the fifth
777 alternative graph in **Figure S4** that is fitting nominally better than the published model (p -value =
778 0.464) includes no basal Eurasian ancestry in early European farmers (LBK Stuttgart), and instead
779 models Onge as having ~50% West Eurasian-related ancestry and MA1 as having ~25% Asian
780 ancestry. According to that graph, the present-day European population was formed by admixture
781 of a MA1-related lineage and a European Neolithic-related lineage, with no West European hunter-
782 gatherer (WHG) contribution. Of course, other lines of evidence make it clear that LBK Stuttgart is a
783 mixture of Anatolian farmer-related ancestry and WHG Loschbour-related ancestry (Lazaridis et al.
784 2016, Lipson et al. 2017), thus providing external information in favor of the Lazaridis *et al.* 2014
785 model, and the use of such ancillary information in concert with graph exploration is important in
786 order to obtain more confident inferences about population history taking advantage of admixture
787 graphs.

788
789 The second alternative graph in **Figure S4** that fits just negligibly worse than the highest-ranking
790 model has another distinctive feature: LBK Stuttgart is modeled as a mixture of a WHG-related and a
791 basal Eurasian lineage, but modern Europeans receive a gene flow not from the LBK-related lineage,
792 but from its basal Eurasian source. Although temporally plausible, this model is much less plausible
793 from the archaeological point of view than the published model, and thus in this case too we can
794 reject it as unlikely based on non-genetic evidence. We note, however, that a large group of newly
795 found models (247 graphs) fits not significantly worse than the published one (**Table S1**), and those
796 are topologically diverse. Thus, strictly speaking, the admixture graph method on the given dataset
797 cannot be used to prove that the published model is the only one fitting the data.

798
799
800 **Shinde et al. 2019**. The skeleton admixture graph in the original study (Shinde et al. 2019) was
801 constructed manually on the basis of a SNP set derived from the 1240K enrichment panel, and
802 subsequently all possible branching orders (105) within the five-population Iranian farmer-related
803 clade were tested. The published model (Fig. 3 in that study) included 9 groups and 3 admixture
804 events, but one group (Belt Cave Mesolithic) had a very high missing data rate, and as a result model
805 fitting relied not just on the merged dataset which included 19,000 polymorphic sites without
806 missing data across groups, but also on a dataset with approximately 470,000 sites that excluded the
807 Belt Cave individual. The topological inferences were consistent for both analyses (Table S3 of that
808 study). Following the approach of the published paper, we repeated *findGraphs* analysis both with
809 and without the Belt Cave individual. Thus, we initially explored the following topology classes: 9
810 groups with 3 admixture events on ca. 19,000 polymorphic sites and 8 groups with 3 admixture
811 events on ca. 470,000 sites (**Table S1**). The sample composition of the groups and the SNP dataset

812 matched that in the original study. We summarize results across 4,000 independent iterations of the
813 *findGraphs* algorithm for each topology class.

814

815 For the nine-population graph we found 89 models with LL nominally better than that of the
816 published model (**Table S1**). For the eight-population graph, we found 61 nominally and 4
817 significantly better fitting models (**Table S1**), and their topological diversity was high (**Figure S5**). We
818 note that the following groups were admixed by default in the graph models compared in the
819 original study: Hajji Firuz Neolithic (labeled “Chalcolithic” in that study but the dates are Neolithic)
820 and Tepe Hissar Chalcolithic were considered as mixtures of an Anatolian farmer-related lineage and
821 an Iranian farmer-related lineage; Indus Periphery was considered as a mixture of an Andamanese-
822 related lineage representing ancient South Indians (ASI) and an Iranian farmer-related lineage.
823 However, calculation of negative “admixture” f_3 -statistics for these target groups is impossible using
824 the original dataset and the original model fitting algorithm for several reasons. First, the Indus
825 Periphery group was represented by a single pseudo-diploid individual (I8726) from the “Indus
826 Valley cline” for whom the best-quality data were available. But direct calculation of “admixture” f_3 -
827 statistics for such a group as a target is impossible since its heterozygosity cannot be estimated.
828 Second, as discussed above, in Classic *ADMIXTOOLS* it is impossible to apply a correction intended
829 for accurate calculation of f_3 -statistics on pseudo-diploid data (the “inbreed: YES” option) if there is
830 at least one population composed of one individual only (a singleton population). Third, the original
831 Hajji Firuz Neolithic group composed of five individuals included a family of three 2nd or 3rd-degree
832 relatives, and that artificially inflated the drift on the Hajji Firuz branch and made detecting a
833 negative statistic f_3 (Hajji Firuz Neolithic; X, Y), even if present, highly unlikely. Indeed, no f_3 -statistic
834 turned out to be nominally negative for the groups on the eight-population graph when statistics
835 were calculated according to the original settings (we used settings equivalent to “useallsnps: NO”
836 and “inbreed: NO” in classic *ADMIXTOOLS*, 470,389 polymorphic sites were available). Considering
837 this fact, it is not surprising that our automated topology space search is not well constrained. The
838 original study differed from ours since the constraints were introduced manually, but we wanted our
839 topology search to be automatic and to explore a wider range of parameter space.

840

841 In order to provide power to detect negative f_3 -statistics useful for constraining the model search,
842 we 1) removed two members of the family from the Hajji Firuz Neolithic group, 2) extended the
843 number of individuals and sites available for the Indus Periphery group by generating new shotgun-
844 sequencing data for a previously published library (Narasimhan et al. 2019) derived from individual
845 I8726 (see **Table S2**) and by adding published data for three other individuals from the Indus Valley
846 cline (from Gonur in Turkmenistan and Shahr-i-Sokhta in Iran; Narasimhan et al. 2019, Shinde et al.
847 2019), 3) removed from other groups two individuals based on 2nd-3rd degree relatedness, and 4)
848 removed two individuals from other groups based on evidence of contamination with modern
849 human DNA. All the changes to the dataset are shown in **Table S3**. In addition to these dataset
850 adjustments, the new algorithm for calculating f -statistics makes it possible to compute negative f_3 -
851 statistics on pseudo-diploid data, but at a cost of removing sites with only one chromosome
852 genotyped in any non-singleton population (see Methods). We eventually detected significantly
853 negative “admixture” f_3 -statistics f_3 (Tepe Hissar Chalcolithic; Ganj Dareh Neolithic, Anatolia
854 Neolithic), f_3 (Indus Periphery; Ganj Dareh Neolithic, Onge), and other similar statistics for the same
855 target groups. We also observed a nominally negative (Z-score = -0.6) statistic f_3 (Hajji Firuz Neolithic;
856 Ganj Dareh Neolithic, Anatolia Neolithic), which is suggestive but does not by itself prove admixture
857 in Hajji Firuz Neolithic. For this updated analysis, 249,009 variable sites without missing data at the
858 group level were available for the eight populations.

859

860 We repeated topology search with this set of f -statistics providing additional constraints, performing
861 4,000 runs of the *findGraphs* algorithm. The Mota ancient African individual was set as an outgroup
862 and 3 admixture events were allowed in the eight-population graph. Among 4,000 resulting graphs
863 (one from each *findGraphs* run), 144 were distinct topologically, and the published model was
864 recovered in 13 runs of 4,000 (**Table S1**). Only 4 distinct topologies fitting nominally better than the
865 published one were found, and those had LL scores almost identical to that of the published eight-

866 population model (16.97 and 17.66 vs. 17.85). These four alternative models (**Figure S6b**) shared all
867 topologically important features of the published model (**Figure S6a**). Five other topologies differed
868 in important ways from the published one and emerged as fitting the data worse, but not
869 significantly worse, than the published one (**Figure S6c**): two-tailed empirical *p*-values reported by
870 our bootstrap model comparison method ranged between 0.060 and 0.112. Three of these
871 topologies included a trifurcation of Iranian farmer-related lineages leading to the Indus Periphery,
872 Hajji Firuz Neolithic, and Ganj Dareh Neolithic groups. The other two topologies included Hajji Firuz
873 Neolithic as an unadmixed Anatolian-related lineage. In both cases, Indus Periphery was modeled as
874 receiving a gene flow from either the Onge lineage (a proxy for ASI) or a deep Asian lineage.

875

876 The finding that the predominant ancestry component of the Indus Periphery group was the most
877 basal branch in the Iranian farmer clade was a prominent claim of the original study (Shinde et al.
878 2019); for example, the abstract stated: “The Iranian-related ancestry in the IVC derives from a
879 lineage leading to early Iranian farmers, herders, and hunter gatherers before their ancestors
880 separated”. Our finding that the Hajji Firuz Neolithic lineage may be as deep within the Iranian clade
881 as the Indus Periphery lineage or may even diverge from the Anatolian branch shows that this
882 statement cannot be confidently made based on admixture graph analysis alone.

883

884 However, the findings we have described up to this point do not invalidate the broader conclusion
885 that the admixture graph modeling in Shinde *et al.* was used to support; namely (using the phrasing
886 from the abstract) that the genetic data “contradict... the hypothesis that the shared ancestry
887 between early Iranians and South Asians reflects a large-scale spread of western Iranian farmers
888 east.” This finding if correct is important, since it implies that the Iranian-related ancestry in the IVC
889 (Indus Valley Civilization genetic grouping, which is the same group as IP), split from the Iranian-
890 related ancestry in the first Iranian plateau farmers before the date of the Hajji Firuz farmers, who at
891 ~8000 years ago are among the earliest people living on the Iranian plateau known to have grown
892 West Asian crops. The ancient DNA record combined with radiocarbon dating evidence suggests that
893 beginning around the time of the Hajji Firuz farmers, both West Asian domesticated plants such as
894 wheat and barley, and Anatolian farmer-related admixture, began spreading eastward across the
895 Iranian-plateau. If the Iranian-related ancestry in IP was spread eastward into the Indus Valley across
896 the Iranian-plateau as part of the same agriculturally-associated expansion—perhaps brought by
897 people speaking Indo-European languages as well as introducing West Asian crops—then we would
898 expect to see at least some of the Iranian-related ancestry in IP being a clade with that in Hajji Firuz
899 relative to Ganj Dareh. The fact that we do not find any models compatible with this scenario is thus
900 a potentially important finding. In summary, there are two reasons the genetic analyses we have
901 reported up to this point continue to support the finding that the Iranian-related ancestry in IP is not
902 a clade with the Iranian-related ancestry in Hajji Firuz (and Tepe Hissar) and thus is unlikely to reflect
903 the same eastward movement of agriculturalists. First, in *findGraphs* analysis, all models specifying
904 IP and Tepe Hissar and/or Hajji Firuz as a clade relative to Ganj Dareh were significantly worse-fitting
905 than the published one. Instead, either the Iranian-related ancestry in IP definitively splits off first
906 (the topology from Shinde *et al.*), or the branching order of IP, Ganj Dareh, and the Hajji Firuz / Tepe
907 Hissar lineages cannot be determined, or IP, Ganj Dareh, and Tepe Hissar are a clade relative to Hajji
908 Firuz. In all these fitting topologies, the ~10,000-year-old radiocarbon date of the Ganj Dareh
909 individuals sets a lower bound on the split time between IP and Hajji Firuz / Tepe Hissar, which is
910 pre-agriculturalist. This suggests that the Iranian-related ancestry in IP is not due to an eastward
911 agriculturalist expansion.

912

913 But in fact, the admixture graph analysis reported above is not an adequate exploration of the
914 problem. Although absolute fits of the best models found are good (WR = 2.5 SE), the parsimony
915 constraint allowing only 3 admixture events precluded correct modeling of basal Eurasian ancestry
916 shared by all Middle Eastern groups (Lazaridis et al. 2016) or of the Indus Periphery group itself, for
917 which a more complex 3-component admixture model was proposed (Narasimhan et al. 2019).
918 Concerned that this oversimplification could be causing our search to miss important classes of
919 models, we explored *qpAdm* models for the Indus Periphery group further, following the “distal”

920 protocol with “rotating” outgroups outlined by Narasimhan *et al.* (2019) and using the dataset and
921 outgroups (“right” populations) from that study. All sites available for analyses were used, following
922 Narasimhan *et al.* (the “useallsnps: YES” option). The combined Indus Periphery group we analyzed
923 included 7 individuals from Shahr-i-Sokhta and 3 individuals from Gonur (3 individuals were
924 removed from the Narasimhan *et al.* 2019 dataset due to potential contamination with modern
925 human DNA and low coverage). We removed one individual from the Ganj Dareh Neolithic group as
926 potentially contaminated, and one 2nd or 3rd degree relative was removed from the Anatolia
927 Neolithic group, see the dataset composition in **Table S4**. We note that no “distal” *qpAdm* models
928 were tested for the combined Indus Periphery group by Narasimhan *et al.* (2019), and individuals
929 from this group were modeled one by one (Table S82 from Narasimhan *et al.* 2019), which
930 potentially reduced the sensitivity of the method.

931

932 A model “Indus Periphery = Ganj Dareh Neolithic + Onge (ASI)” was strongly rejected for the Indus
933 Periphery group of 10 individuals with a p -value = 2×10^{-15} , and a model that was shown to be fitting
934 for all Indus Periphery individuals modeled one by one by Narasimhan *et al.* (Ganj Dareh Neolithic +
935 Onge (ASI) + West Siberian hunter-gatherers (WSHG)) was rejected for the grouped individuals with
936 a p -value = 0.0044. In contrast, a model “Indus Periphery = Ganj Dareh Neolithic + Onge (ASI) +
937 WSHG + Anatolia Neolithic” was not rejected based on the $p > 0.01$ threshold used in Narasimhan *et al.*
938 (p -value was marginal but passing at 0.03) and produced plausible admixture proportions for all
939 four sources that are confidently above zero: $53.2 \pm 5.3\%$, $28.7 \pm 2.1\%$, $10.5 \pm 1.3\%$, $7.7 \pm 2.9\%$,
940 respectively (**Table S5**). The same “distal” model albeit with Anatolia Neolithic always in higher
941 proportion was found as one of the simplest models (or the only simplest model) fitting the data for
942 many other groups from Iran and Central Asia explored by Narasimhan *et al.* (2019): Aligrama2_IA
943 (13% Anatolia Neolithic), Barikot_H (21%), BMAC (26%), Bustan_BA_o2 (15%), Butkara_H (24%),
944 Saidu_Sharif_H_o (12%), Shahr_I_Sokhta_BA1 (19%), SPGT (23%) (**Table S5** gives a compendium of
945 distal modeling results by Narasimhan *et al.*). When we modeled Indus Periphery individuals
946 separately, as in Narasimhan *et al.* (2019), the simplest two-component model “Ganj Dareh Neolithic
947 + Onge (ASI)” was rejected for 5 of 10 individuals (at least ~315,000 sites were genotyped per
948 individual), including the individual I8726 used for the admixture graph analysis in Shinde *et al.*
949 (**Table S6**). The model “Ganj Dareh Neolithic + Onge (ASI)” was not rejected only for individuals with
950 fewer than 141,000 sites genotyped, suggesting that this result is attributed not to population
951 heterogeneity, but to lack of power.

952

953 These *qpAdm* results show that the parsimony assumption that was made when constructing the
954 admixture graph analysis in Shinde *et al.* (2019) is contradicted by f -statistic evidence, and indeed
955 Narasimhan *et al.* themselves showed this when they presented a distal *qpAdm* model that was
956 more complex (Ganj Dareh Neolithic + Onge (ASI) + WSHG) than the one used for constraining the
957 admixture graph model comparison (Ganj Dareh Neolithic + Onge (ASI)). Another line of evidence
958 used to support the principal historical conclusion by Shinde *et al.* was a series of f_4 -statistic cladality
959 tests following correction of allele frequencies using an admixture model “target group = Iranian
960 farmer + Anatolia Neolithic + Onge (ASI)”, with a great majority of tests supporting the deepest
961 position of the Iranian farmer ancestry component in the Indus Periphery group within the Iranian
962 farmer clade (Shinde *et al.* 2019). However, the model used for allele frequency correction (Ganj
963 Dareh Neolithic + Onge (ASI) + Anatolia Neolithic) was simpler than the 4-component model and
964 different from the 3-component model for the Indus Periphery group suggested by Narasimhan *et al.*
965 (2019), which is a weakness of that analysis. A valuable direction for future work would be to
966 repeat this analysis with a 4-component allele frequency correction model (Ganj Dareh Neolithic +
967 Onge (ASI) + WSHG + Anatolia Neolithic), although that is beyond the scope of the present study,
968 which simply aims to re-visit the reported analyses and test if they fully support their inferences by
969 ruling out alternative explanations.

970

971 To explore how the parsimony constraint influences results, we allowed 4 admixture events in the
972 eight-population graph (**Table S1**). Among 4,000 resulting graphs (one from each *findGraphs* run),
973 443 were distinct topologically, and 270 had WRs between 2 and 3 SE, i.e., fitted the data well. We

974 explored 35 topologies with LL scores in a narrow range between 9.3 (the best value) and 13.3. In
975 **Figure S7b** we show four graphs with four admixture events that model the Indus Periphery group as
976 a mixture of three or four sources, with a significant fraction of its ancestry derived from the Hajji
977 Firuz Neolithic or Tepe Hissar Chalcolithic lineages including both Iranian and Anatolian ancestries.
978 The fits of these models are just slightly different (e.g., LL = 11.7 vs. 9.3, both WRs = 2.4 SE) from that
979 of the best-fitting model (**Figure S7a**), and similar to that of the published graph. Besides these four
980 illustrative graphs, dozens of topologies with very different models for the Indus Periphery group fit
981 the data approximately equally well, suggesting that there is no useful signal in this type of
982 admixture graph analysis when the parsimony constraint is relaxed (this finding is similar to that in
983 our reanalysis of the dog admixture graph in Bergström *et al.* 2020, where relaxation of the
984 parsimony constraint identified equally well fitting admixture graphs that were very different with
985 regard to their inferences about population history). These results show that at least with regard to
986 the admixture graph analysis, a key historical conclusion of the study (that the predominant genetic
987 component in the Indus Periphery lineage diverged from the Iranian clade prior to the date of the
988 Ganj Dareh Neolithic group at ca. 10 kya and thus prior to the arrival of West Asian crops and
989 Anatolian genetics in Iran) depends on the parsimony assumption, but the preference for three
990 admixture events instead of four is hard to justify based on archaeological or other arguments.

991
992 Why did the Shinde *et al.* 2019 admixture graph analysis find support for the IP Iranian-related
993 lineage being the first to split, while our *findGraphs* analysis did not? The Shinde *et al.* 2019 study
994 sought to carry out a systematic exploration of the admixture graph space in the same spirit as
995 *findGraphs*—one of only a few papers in the literature where there has been an attempt to do so—
996 and thus this qualitative difference in findings is notable. We hypothesize that the inconsistency
997 reflects the fact that the deeply-diverging WSHG-related ancestry (Narasimhan *et al.* 2019) present
998 in the IVC (Indus Valley Civilization genetic grouping, which is the same group as Indus Periphery) at
999 a level of ca. 10% was not taken into account explicitly neither in the admixture graph analysis nor in
1000 the admixture-corrected f_4 -symmetry tests also reported in Shinde *et al.* (2019). The difference in
1001 qualitative conclusions may also reflect the fact that the Shinde *et al.* study was distinguishing
1002 between fitting models relying on a LL difference threshold of 4 units (based on the AIC). As
1003 discussed above, AIC is not applicable to admixture graphs where the number of independent model
1004 parameters is topology-dependent even if the numbers of groups and admixture events are fixed,
1005 and models compared with AIC should have the same number of parameters. Moreover, model
1006 comparisons with AIC do not account for the variability across SNPs, unlike the bootstrap model-
1007 comparison method we use. Thus, the analysis by Shinde *et al.* was over-optimistic about being able
1008 to reject models that were in fact plausible using its admixture graph fitting setup.

1009
1010 The archaeological and linguistic implications of the Shinde *et al.* study are important, and there are
1011 several avenues available for further attempting to distinguish historical scenarios using f -statistics
1012 that are outside the scope of a methodological study like this one. Some of our observations that are
1013 most challenging for the conclusions of Shinde *et al.* are those related to the graphs with four
1014 admixture events in **Figure S7b** that fit the Iranian farmer-related ancestry in the Indus Periphery
1015 group as deriving partially from the Hajji Firuz Neolithic or Tepe Hissar Chalcolithic-related lineages.
1016 *qpAdm* is able to use information from distal outgroups (such as WSHG) not included in the
1017 admixture graph modeling exercise revisited here. Leveraging this information might be able to
1018 obtain constraints that would further test the key historical conclusions from Shinde *et al.* Non- f -
1019 statistic-based methods could also be informative. Finally, we emphasize that the f_4 -statistic cladality
1020 tests correcting for the Anatolian farmer-related and Onge-related admixture in the Indus Periphery
1021 grouping do continue to provide support for the historical conclusion of Shinde *et al.* (these analyses
1022 reject models where the Tepe Hissar or Hajji Firuz groups share genetic drift with the Indus
1023 Periphery individuals), with the caveat that they do not correct for the WSHG admixture.

1024
1025
1026 **Librado *et al.* 2021.** In contrast to the other studies revisited in our work, the admixture graph
1027 published by Librado *et al.* (2021) was inferred automatically using *OrientAGraph*. Models with three

1028 (Fig. 3b in that study) and zero to five (Ext. Data Fig. 5a-d) admixture events were shown. The
1029 dataset included 10 populations (9 horse populations and donkey as an outgroup) and was based on
1030 7.4 million polymorphic transversion sites with no missing data at the group level. We observed that
1031 some groups used for the *OrientAGraph* and *qpAdm* analyses were very broad geographically and
1032 temporally (see Table S1 in the original study), and thus we tested two alternative group
1033 compositions: the original one and a streamlined one. In the latter case we included individuals from
1034 one archaeological site and one archaeological period per group: the Botai, C-PONT, DOM2, ELEN,
1035 and NEO-ANA groups were modified in this way, and the CWC, LP-SFR, Tarpan, and TURG groups
1036 were left with the composition used in the original paper (Table S7). In addition, 7 individuals with
1037 missing data proportion exceeding 80% were removed from the analysis, affecting the donkey
1038 outgroup, DOM2, and NEO-ANA groups (Table S7). Since among all possible f_3 -statistics for the 10
1039 populations three were negative (using all sites available for each population triplet, “useallsnps:
1040 YES”), we applied the upgraded algorithm for calculating f -statistics, which removed sites with only
1041 one chromosome genotyped in any non-singleton population, resulting in the following site counts
1042 for the original and modified population compositions: 11,092 and 1,767,419 sites, respectively. The
1043 very low number of sites available in the former case is due to the fact that all individuals are
1044 pseudohaploid, and that two groups (the donkey outgroup and NEO_ANA) are composed of two
1045 individuals, a high-coverage one and a low-coverage one. Thus, just sites genotyped in both donkey
1046 individuals and in both NEO_ANA individuals were kept. Considering this problem, we focused on
1047 the modified group composition only. We tested a range of model complexities (from 3 to 9 gene
1048 flows) and performed 1,000 *findGraphs* topology search iterations per model complexity class.
1049

1050 Unlike all the other admixture graphs we re-evaluate in this study whose fits to the data were
1051 evaluated in the published studies using *qpGraph*, the topologies published in Librado *et al.* 2021
1052 (with 3 to 5 admixture events) were not evaluated for statistical goodness-of-fit, and in fact fit the f -
1053 statistic data so poorly that even simple statistics show they cannot be correct (Figure 5b, Figure
1054 S8a,c,e, Table S1). In this case, the approach of using *findGraphs* to identify alternative topologies
1055 with the same number of admixture events that fit the data better is meaningless, as both the
1056 published models and the alternative models do not have enough degrees of freedom to
1057 accommodate the complexity present in the real data; all models are guaranteed to be wrong. In
1058 particular, we found that WR of the published model with 3 admixture events is 23.9 SE (Figure S8a).
1059 In this complexity class *findGraphs* found 22 topologically diverse models that fit significantly better
1060 than the published one (Table 1, Table S1), but nevertheless have extremely poor absolute fits (from
1061 16.2 to 21.3 SE, see a temporally plausible example in Figure S8b). In the complexity class with 4
1062 admixture events, no model fitting better than the published one was found; however, five
1063 alternative models fitting not significantly worse than the published one had lower WR (10 or 12 SE
1064 vs. 14.1 SE, Figure S8d). The WR of the published model with 5 admixture events was 6.9 SE (Figure
1065 S8e); just two models fitting nominally better and 223 models fitting non-significantly worse than
1066 the published model and having similar or higher WR were found (Table 1, Table S1). These results
1067 suggest that while *OrientAGraph* was often (but not always) able to find the same tentative global
1068 likelihood optimum as *findGraphs*, neither 3 nor 5 admixture events are enough to explain the data
1069 since nearly all the groups are admixed.
1070

1071 For this reason, we moved to topology searches in more complex model spaces incorporating 6 to 9
1072 admixture events. Temporally plausible models with even a modest fit (WR between 3 and 4 SE)
1073 were encountered only among models with 8 and 9 admixture events (Figure S8j-r). In the
1074 complexity class with 8 admixture events, 5 such temporally plausible fitting models were found,
1075 with WRs ranging from 3.4 to 3.9 SE (all these models are shown in Figure S8j-l). In the complexity
1076 class with 9 admixture events, 11 such models were found, with WRs ranging from 3.4 to 4.0 SE (all
1077 these models are shown in Figure S8m-r).
1078

1079 Librado *et al.* 2021 discussed the following inferences relying fully or partially on their published
1080 admixture graphs reported in that study (Table 2): 1) NEO-ANA-related admixture is absent in
1081 DOM2; 2) DOM2 and C-PONT are sister groups (they form a clade); 3) there is no gene flow

1082 connecting the CWC and the cluster associated with Yamnaya horses and horses of the later
1083 Sintashta culture whose ancestry is maximized in the Western Steppe (DOM2, C-PONT, TURG); 4)
1084 there was gene flow from a deep-branching ghost group to NEO-ANA; and 5) Tarpan is a mixture of a
1085 CWC-related and a DOM2-related lineage.

1086
1087 The simplest temporally plausible and best-fitting (WR = 3.4 SE) model we found (modified group
1088 composition, 8 admixture events, **Figure 5b**, **Figure S8j**, upper panels) supports inferences 2 and 4,
1089 and is incompatible with inferences 1, 3, and 5 (**Table 2**). This newly found model can be interpreted
1090 as follows. There is a trifurcation of three deep lineages: a lineage maximized in Western and Central
1091 Europe (up to 100% of ancestry in a Late Paleolithic group from France, LP_SFR), a Western Steppe-
1092 specific lineage (up to 55% in TURG), and a Tarpan-specific lineage (22% in Tarpan). Western and
1093 Central European horses, represented by LP-SFR, by the majority ancestry in horses found in the
1094 Corded Ware culture context (CWC), and by the majority ancestry in wild Neolithic Anatolian horses
1095 (NEO_ANA), contributed about half of the ancestry in the Western Steppe groups TURG, C-PONT,
1096 and DOM2. The other half of ancestry in the Western Steppe groups is represented by the Western
1097 Steppe-specific lineage. That lineage also contributed about 50% of ancestry in wild horses from the
1098 Yana Upper Paleolithic site (ELEN), and the other half of ELEN's ancestry is derived from an even
1099 deeper lineage. The Botai group is modeled as a mixture of European horses (69%) and Siberian
1100 horses (31% ELEN-related ancestry). In contrast to Librado *et al.* (2021), Tarpan is modeled as a
1101 mixture of its specific lineage (22%) and a DOM2-related group (78%), and CWC also received
1102 ancestry (21%) from a DOM2-related group. All the populations included in the model except for
1103 LP_SFR are admixed, and there is evidence of substantial genetic influence from a lineage that was
1104 eventually maximized in the Western Steppe (although it did not necessarily originate there) in the
1105 ELEN and Botai groups. We consider this model to be plausible from both temporal and geographical
1106 perspectives.

1107
1108 We are not arguing here that our 8-admixture-event model represents the true history; in fact, it is
1109 highly unlikely to be the truth, given how large the space of all possible admixture events is and how
1110 much admixture evidently occurred relating all these groups (which makes finding the unique truly
1111 fitting model extremely unlikely based on *f*-statistic fitting, see the results on simulated data in
1112 **Figure 2b**). We have also not attempted in any way to replicate the admixture graph exploration
1113 procedure performed in the Librado *et al.* study; the graph fitting procedure was quite different
1114 from ours, based on *OrientAGraph* optimization rather than *findGraphs* optimization, and a Block
1115 Jackknife procedure with a different genome block size for determining standard errors (4 Mbp in
1116 our protocol and ca. 500 kbp in the Librado *et al.* study). Regardless of how the graph was obtained,
1117 it is valuable for providing readers with guidance about which topological features of the graphs are
1118 meaningful and stable, and which are less certain, especially—as in the case of the admixture graph
1119 presented in the paper—when some features of the presented model cannot be right, as evident by
1120 the WR of 6.9 in the published model for 5 admixture events. Our set of 16 temporally plausible and
1121 fitting (WR < 4 SE) models with 8 or 9 admixture events (**Figure S8j-r**) is consistent with some
1122 features of the published graph being stable: the features (2) that DOM2 and C-PONT are sister
1123 groups, and (4) that there was a gene flow from a deep-branching ghost group to NEO-ANA (**Table**
1124 **2**).

1125
1126 Equally important, however, is our finding that there are plausible models that are inconsistent with
1127 other inferences in Librado *et al.* (**Table 2**). For example, 13 of these 16 models are inconsistent with
1128 the suggestion that there was no gene flow connecting the CWC and the cluster maximized in the
1129 Western steppe (DOM2, C-PONT, TURG) (**Figure S8j-r**). In the 8-admixture-event best-fitting
1130 plausible model (**Figure 5b**, **Figure S8j**, upper panels), CWC actually derives appreciable ancestry
1131 from the early domestic horse lineage DOM2 associated with the Sintashta culture to the exclusion
1132 of the more distant Yamnaya-associated TURG and C_PONT horses. This scenario presents a parallel
1133 to the one observed in humans, with individuals associated with the CWC receiving admixture from
1134 Steppe pastoralists albeit in different proportions: ~75% for humans, versus ~20% in horses. These
1135 models specifying a substantial Steppe horse contribution to CWC horses would weaken support for

1136 the inference in Librado *et al.* that “Our results reject the commonly held association between
1137 horseback riding and the massive expansion of Yamnaya steppe pastoralists into Europe around
1138 3000 BC”. We are not aware of other lines of evidence in the paper (apart from the fitted admixture
1139 graph) that support the claim of no Yamnaya horse impact on CWC horses.

1140

1141 Another example of a feature of the published graph that turned out to be unstable is the model for
1142 the Tarpan horse. Only 8 of 16 temporally plausible and fitting models (**Figure S8j-r**) support the
1143 conclusion by Librado *et al.* that the Tarpan is a mixture of a DOM2-related and a CWC-related
1144 lineage. The other 8 models suggest that Tarpan is a mixture of a deep lineage and a DOM2-related
1145 lineage (**Figure 5b, Figure S8j**, upper panels), echoing a hypothesis that Tarpan may be a hybrid with
1146 Przewalski’s horses not represented in the admixture graph (Librado *et al.* 2021).

1147

1148 Again, we are not arguing here that our fitting alternative model is right—indeed we are nearly
1149 certain it is wrong in important aspects—but we are merely pointing out that the complexity of the
1150 admixture graph space means that qualitatively quite different conclusions are compatible with the
1151 genetic data. Other aspects of the Librado *et al.* study, most notably the dramatic geographic
1152 expansion of the DOM2 modern domestic horse lineage after 4000 years ago in association with the
1153 Sintashta culture which is the most extraordinary finding of Librado *et al.*, are in no way challenged
1154 by our results.

1155

1156

1157 **Hajdinjak *et al.* 2021.** The admixture graph inferred by Hajdinjak *et al.* was constructed manually on
1158 the basis of a SNP set derived from in-solution enrichment of two SNP panels (1240K + a further
1159 million transversion polymorphisms discovered as polymorphic within one or two sub-Saharan
1160 African individuals or among archaic humans) and incorporated 11 groups and 8 admixture events
1161 (Figure 2d in the original study). The published graph has no clear outgroup since the deepest
1162 branch (Denisovan) is admixed. This property of the graph makes automated graph space
1163 exploration difficult. We explored two topology classes: 1) 11 groups with 8 admixture events, the
1164 original SNP set, Denisovan assigned as an outgroup only at the stage of generating random starting
1165 graphs (gene flows to/from the Denisovan branch were allowed at the topology optimization step);
1166 and 2) 12 groups with 8 admixture events, chimpanzee added and the original SNP set changed due
1167 to the zero missing rate condition, and chimpanzee assigned as an outgroup at both algorithm
1168 stages (**Table S1**). For both graph complexity classes, two topology search settings were tested: 1)
1169 either no additional constraints were applied beyond the outgroup constraints described above, or
1170 2) the Vindija Neanderthal and Mbuti were allowed to have no admixture events in their history, and
1171 the Denisovan lineage was allowed to have up to one admixture event in its history (these
1172 constraints were in line with the model in the original study and with literature on the genetic
1173 history of archaic humans, e.g., Prüfer *et al.* 2014). The composition of the groups matched that in
1174 the original study, as did the parameter settings for *qpGraph*, with the exception of “least squares
1175 mode”, which was used in the original study, but not in our analysis. “Least squares mode”
1176 computes LL scores without taking into account the f -statistic covariance matrix, and we confirmed
1177 that changing this parameter does not qualitatively change our results. Since no f_3 -statistics were
1178 negative when all sites available for each population triplet were used (the “useallsnps: YES” option),
1179 we did not use the upgraded algorithm for calculating f_3 -statistics on pseudo-diploid data. We
1180 summarize results across 2,000 to 4,000 independent iterations of the *findGraphs* algorithm (**Table**
1181 **S1**).

1182

1183 When chimpanzee was not included into the analysis and no topology constraints were applied,
1184 nearly all newly found models turned out to be distinct (3,996 of 4,000), nearly all (96.8%) fit
1185 nominally better and 15.9% fit significantly better than the published model (**Table S1**), and absolute
1186 fits of 91.3% of novel models are good (WR < 3 SE). Similar results were obtained when the topology
1187 search algorithm was constrained: nearly all (89.5%) of 1,999 newly found models fit nominally
1188 better and 26% fit significantly better than the published model (**Table S1**).

1189

1190 When chimpanzee was set as an outgroup and no topology constraints were applied, the picture
1191 remained similar. Nearly all newly found models turned out to be distinct (1,996 of 2,000), and a
1192 very large fraction of them (56.8%) fit significantly better than the published model (**Table S1**);
1193 16.4% of novel models demonstrated $WR < 3 SE$. Similar results were obtained when the topology
1194 search algorithm was constrained: most (71.4%) newly found models fit nominally better and 15.7%
1195 fit significantly better than the published model (**Table 1, Table S1, Figure 4**), which has a poor
1196 absolute fit on this set of sites and groups ($WR = 4.8 SE$, **Figure 5c, Figure S9**). The statistics
1197 described above and the fact that LL scores on all sites lie outside of the LL distribution on resampled
1198 datasets (**Figure 4**) suggest that models in this complexity class are overfitted, but the published
1199 topology emerged as fitting relatively poorly.

1200

1201 Overfitting arises naturally during manual graph construction as performed in many studies (not
1202 only in Hajdinjak *et al.* (2021), but also in Fu *et al.* 2016, Skoglund *et al.* 2016, Yang *et al.* 2017, Posth
1203 *et al.* 2018, McColl *et al.* 2018, Moreno-Mayar *et al.* 2018, Tambets *et al.* 2018, van de Loosdrecht *et al.*
1204 2018, Flegontov *et al.* 2019, Sikora *et al.* 2019, Wang *et al.* 2019, Lipson *et al.* 2020, Shinde *et al.*
1205 2019, Yang *et al.* 2020, and Wang *et al.* 2021). The graph grew one group at a time, and each newly
1206 added group was mapped on to the pre-existing skeleton graph as unadmixed or as a 2-way mixture.
1207 This imposed constraints on the model-building process. Another constraint imposed was the
1208 requirement that all intermediate graphs have good absolute fits (WR below 3 or 4 SE). When the
1209 model-building process is constrained in a particular path and fits of all intermediates are required
1210 to be good, unnecessary admixture events are often added along the way, and the resulting graph
1211 belongs to a complexity class in which models are overfitted and many alternative models fit equally
1212 well. There is no single obviously correct order of adding branches to a growing graph. For example,
1213 the Kostenki and Sunghir lineages were included into the initial graph (Fig. S6.1 in the original study)
1214 as unadmixed lineages, and their admixture status was not revisited at subsequent steps (unlike that
1215 of Tianyuan and Ust'-Ishim), except for adding the archaic gene flow common for non-Africans. For
1216 that reason, the published graph differs from many alternative better-fitting and temporally
1217 plausible graphs where the Kostenki and Sunghir lineages are modeled as more complex mixtures
1218 (**Figure S9**).

1219

1220 Hajdinjak *et al.* (2021)'s published graph had the following notable features that were interpreted by
1221 the authors and used to support some conclusions of the study (**Table 2**): 1) there are gene flows
1222 from the lineage found in the ~45,000-43,000-years-old Bacho Kiro Initial Upper Paleolithic (IUP)
1223 individuals to the Ust'-Ishim, Tianyuan, and GoyetQ116-1 lineages; 2) the ~35,000-years-old Bacho
1224 Kiro Cave individual BK1653 belonged to a population that was related, but not identical, to that of
1225 the GoyetQ116-1 individual; and 3) the Vestonice16 lineage is a mixture of a Sunghir-related and a
1226 BK1653-related lineage. To assess if these features are supported by our re-analysis, we focused on
1227 our most constrained *findGraphs* run: with chimpanzee set as an outgroup and with the topology
1228 constraints applied at the topology search step. We identified 1,421 topologies fitting nominally or
1229 significantly better than the published model and satisfying the constraints and moved on to inspect
1230 50 best-fitting topologies for temporal plausibility (all of them fitting significantly better than the
1231 published model). All non-African individuals included in the model are Upper Paleolithic and their
1232 dates are not drastically different in relative terms: from ca. 45 kya (thousand years before present)
1233 for some Bacho Kiro IUP individuals (Hajdinjak *et al.* 2021) to ca. 30 kya for the Vestonice16
1234 individual (Fu *et al.* 2016). Nevertheless, we considered most gene flows from later-attested to
1235 earlier-attested lineages as temporally implausible (for instance, GoyetQ116-1 (~35 kya) → Ust'-
1236 Ishim (~44 kya), GoyetQ116-1 (35 kya) → Bacho Kiro IUP (45-43 kya), Kostenki14 (38 kya) → Ust'-
1237 Ishim (44 kya), Sunghir III (34.5 kya) → Tianyuan (40 kya), Vestonice16 (30 kya) → Tianyuan (40 kya))
1238 since they imply great antiquity of the later-attested lineages, e.g., >40 kya for the Vestonice16
1239 lineage, and even greater antiquity for the related lineages such as Sunghir III and Kostenki14.

1240

1241 Of the 50 topologies inspected, 32 were considered temporally plausible. Of those topologies, none
1242 supported feature 1 of the published admixture graph (there is no replication of the finding of gene

1243 flows from the Bacho Kiro IUP lineage specifically into all three of the Ust'-Ishim, Tianyuan, and
1244 GoyetQ116-1 lineages). One topology supported features 2 and 3, and partially supported feature 1
1245 (there was Bacho Kiro → GoyetQ116-1 gene flow, but no Bacho Kiro → Tianyuan and Bacho Kiro →
1246 Ust'-Ishim gene flows). A total of 17 topologies supported features 2 and 3 but were inconsistent
1247 with feature 1; and 14 topologies supported feature 3 only (**Table 2**). Best-fitting representatives of
1248 each of these topology classes are shown along with the published model in **Figure S9**. Considering
1249 topological diversity among models that are temporally plausible, conform to current knowledge
1250 about relationships between modern and archaic humans, and fit significantly better than the
1251 published model, we conclude that feature 3 is probably robust but other details of the fitted
1252 admixture graph in Hajdinjak *et al.* (Figure 2d of that study)—for example, gene flows to the Ust'-
1253 Ishim, Tianyuan and Goyet Q116-1 lineages from sources sharing drift exclusively with the Upper
1254 Paleolithic Bacho Kiro lineage—should not be interpreted as providing meaningful inferences about
1255 population history of Upper Paleolithic modern humans. For example, the upper right-hand
1256 alternative model plotted in **Figure S9c** supports features 2 and 3 but includes no gene flows from
1257 the Bacho Kiro IUP lineage.

1258
1259 A central finding of Hajdinjak *et al.* is that the Bacho Kiro IUP group shares more alleles with present-
1260 day East Asians than with Upper Paleolithic Holocene Europeans despite coming from Europe.
1261 Specifically, the study documents significantly positive statistics of the form $D(\text{an Asian group,}$
1262 $\text{Kostenki14; Bacho Kiro IUP, Mbuti})$ (Fig. 2b and Extended Data Fig. 5 in the original study). For
1263 example, $D(\text{Tianyuan, Kostenki14; Bacho Kiro IUP, Mbuti})$ is significantly positive ($D = 0.0032$, $SE =$
1264 0.0010 , $Z = 3.2$) on the dataset used for testing the twelve-population graphs (263,698 sites without
1265 missing data across all 12 groups). The same statistic is also significantly positive ($D = 0.0029$, $SE =$
1266 0.0006 , $Z = 4.4$) when all 1,312,292 non-missing sites in the population quadruplet are analyzed.
1267 Hajdinjak *et al.*'s interpretation of this observation, using the language from the abstract, is that
1268 "there was at least some continuity between the earliest modern humans in Europe [Bacho Kiro IUP]
1269 and later people in Eurasia [East Asians]".

1270
1271 However, a significant D -statistic can have multiple explanations. The statistic $f_4(\text{Tianyuan,}$
1272 $\text{Kostenki14; Bacho Kiro IUP, Mbuti})$ is fitted equally well by the published twelve-population
1273 admixture graph (Z-score for the difference between the observed and fitted statistics = 0.64) and
1274 by, for example, the lower left-hand graph in **Figure S9c** (Z-score = 0.94) reproduced in **Figure 5c**.
1275 Under the latter model that fits the data significantly better than the published model (p -value =
1276 0.02), the Bacho Kiro IUP and Tianyuan branches are not connected by a gene flow and do not
1277 receive gene flows from a third common source, but the common ancestor of Ust'-Ishim and all
1278 European Paleolithic lineages receives an 8% gene flow from a divergent modern human lineage
1279 splitting deeper than Bacho Kiro IUP and Tianyuan (**Figure 5c, Figure S9c**). This scenario or some
1280 version of it seems archaeologically and geographically plausible and is not disproven by any other
1281 line of genetic or non-genetic evidence of which we are aware. It could correspond to a scenario
1282 where a primary modern human expansion out of the Near East contributed serially to the major
1283 lineages leading to Bacho Kiro, then later East Asians, then Ust'-Ishim, and finally the primary
1284 ancestry in later European hunter-gatherers. This has a very different interpretation from the
1285 scenario of distinctive shared ancestry between the earliest modern humans in Europe such as
1286 Bacho Kiro IUP and later people in East Asia—to the exclusion of later European hunter-gatherers—
1287 that is suggested by the Hajdinjak *et al.* published graph.

1288
1289 We are not claiming that this specific alternative model is correct—indeed, it is almost certainly not
1290 the correct one given the topological complexity of the set of all admixture graphs consistent with
1291 the data—but the existence of it and many other models that fit the data makes it clear that we do
1292 not yet have a unique historical explanation for the excess sharing of alleles that has been
1293 documented between some Upper Paleolithic European groups (Bacho Kiro IUP, Hajdinjak *et al.*
1294 2021, and GoyetQ116-1, Yang *et al.* 2017 and Hajdinjak *et al.* 2021) and all East Asians.

1295
1296

1297 **Lipson et al. 2020.** The admixture graph in the original study (Lipson et al. 2020) was constructed
1298 manually based on a SNP set derived from the 1240K enrichment panel, and the final model was
1299 alternatively tested on the combined HumanOrigins sub-panels 4 and 5 (each ascertained on one
1300 African individual) or on sites ascertained as polymorphic in archaic humans. The final published
1301 model (Extended Data Fig. 4 in that study) is very complex (12 groups and 12 admixture events): it
1302 exists in a space of $\sim 10^{44}$ topologies of this complexity. We note that one admixture event was
1303 added by Lipson et al. (2020) to account for potential modern DNA contamination in ancient Shum
1304 Laka individuals, and removing it caused a negligible difference in the fit of the published model
1305 (**Table S1**). Thus, to decrease the complexity of the graph search space, we considered graphs with
1306 12 groups and 11 admixture events. Twenty-two f_3 -statistics for these 12 groups turned out to be
1307 negative (when the “useallsnps: YES” setting was used), and thus for exploring this graph complexity
1308 class we had to remove sites with only one chromosome genotyped in any non-singleton population
1309 (**Table S1**). The following constraints were applied during the topology search: chimpanzee was
1310 assigned as an outgroup at both stages of the process (while generating random starting graphs and
1311 while searching the topology space); Altai Neanderthal was required to be unadmixed; and non-
1312 Africans (French) were required to have at least one admixture event in their history. The
1313 composition of the groups we analyzed matched that in the original study. We summarize results
1314 across 2,000 independent iterations of the *findGraphs* algorithm.

1315
1316 All newly found models turn out to be distinct (2,000), and 11.9% fit nominally (but not significantly)
1317 better than the published model (**Table 1, Table S1, Figure 4**). Absolute fits of 36.7% of novel models
1318 are good (WR < 3 SE). Fits of the highest-ranking model and the published model are not significantly
1319 different according to the bootstrap model comparison method (p -value = 0.176). These metrics,
1320 along with the fact that LL scores on all sites lie outside of the LL distribution on resampled datasets
1321 (**Figure 4**), suggest that models in this complexity class, including the published model, are
1322 overfitted. Of the admixture graphs we re-evaluate in this study, Lipson *et al.* 2020 shares with
1323 Hajdinjak *et al.* 2021, Sikora *et al.* 2019, and Wang *et al.* 2021, evidence of being overfitted (**Figure**
1324 **4**).

1325
1326 We also wanted to check if overfitting would be found in the graph complexity classes
1327 corresponding to two simpler intermediate graphs from the original study (**Table S1**): 7 groups and 4
1328 admixture events (Figure S3.24 in that study) and 10 groups and 8 admixture events (Figure S3.25 in
1329 that study). The population composition of the dataset we used for this analysis was slightly
1330 different from the dataset used by Lipson *et al.*: the ancient South African hunter-gatherer group
1331 was replaced by a related group (present-day Juǀǀhoan North), and instead of the Shum Laka
1332 ancient group, only one high-coverage individual from the same group (I10871) was used. We
1333 summarize results across 2,000 or 10,000 independent iterations for each SNP set, for the small and
1334 large graphs, respectively. For 7 groups, we found 201 novel topologies fitting better than the
1335 published one, and for 10 groups we found nearly 9,000 such topologies (**Table S1**). In the latter case
1336 6.8% of newly found topologies fit significantly better than the published topology. For the more
1337 complex graph class with 10 groups and 8 admixture events we also found evidence of overfitting:
1338 the LL score of the published graph run on the full data is better than almost all the bootstrap
1339 replicates on the same data (it falls below the 5th percentile).

1340
1341 Below we discuss selected prominent features of the admixture graph published in the original study
1342 (that were interpreted by the authors and used to support some conclusions of the study) and the
1343 extent to which these features consistently replicate across the large number of fitting 12-
1344 population graphs with 11 admixture events (**Table 2**): 1) A lineage maximized in present-day West
1345 African groups (Lemane, Mende, and Yoruba) also contributed some ancestry to the ancient Shum
1346 Laka individual and to present-day Biaka and Mbuti; 2) another ancestry component in Shum Laka is
1347 a deep-branching lineage maximized in the rainforest hunter-gatherers Biaka and Mbuti; 3) “super-
1348 archaic” ancestry (i.e., diverging at the modern human/Neanderthal split point or deeper)
1349 contributed to Biaka, Mbuti, Shum Laka, Lemane, Mende, and Yoruba; and 4) a ghost modern
1350 human lineage (or lineages) contributed to Agaw, Mota, Biaka, Mbuti, Shum Laka, Lemane, Mende,

1351 and Yoruba. We identified 232 twelve-population topologies that fit nominally better than the
1352 published one, 34 best-fitting topologies (of 232) were manually assessed for temporal plausibility,
1353 and we focus on 30 topologies identified as temporally plausible and including a low-level
1354 Neanderthal contribution ($\leq 10\%$) in non-Africans (French). These 30 topologies are shown along with
1355 the published model in **Figure S10**.

1356
1357 In this set of alternative models, high topological diversity is observed (see an example in **Figure 5d**
1358 and further topologies in **Figure S10**). We classified the topologies as follows. If an ancestral lineage
1359 defined above (for example, a deep-branching lineage maximized in rainforest hunter-gatherers
1360 Biaka and Mbuti) exists in the graph, we compared the sets of populations where it is found in the
1361 published model and in the model examined. If there was no more than one population where the
1362 ancestry is expected according to the published model but not present, or present but not expected,
1363 we considered this feature of the published graph supported by the alternative graph. If no ancestral
1364 lineage meets the definition above, the feature of the published graph was considered not
1365 supported. In all other cases partial support for the feature was declared (**Figure S10**). Considering
1366 extreme cases, two alternative graphs completely lacked support for three features of the published
1367 graph (**Figure 5d**, **Figure S10c**), and one graph supported all four features of the published graph
1368 fully (**Figure S10q**, bottom panels). There are some graphs where defining two distinct ancestral
1369 lineages maximized in West Africans and in Mbuti and Biaka (features 1 and 2) is essentially
1370 impossible since all or nearly all Africans are modeled as a mixture of at least two deep lineages (see
1371 graph no. 4, **Figure S10d**). In some graphs there is no single lineage specific to rainforest hunter-
1372 gatherers (Biaka, Mbuti, and Shum Laka) since the primary ancestries in these groups form
1373 independent deep branches in the African graph (see graph no. 2 in **Figure 5d** and graph no. 16 in
1374 **Figure S10j**, bottom panels). The ghost modern and super-archaic gene flows to Africans also had no
1375 universal support in the set of alternative graphs we examined (see, for example, **Figure 5d** and
1376 **Figure S10c**).

1377
1378 Considering the high degree of topological diversity among models that are temporally plausible,
1379 conform to known findings about relationships between modern and archaic humans, and fit
1380 nominally better than the published model, we conclude that the features from the original study
1381 are not supported by our re-analysis (**Table 2**). As in the case study above, the published manually
1382 constructed model is a representative of a large class of models that are equally well fitting to the
1383 limits of our resolution. This situation may be attributed to 1) overfitting and/or to 2) the lack of
1384 information in the dataset (in the combination of groups and SNP sites) and/or to 3) inherent
1385 limitations of f -statistics, when distinct topologies predict identical f -statistics.

1386
1387 In reconsidering the findings of Lipson *et al.* 2020 it is important to keep in mind that analysis of
1388 allele frequency correlation statistics is not the only type of information that can be used to make
1389 inferences about population relationships in deep time. Other methodologies have provided
1390 important insights into deep African population history, and the model-building in Lipson *et al.* 2020
1391 was guided in an informal way by these other lines of evidence. For example, unknown archaic
1392 lineages admixing into some African populations were hypothesized through identification of deeply
1393 splitting haplotypes that are too long to have been freely mixing with other haplotypes in present-
1394 day populations for all of their history (Hammer *et al.* 2011, Lachance *et al.* 2012, Speidel *et al.*
1395 2019). Similarly, analysis of haplotype divergence times of pairs of populations has been used to
1396 provide evidence of an early radiation of modern human lineages maximized today in southern
1397 African hunter-gatherers, Mbuti rainforest hunter-gatherers, and the great majority of other
1398 present-day populations; and a later split of lineages related to East African hunter-gatherers, West
1399 African agriculturalists, and non-Africans, which is a feature of the Lipson *et al.* model (Campbell and
1400 Tishkoff 2008, Mallick *et al.* 2016). Notably, some alternative models we found do not contradict the
1401 above-mentioned results and are profoundly different from the published model at the same time
1402 (see, for example, **Figure 5d**). These constraints are not enough, however, to provide evidence for all
1403 the topological details of the Lipson *et al.* 2020 admixture graph highlighted in this section, or for
1404 other features of Lipson *et al.* 2020 that were not invoked in the previous literature and newly

1405 proposed in that study, such as the “ghost modern” lineage splitting around the same time as the
1406 lineages leading to southern African hunter-gatherers and central African rainforest hunter-
1407 gatherers and mixing in highest proportion to Ethiopian hunter gatherers and to a lesser proportion
1408 to West Africans, and the “basal West African” lineage that contributes uniquely to Shum Laka.
1409 Many of the models that emerged as good fits in our admixture graph building exercise as the
1410 published one did not share some of these features (**Figure S10**).

1411
1412 The high diversity of well-fitting admixture graph models that satisfy known constraints relating
1413 diverse African populations highlights the need for further research based on multiple lines of
1414 genetic analysis (in addition to allele frequency correlation patterns) to obtain further insights into
1415 deep African history. Our results particularly highlight the mystery around the highly distinctive
1416 genetic ancestry of the Shum Laka individuals themselves, who represent the newly reported data in
1417 the Lipson *et al.* 2020 study, and a highly important set of genetic datapoints that was not available
1418 prior to the study. The ancestral relationships of these four individuals to both rainforest hunter-
1419 gatherers, and to the primary lineage in present-day West Africans, remains an open question, one
1420 whose resolution promises meaningful new insights into modern human population history.

1421
1422
1423 **Wang et al. 2021.** The admixture graph inferred by Wang *et al.* 2021 was constructed manually on
1424 the basis of a SNP set derived from the 1240K enrichment panel. We focused our analysis on the
1425 final graph (Extended Data Figure 6 in Wang et al. 2021, 12 groups and 8 admixture events) and on
1426 two simpler intermediates in the model building process (Figures S13-9 and S13-10a in Wang et al.
1427 2021). To simplify the latter two models further, we removed a low-level gene flow (1%) from a
1428 West European hunter-gatherer-related lineage (Loschbour) to the Mongolia Neolithic group, which
1429 resulted in negligible LL differences (0.5 and 2.4 log-units, respectively). Thus, using *findGraphs* we
1430 explored the following topology classes: 9 groups with 4 admixture events, 10 groups with 5
1431 admixture events, and 12 groups with 8 admixture events (**Table S1**). The composition of the groups
1432 matched that in the original study. We summarized results across 2,000 independent iterations of
1433 the *findGraphs* algorithm for each topology class. In the case of the most extensive population set
1434 (12 groups), three f_3 -statistics turned out to be negative (when the “useallsnps: YES” setting was
1435 used), and thus for exploring this graph complexity class we had to remove sites with only one
1436 chromosome genotyped in any non-singleton population (**Table S1**). For this complexity class, we
1437 also applied several constraints on the graph space exploration process all of which were shared
1438 with the Wang *et al.* graphs: the Denisovan genome was assigned as an outgroup in the random
1439 starting graphs, but not at the topology search stage; up to one admixture event was allowed in the
1440 history of the Denisovan group; no admixture events were allowed in the history of Mbuti,
1441 Loschbour, and Onge; and the (Denisovan, (Mbuti, (Loschbour, Onge))) branching order was
1442 required.

1443
1444 For each topology class we found hundreds to thousands of topologically unique graphs fitting
1445 nominally better than the published models (**Table 1, Table S1**). For both simple topology classes, no
1446 model fitting significantly better than the published one was found (**Table S1**). However, the final
1447 published model fits the data significantly worse than 12.6% of newly found models of the same
1448 complexity (**Table 1, Table S1**). The fact that many topologically diverse models had good absolute
1449 fits (65%, 55%, and 15% of distinct newly found graphs with 9, 10, and 12 groups, respectively, had
1450 $WR < 3 SE$) suggests that admixture graph models in these complexity classes are overfitted. Further
1451 evidence of overfitting comes from the poor fits of the published model on bootstrap-resampled
1452 datasets as compared to their fits on all sites (**Figure 4**).

1453
1454 An important feature of the published graphs in Wang *et al.* 2021 that was remarked upon in the
1455 study is admixture from a source related to Andamanese hunter-gatherers that is almost universal in
1456 East Asians, occurring in the Jomon, Tibetan, Upper Yellow River Late Neolithic, West Liao River Late
1457 Neolithic, Taiwan Iron Age, and China Island Early Neolithic (Liangdao) groups (**Table 2**). For
1458 example, the abstract states “Hunter-gatherers from Japan, the Amur River Basin, and people of

1459 Neolithic and Iron Age Taiwan and the Tibetan Plateau are linked by a deeply splitting lineage that
1460 probably reflects a coastal migration during the Late Pleistocene epoch.” We performed 2,000
1461 *findGraphs* iterations and obtained 1,778 distinct topologies satisfying all the constraints, nearly all
1462 of them (1,724) fitting nominally better than the published model, and 12.6% fitting significantly
1463 better (**Table S1**). The models were ranked by LL, and 56 highest-ranking topologies, all of them
1464 fitting significantly better than the published one, were assessed for temporal plausibility (models
1465 with gene flows from a later group to Tianyuan dated to 40 kya were removed), and 20 topologies
1466 were considered temporally plausible (all of them are shown in **Figure S11**). According to these
1467 topologies, 0 to 2 East Asian groups had a fraction of their ancestry derived from a source
1468 specifically related to Onge, and 19 topologies included gene flows from the European (Loschbour)-
1469 related branch to all 8 East Asian groups (**Figure S11**). The inferred topological relationships among
1470 East Asians are variable in this group of 20 models, and we decided to apply further constraints that
1471 guided model ranking and elimination by Wang *et al.*, based on considerations from archaeological
1472 evidence, Y chromosome haplogroup divergence patterns, and population split time estimation.

1473
1474 The constraints that are not based on correlation of allele frequencies across populations that Wang
1475 *et al.* applied and that we applied in our re-examination are as follows. First, combined evidence
1476 from archaeology, linguistics, and genetics (a closely shared Y chromosome haplogroup) suggests
1477 that the present-day Tibetan Plateau population harbors a substantial proportion of ancestry from a
1478 large-scale migration from the Neolithic farming groups from the Upper and Middle Yellow River
1479 (Chen *et al.* 2015, Lu *et al.* 2016, Zhang *et al.* 2019). These arguments and radiocarbon dates favor
1480 the following branching order of predominant ancestry components: (Mongolia East N, (China Upper
1481 YR LN, Nepal Chokhopani)). Second, evidence from archaeology, linguistics and genetics suggests
1482 that the expansion of Austronesian speakers and the peopling of Taiwan was from southeast coastal
1483 China to Taiwan and Southeast Asia, but not from Taiwan to mainland China (Bellwood 2011, Gray
1484 and Jordan 2000, Ko *et al.* 2011). These arguments make a China Island EN → Taiwan IA gene flow
1485 direction plausible and make the opposite direction of flow less likely. Third, in the original study
1486 MSMC cross-coalescence rates were computed for a few pairs of present-day proxies for the ancient
1487 groups, and it was argued that they impose constraints on the graph topology. The inferred
1488 coalescence date for the Tibetan and Ulchi groups was slightly younger than the Tibetan-Ami and
1489 Tibetan-Atayal dates (Fig. S13-1 in the original study), suggesting that the Nepal Chokhopani and
1490 Mongolia East N group may share ancestral source populations more recently than these two groups
1491 and Taiwan IA. We note that it was not clear in the original paper if the difference in coalescence
1492 dates is statistically significant, the finding was clearer in MSMC than in MSMC2 analysis, and there
1493 was no attempt to calculate expected cross-coalescence profiles using these methods from models
1494 incorporating many gene flows. Nevertheless, we applied this constraint as well in an attempt to
1495 understand whether, if we used a constraint system similar to that in Wang *et al.*, we would obtain
1496 results that agreed with respect to the finding of Onge-related admixture ubiquitous among East
1497 Asian groups.

1498
1499 Applying these three additional constraints, we identified two models (among the 56 ones subjected
1500 to manual inspection) that satisfied all of them. The highest-ranking of those models is shown in
1501 **Figure 5e** and **Figure S11c** (lower panels), and it includes a 13% (deeply) European-related gene flow
1502 to the common ancestor of all East Asians, and gene flows from the Onge-related branch to just two
1503 East Asian groups: Nepal Chokhopani and China WLR LN. This model fits the data significantly better
1504 than the published model (p -value = 0.028). We do not claim that this is the correct model (indeed
1505 we are almost certain that it is not given the high degeneracy of fitting models), but it is not
1506 obviously wrong and differs in qualitatively important ways from the published one.

1507
1508 The Wang *et al.* 2021 admixture graph provides an illuminating example that helps us to understand
1509 the value added by admixture graph construction. The admixture graph construction process in
1510 Wang *et al.* followed a philosophy of not relying entirely on the allele frequency correlation data
1511 (not treating the genetic data as independent to explore how much new insight could come from
1512 genetic data alone). Instead, the study integrated other lines of genetic evidence as well as linguistic

1513 and archaeological insights explicitly into the admixture graph construction process, with the goal of
1514 identifying models consistent with multiple lines of evidence. The fact that after this procedure a
1515 fitting graph was obtained is not of great interest, as it is essentially always possible to obtain a fit to
1516 allele frequency correlation data when enough admixture events are added. The important question
1517 is whether any of the emergent features of the graph that were not applied as constraints in the
1518 construction process—for example the evidence of ubiquitous Andamanese-related gene flow
1519 throughout East Asia suggesting a coastal route expansion that admixed with an interior route
1520 expansion proxied by Tianyuan—were stably inferred. Our analysis does not come to this finding
1521 consistently among well-fitting and plausible admixture graphs. We conclude that an important
1522 feature of the published graph, i.e. variable levels of Andamanese-related ancestry found in all East
1523 Asians except for Siberians (Mongolia Neolithic) and the Upper Paleolithic Tianyuan (Fig. 2 in (Wang
1524 et al. 2021)), is not supported by f -statistic analysis alone (**Table 2**), and indeed we are not aware of
1525 a single feature of the Wang *et al.* 2021 admixture graph that is stably inferred beyond the
1526 constraints applied to build it.

1527

1528

1529 **Sikora et al. 2019.** Two admixture graphs inferred by Sikora *et al.* (2019) were constructed manually
1530 based on a SNP set derived from whole-genome shotgun data and incorporated 12 or 13 groups and
1531 10 admixture events (Extended Data Figure 3f in the original study). One graph was focused on West
1532 Eurasians, and the other one on East Eurasians, and both included a Neanderthal, a Denisovan, and
1533 an African group (Dinka). Although the chimpanzee outgroup was not included in the original graphs,
1534 we added it as it drastically constrains the topology search space. The following additional
1535 constraints were applied at the *findGraphs* model optimization stage: the Neanderthal and African
1536 groups were unadmixed and the Denisovan group had no more than one admixture event in its
1537 history. These three constraints match the features of the published graph. We also repeated
1538 topology searches without constraining the admixture status of the Neanderthal, Denisovan, and
1539 Dinka. Since no f_3 -statistics were negative when all sites available for each population triplet were
1540 used (the “useallsnps: YES” option), we did not apply the algorithm that allows unbiased calculation
1541 of f_3 -statistics on pseudo-diploid data at the expense of loss of analyzed SNPs.

1542

1543 In contrast to most other published graphs discussed above, gene flows in the graphs inferred by
1544 Sikora *et al.* do not have equal standing: four low-level gene flows (0-1%) connect the Neanderthal
1545 lineage to Upper Paleolithic lineages (Kostenki, Sunggir, Yana, Ust'-Ishim in the 'Western' graph and
1546 Sunggir, Yana, Mal'ta, Ust'-Ishim in the 'Eastern' graph). We repeated each topology search under
1547 two alternative settings: either keeping the number of admixture events at 10 to match the
1548 published graphs, or at 6 to match simplified versions of the published graphs lacking these low-level
1549 Neanderthal gene flows. We performed that modification to simplify the search space and to
1550 alleviate the overfitting problem which becomes severe if 10 gene flows across the graph are
1551 allowed (**Table S1**). Here we compare LL and WR for the original published models and their
1552 simplified versions: the Western graph including chimpanzee (LL = 65.7, WR = 3.32 SE) vs. its
1553 simplified version (LL = 76.5, WR = 3.78 SE) and the Eastern graph including chimpanzee (LL = 85.3,
1554 WR = 3.11 SE) vs. its simplified version (LL = 102.4, WR = 4.16 SE). In both cases we found no
1555 statistically significant differences in model fits (relying on the bootstrap model comparison
1556 method). In summary, topology search was repeated under 8 settings: for the Western or Eastern
1557 graphs, with no constraints on the admixture status or with the constraints specified above, and with
1558 10 or 6 gene flows (**Table S1**). Below we focus on results for constrained models with 6 admixture
1559 events. In contrast, **Figure 4** and **Table 1** show results for constrained Western graphs with 10
1560 admixture events.

1561

1562 In the case of the constrained Western graphs with 6 admixture events, 1,000 topology search
1563 iterations were performed, 894 distinct topologies were found, 4 models fit significantly better, and
1564 151 models fit nominally better than the published one (**Table 1**, **Table S1**). We inspected those 155
1565 topologies and identified 29 topologies (**Figure S12**) that are temporally plausible and include no
1566 non-canonical gene flows from archaic groups such as Denisovan or gene flows ghost archaic →

1567 non-Africans. Sikora *et al.* came to the following striking conclusion relying on the Western
1568 admixture graph (**Table 2**): the Mal'ta (MA1_ANE) lineage received a gene flow from the Caucasus
1569 hunter-gatherer (CaucasusHG_LP or CHG) lineage. However, in our *findGraphs* exploration this
1570 direction of gene flow (CHG → Mal'ta) was supported by two of the 29 topologies, and the opposite
1571 gene flow direction (from the Mal'ta and East European hunter-gatherer clade to CHG) was
1572 supported by the remaining 27 plausible topologies (**Figure S12**). The highest-ranking plausible
1573 topology (**Figure 5f**) has a fit that is not significantly different from that of the simplified published
1574 model (p -value = 0.392). We note that the gene flow direction contradicting the graph by Sikora *et al.*
1575 was supported by a published *qpAdm* analyses (Lazaridis *et al.* 2016, Narasimhan *et al.* 2019), and
1576 *qpAdm* is not affected by the same model degeneracy issues that are the focus of this study.
1577 Considering the topological diversity among models that are temporally plausible, conform to robust
1578 findings about relationships between modern and archaic humans, and fit nominally better than the
1579 published model, we conclude that the direction of the Mal'ta-CHG gene flow cannot be resolved by
1580 admixture graph analysis (**Table 2**).

1581
1582 Some important conclusions based on the Eastern graph also do not replicate across all plausible
1583 admixture graphs (**Table 2**). In the case of the constrained Eastern graphs with 6 admixture events,
1584 4,446 topology search iterations were performed, and 2,785 distinct topologies were found. Only 3
1585 topologies fit significantly and 13 nominally better than the published one (p -value for the highest-
1586 ranking newly found model vs. the simplified published model = 0.112), and 9.8% of topologies fit
1587 not significantly worse than the published one (**Table 1, Table S1**). Of the topologies belonging to
1588 these groups, we inspected 116 best-fitting ones and identified 97 topologies that are temporally
1589 plausible and include no gene flows from archaic groups such as Denisovan or ghost archaic → non-
1590 Africans that are qualitatively different from the gene flows that are currently widely accepted. The
1591 Sikora *et al.* Eastern admixture graph had the following distinctive features that were used to
1592 support some conclusions of the study (**Table 2**): 1) the Mal'ta (MA1_ANE) and Yana (Yana_UP)
1593 lineages receive a gene flow from a common East Asian-associated source diverging before the ones
1594 contributing to the Devil's Cave (DevilsCave_N), Kolyma (Kolyma_M), USR1 (Alaska_LP), and Clovis
1595 (Clovis_LP) lineages; 2) European-related ancestry in the Kolyma, USR1, and Clovis lineages is closer
1596 to Mal'ta than to Yana; 3) the Devil's Cave lineage received no European-related gene flows, and
1597 Kolyma has less European-related ancestry than ancient Americans (USR1 and Clovis). Only feature 2
1598 was universally supported by all the 97 plausible alternative models fitting significantly better,
1599 nominally better, or not significantly worse than the simplified published model, while feature 3 was
1600 supported by 83 of 97 plausible models, and feature 1 was supported by 28 of 97 plausible models
1601 (**Table 2**). We plotted 14 plausible graphs as examples of topologies supporting all three features,
1602 two features, or one feature of the published graph (**Figure S13**). We note that all the Eastern graphs
1603 discussed here, both the published and alternative ones, have relatively poor absolute fits with WR
1604 above 4 or 5 SE. Increasing the number of gene flows to 10 allowed us to reach much better
1605 absolute fits (with WR as low as 2.42 SE), but that resulted in high topological diversity (on a par with
1606 some other case studies discussed above). In the case of the constrained Eastern graphs with 10
1607 admixture events, 1,000 topology search iterations were performed, and 1,000 distinct topologies
1608 were found. Of these topologies, 13.2% fit significantly better, 30% nominally better, and 17.6%
1609 non-significantly worse than the published model (p -value for the highest-ranking newly found
1610 model vs. the published model < 0.002) (**Table S1**).

1611

1612

1613 **A Proposed Protocol for Using Admixture Graph Fitting in Genetic Studies**

1614

1615 Admixture graphs represent a conceptually powerful framework for thinking about demographic
1616 history, but the practice of manually constructing a small number of complex models without
1617 exploring admixture graph space in an automated way can lead to overconfidence in the validity of
1618 these models. An ideal outcome of an admixture graph model exploration exercise would be the
1619 identification of a model or a group of topologically very similar models which fits the data well and
1620 significantly better than all alternative models with the same number of admixture events; however,

1621 this is almost never achieved for graphs with more than eight populations and three admixture
1622 events in our experience, and even this approach can lead to potentially unstable results as relaxing
1623 the assumption of parsimony (that fewer admixture events is more likely) can lead to qualitatively
1624 quite different equally well fitting topologies as in our reanalysis of the Bergström *et al.* and Shinde
1625 *et al.* datasets. Most of the examples of admixture graphs in eight recently published studies we
1626 revisited do not fit this ideal pattern, as we were able to identify many topologically different
1627 alternative models that could not easily be rejected based on temporal plausibility or other
1628 constraints (**Figures S5-S13**). In particular, for all studies except Shinde *et al.* 2019 (under a strict
1629 parsimony assumption however), we identified admixture graphs that were not significantly worse
1630 fitting than the published ones, and with topological features that were different in qualitatively
1631 important ways. There were also some more encouraging findings of the exercise we performed to
1632 re-evaluate published models. For example, at least one of the key inferences about population
1633 history relying on the admixture graph modeling were stable for all analyzed models for the Lazaridis
1634 *et al.*, Librado *et al.*, Hajdinjak *et al.*, Shinde *et al.* 2019 (under the parsimony assumption), and
1635 Sikora *et al.* (simplified Eastern graph) studies. The existence of some stable features in these graphs
1636 helps to point the way toward a protocol that we believe should be applied in all future studies that
1637 use admixture graph fitting exercises to support claims about population history.

1638
1639 We propose the following tentative protocol to identify features of fitting admixture graphs that are
1640 stable enough to be used to make inferences about population history.

- 1641
1642 1. For a given combination of populations, carry out an initial scan using *findGraphs* to identify
1643 reasonable parameter values for the number of allowed admixture events (the graph complexity
1644 class). For example, run *findGraphs* allowing between zero and eight admixture events (100
1645 algorithm iterations per graph complexity class), each iteration saving one or a few best-fitting
1646 outcomes. The smallest number of admixture events that yields models where the (negative) LL
1647 score or the worst *f*-statistic residual is lower than some threshold can then be explored more
1648 deeply by running more iterations of *findGraphs*.
- 1649
1650 2. Run *findGraphs* on the determined complexity class, where some of the resulting graphs should
1651 be inspected manually to determine whether they could in principle be historically plausible models.
1652 Implausible models (for example, models where a very ancient population appears to be admixed
1653 between two modern populations) can be filtered out by imposing topological constraints. If no or
1654 only a few graphs remain, *findGraphs* can be run again under these constraints. This can be repeated
1655 until one or more graphs with an acceptable LL score or worst residual has been identified. At this
1656 stage we apply the bootstrap method to determine whether the best-fitting graph is significantly
1657 better than the next best-fitting graph. If it is not, we identify a set of graphs which are not clearly
1658 worse than the best-fitting graph by performing the bootstrap model comparison for many model
1659 pairs.
- 1660
1661 3. Researchers should compare the resulting graphs to each other with the goal of identifying
1662 common features. Although *ADMIXTOOLS 2* includes automated tools for cataloguing common
1663 topological features (Suppl. Methods), we found a manual approach to be valuable, as the fitted
1664 parameters (especially admixture proportions) are as important for this task as graph topology.
- 1665
1666 4. Once a set of fitting graphs and stable topological features shared between them is identified,
1667 researchers should carry out a *findGraphs* exploration of the space of graphs with one additional
1668 admixture event. If inferences are stable even when fitting graphs with one more level of complexity
1669 than the graphs with the minimal number of admixture events needed to fit the data, this increases
1670 confidence in the inferences. Furthermore, the addition of a new population may introduce crucial
1671 information to an existing set of populations, which can change the space of fitting topologies in a
1672 profound way, as in our reanalysis of the data from Bergström *et al.* 2020 (**Figure 5a**, **Figure S2**).
1673 Thus, it is advisable that the topology optimization procedure is repeated on several alternative

1674 population sets, in addition to considering models that allow an additional admixture event beyond
1675 the minimum required for parsimony, to explore if inferences about topology change qualitatively.

1676

1677 5. Admixture graphs fitted with f -statistics do not distinguish between time and population size as
1678 the two sources of genetic drift, and many different complex genetic histories for a set of
1679 populations can result in the exact same expected f -statistics. This provides an important
1680 opportunity to further constrain a model fitting procedure. Methods that take advantage of
1681 information from the site frequency spectra (*mom2*, *fastsimcoal*, Kamm et al. 2019, Excoffier et al.
1682 2013) or derived site patterns, a special case of site frequency spectra (*Legofit*, Rogers 2019), can
1683 supply alternative information not captured by f -statistics (further information can come from
1684 methods that fit haplotype divergence patterns such as *MSSMC* (Schiffels and Durbin 2014) and
1685 *SMC++* (Terhorst et al. 2017), or inferences based on fitted gene trees such as *RELATE* (Speidel et al.
1686 2019), and *ARGweaver* (Hubisz et al. 2020, Hubisz and Siepel 2020)). These tools are too
1687 computationally intensive to explore a large number of models, but the advantages of the different
1688 approaches can be combined by first identifying a set of candidate models using *findGraphs*, and
1689 then testing these candidate models with other methods. This approach is also expected to deal
1690 with overfitting since different data types almost always include different variable site sets.

1691

1692 We believe that researchers should only begin to make strong claims about population history with
1693 admixture graphs once a protocol such as we propose is applied.

1694

1695 We see the guidelines above as analogous in spirit to the protocols that were introduced in medical
1696 genetics at a time when a reproducibility crisis was found in the field of candidate gene association
1697 studies. Many studies looking for risk factors for common complex diseases resulted in publications
1698 with marginally significant p -values without correcting for the multiple hypothesis testing that was
1699 implicitly performed due to many candidate genes being tested and only those with significant
1700 findings being published. Unsurprisingly, most of these claims failed to replicate in follow-up studies
1701 in independent sets of samples (Ioannidis 2005, Border et al. 2019, Collins et al. 2012, Duncan et al.
1702 2019). The human medical genetic community addressed this challenge by coming together to
1703 support a rigorous set of commonly accepted standards for declaring genome-wide statistical
1704 significance, such as the requirement that p -values be corrected for the effective number of
1705 independent common variants in the genome and requiring correcting for the known confounders
1706 of population structure and undocumented relatedness among individuals (Hirschhorn and Daly
1707 2005).

1708

1709

1710 Conclusion

1711

1712 Sampling admixture graph space is a useful method for modeling population histories, but finding
1713 robust and accurate models can be challenging. As we demonstrated by revisiting a handful of
1714 published admixture graphs and re-analyzing the same datasets used to fit them, f -statistic and,
1715 more generally, allele frequency data alone are usually insufficient for building accurate graph
1716 models, making it necessary to incorporate other sources of evidence. This provides a challenge to
1717 previous approaches for automated model building. We investigated several published admixture
1718 graph models and, in nearly all cases, found many alternative models, some of which are historically
1719 and geographically plausible but contradict conclusions that were derived from the published
1720 models. To conduct these analyses, we developed a method for automated admixture graph
1721 topology optimization which can incorporate external sources of information as topological
1722 constraints. This method is developed in a framework called *ADMIXTOOLS 2*, which aside from
1723 admixture graph modeling, implements many other methods for population history inference based
1724 on f -statistics. In the process of revisiting published admixture graphs we found a well-fitting model
1725 for the history of dogs that is substantially different from the published model (Bergström et al.
1726 2020) and is strikingly congruent with the known history of relevant human lineages. We also found
1727 a novel admixture graph for domestic and wild ancient horses (Librado et al. 2021) that is

1728 substantially different from the published model, fits the data significantly better, and is
1729 geographically and historically plausible. These alternative graphs, however, have some of the same
1730 challenges as the published ones: they are almost certainly oversimplified relative to the true
1731 histories, and they exist in a large space of admixture graphs with meaningful topological differences
1732 that fit the allele frequency correlation data equally well. An important topic for future work should
1733 be to test these new alternative models as well as the previously published models as hypotheses
1734 with newly reported ancient samples and additional lines of genetic, archaeological, and other forms
1735 of analysis to obtain further clarity about population history.

1736

1737 It is important to recognize that the key concern we have highlighted in this study—the fact that
1738 there can often be multiple different topologies that are equally good fits to the allele frequency
1739 correlation patterns relating a set of populations—does not invalidate the use of allele frequency
1740 correlation testing in many other contexts in which it has been applied to make inferences about
1741 population history. For example, negative f_3 -statistics (“admixture” f_3 -statistics) continue to provide
1742 unambiguous evidence for a history of mixture in tested populations, and f_4 and D symmetry
1743 statistics remain a powerful way of evaluating whether a tested pair of populations is consistent
1744 with descending from a common ancestral population since separation from the ancestors of two
1745 groups used for comparison. The *qpWave* methodology remains a fully valid generalization of f_4 -
1746 statistics, making it possible to test whether a set of populations is consistent with descending from
1747 a specified number of ancestral populations (which separated at earlier times from a comparison set
1748 of populations). In addition, the *qpAdm* extension of *qpWave*—which allows for estimating
1749 proportions of mixtures for the tested population under the assumption that we have data from the
1750 source populations for the mixture—remains a valid approach, unaffected by the concerns identified
1751 here. Instead of relying on a specific model of deep population relationships, *qpAdm* relies on an
1752 empirically measured covariance matrix of f_4 -statistics for the analyzed populations, which is highly
1753 constraining with respect to estimation of mixture proportions but can be consistent with a wide
1754 range of deep history models. All these methods are implemented in *ADMIXTOOLS 2*.

1755

1756 Finally, approaches that use admixture graphs to adjust for the covariance structure relating a set of
1757 populations without insisting that the particular admixture graph model that is proposed is true with
1758 can be useful, for example for the purpose of analyzing shared genetic drift patterns of a group of
1759 populations that derive from similar mixtures. One example was a study that attempted to test for
1760 different source populations for Neolithic migrations into the Balkans after controlling for different
1761 proportions of hunter-gatherer admixture (Mathieson et al. 2018). Another example was a study
1762 that attempted to study shared ancestry between different East African forager populations after
1763 controlling for different proportions of deeply divergent source populations (Lipson et al. 2022).
1764 However, with respect to the inferences about deep history produced by admixture graphs
1765 themselves, our results highlight the importance of caution in proposing specific models of
1766 population history that relate a set of groups.

1767

1768

1769 **Methods**

1770

1771 **Technical presentation of *ADMIXTOOLS 2* in the context of f -statistic modeling methods**

1772

1773 Much of the content that follows recapitulates theory presented in previous work, notably Reich et
1774 al. 2009, Green et al. 2010 and Patterson et al. 2012, but we summarize it here for coherence.

1775

1776 *f*-statistics

1777

1778 All *ADMIXTOOLS* programs are based on the statistics f_2 , f_3 , and f_4 , for population pairs, triplets,
1779 and quadruples, respectively.

1780

1781 f_2 quantifies the genetic drift separating two populations A and B . For a single SNP, it is given by
1782 $f_2(A, B) = \frac{1}{M} \sum_j (a_j - b_j)^2$, where a_j and b_j are the allele frequencies for SNP j in populations A
1783 and B . When allele frequencies are estimated using a small number of samples, this estimator of f_2
1784 will be biased upwards. An unbiased estimator of f_2 is given by
1785 $f_2 = \frac{1}{M} \sum_j (a_j - b_j)^2 - \frac{a_j(1-a_j)}{n_{A,j}-1} - \frac{b_j(1-b_j)}{n_{B,j}-1}$, where $n_{A,j}$ and $n_{B,j}$ are the observed allele counts in
1786 populations A and B .

1787
1788 $f_3(A; B, C) = \frac{1}{M} \sum_j (a_j - b_j)(a_j - c_j)$ is the covariance of the allele frequency differences between
1789 populations A and B , and the allele frequency differences between populations A and C (assuming
1790 that alleles are coded randomly, so that $a - b$ and $a - c$ are both 0 in expectation). Significantly
1791 negative values of $f_3(A; B, C)$ suggest that A is a mixture of sources related to B and C (although the
1792 converse does not hold: A might be admixed between B and C even if f_3 is positive).

1793
1794 $f_4(A, B; C, D) = \frac{1}{M} \sum_j (a_j - b_j)(c_j - d_j)$ is the covariance of the allele frequency differences
1795 between A and B , and the allele frequency differences between C and D . Significantly positive
1796 values of $f_4(A, B; C, D)$ (or equivalently significantly negative values of $f_4(A, B; D, C)$) reveal that A
1797 and B do not form a clade with respect to C and D , and that some of the drift separating A from C is
1798 shared with the drift separating B from D .

1799
1800 f_3 and f_4 can be written as linear combinations of f_2 statistics:

1801 $f_3(A; B, C) = \frac{1}{2}(f_2(A, B) + f_2(A, C) - f_2(B, C))$ (Eq. 1)

1802 $f_4(A, B; C, D) = \frac{1}{2}(f_2(A, D) + f_2(B, C) - f_2(A, C) - f_2(B, D))$ (Eq. 2)

1803 This implies that all f_3 - and f_4 -statistics can be computed from f_2 -statistics as long as they are
1804 defined on the same SNPs.

1805
1806 For revisiting published studies, we used the “*extract_f2*” function with the “*maxmiss*” argument set
1807 at 0, which corresponds to the “*useallsnps: NO*” setting in classic *ADMIXTOOLS*. It means that no
1808 missing data are allowed (at the level of populations) in the specified set of populations for which
1809 pairwise f_2 -statistics are calculated. For the values of the “*blgsize*”, “*adjust_pseudohaploid*”, and
1810 “*minac2*” arguments we use in our analyses, see **Table S1**. The “*blgsize*” argument sets the SNP
1811 block size in Morgans, and we used either the default value of 0.05 (5 cM), or 4,000,000 bp when a
1812 genetic map was not available. Genotypes of pseudo-haploid samples are usually coded as 0 or 2
1813 (i.e., they are, strictly speaking, pseudo-diploid), even though only one allele is observed. The
1814 “*adjust_pseudohaploid*” argument ensures that the observed allele count increases only by 1 for
1815 each pseudo-haploid sample. If “*TRUE*” (default), samples that do not have any genotypes coded as
1816 1 among the first 1,000 SNPs are automatically identified as pseudo-haploid. This leads to slightly
1817 more accurate estimates of f -statistics. Setting this parameter to “*FALSE*” treats all samples as
1818 diploid.

1819
1820 Another important argument (“*minac2=2*”) of the “*extract_f2*” function removes sites with only one
1821 chromosome genotyped in any non-singleton population and is needed for unbiased estimation of
1822 negative f_3 -statistics in non-singleton pseudo-haploid populations. In the absence of negative f_3 -
1823 statistics or pseudo-haploid populations, this argument has no influence on admixture graph log-
1824 likelihood scores. This algorithm for calculation of f -statistics triggered by the “*minac2=2*” argument
1825 is described below.

1826
1827 For $f_3(a; b, c)$, we compute the uncorrected numerator for each SNP, $(a - b) \times (a - c)$. We then
1828 subtract a bias correction factor at each SNP, $p(1 - p) / (ac - 1)$, which we only need for population a
1829 (because the other factors cancel out); p is the allele frequency, and ac is the observed allele count.
1830 In pseudo-haploid samples, $(ac - 1)$ would be zero and produce an error in any sites with only one
1831 observed allele. With the “*inbreed: NO*” setting in Classic *ADMIXTOOLS*, the smallest non-zero value

1832 for ac is 2, so the division by 0 problem is avoided, but the correction factor is slightly smaller than it
1833 should be. *ADMIXTOOLS2* adds only an allele count of 1 for each site in a pseudo-haploid sample
1834 (with the default option “*adjust_pseudohaploid = TRUE*”), so there can be cases where $ac = 1$. To
1835 imitate what the setting “*inbreed: NO*” in Classic *ADMIXTOOLS* is doing, ac is set to 2 at those sites
1836 (or the denominator is set to 1). There is still a small difference between Classic *ADMIXTOOLS* and
1837 *ADMIXTOOLS2* at other sites because each observed site adds 2 alleles in *ADMIXTOOLS* with the
1838 default setting “*inbreed: NO*”, but only 1 allele in *ADMIXTOOLS2* with the default setting
1839 “*adjust_pseudohaploid = TRUE*”, but for admixture graph fitting that does not matter. One solution
1840 to avoid biased correction factors is to only consider sites with ac of at least two, which is what the
1841 “*inbreed: YES*” setting in Classic *ADMIXTOOLS* does. The problem with this is that we cannot use
1842 populations with a single pseudo-haploid sample, which is often useful, and would only give
1843 misleading results if that population is admixed. The new option “*minac2=2*” in *ADMIXTOOLS2* is
1844 different from the “*inbreed: YES*” setting in Classic *ADMIXTOOLS* since it makes an exception for
1845 populations consisting of a single pseudo-haploid sample in that it sets ac to 2 at each site
1846 (denominator is set to 1) when computing the correction factor of those populations.

1847

1848

1849 *Fitting admixture graphs*

1850

1851 An admixture graph is a directed acyclic graph specifying the topology of the ancestral relationships
1852 among a set of populations. Each node in this graph represents a (present-day or ancient)
1853 population. Terminal nodes (also called leaf nodes) represent observed populations, while internal
1854 nodes represent unobserved ancestral populations. Modeling all observed populations as leaf nodes
1855 confers some robustness to drift specific to single populations and to genotyping errors. The edges
1856 connecting the populations are weighted and correspond either to the magnitude of genetic drift
1857 that has occurred along that branch (drift edges), or to the admixture proportions (admixture edges,
1858 where two edges point to the same node).

1859

1860 The goal of *qpGraph* is to test how well a given graph topology fits the observed f -statistics. This is
1861 achieved by varying the edge weights until the maximum likelihood fit is obtained. The following
1862 section describes the graph fitting in more detail.

1863

1864 First, for k populations, all $\frac{k(k+1)}{2} f_3$ -statistics of the form $f_3(O; X_1, X_2)$ are computed, where O is
1865 one of the k populations (typically an outgroup), and X_1 and X_2 are all pairs formed from the other
1866 populations (including pairs where $X_1 = X_2$). These f_3 -statistics can then be used to fit the graph
1867 and to compute the likelihood. The likelihood score of a graph is the dot product of the differences
1868 between the expected and observed f_3 -statistics, weighted by the inverse covariance matrix of f_3 -
1869 statistics:

1870

$$1871 \quad L(g) = -\frac{1}{2} (f_{3,obs} - f_{3,fit})' Q^{-1} (f_{3,obs} - f_{3,fit}) \quad (\text{Eq. 3})$$

1872

1873 Here, $f_{3,obs}$ are the observed f_3 -statistics and $f_{3,fit}$ are the fitted f_3 -statistics. Both are vectors of
1874 length $q = \frac{k(k+1)}{2}$ for k populations excluding the outgroup. Q is the $q \times q$ covariance matrix of f_3 -
1875 statistics, where the diagonal entries are the f_3 -statistic variances, and the off-diagonal entries are
1876 the covariances for all pairs of f_3 -statistics. Just like the variances (the squared standard errors), the
1877 covariances are estimated from the jackknife leave-one-block-out f_3 -statistics.

1878

1879 Finding the edge weights which maximize the likelihood score involves two nested optimization
1880 steps. The inner optimization finds the drift weights which maximize the likelihood score while fixing
1881 the admixture weights. The outer optimization finds the admixture weights which maximize the
1882 likelihood score, while optimizing the drift weights for each set of admixture weights. The inner

1883 optimization uses a quadratic programming solver to find the optimal drift weights, while the outer
1884 optimization uses a general purpose optimization algorithm to find the optimal admixture weights.
1885 While the gradient function in the outer optimization adjusts the admixture weights, the objective
1886 function iterates over the following steps:

1887

- 1888 1. Optimization of drift weights conditional on admixture weights
- 1889 2. Estimation of fitted f_3 -statistics
- 1890 3. Calculation of the graph likelihood using observed and fitted f_3 -statistics

1891

1892 These steps are repeated until convergence is reached and the likelihood score can no longer be
1893 improved by adjusting the admixture weights.

1894

1895 Step 1 optimizes the drift edge weights, while holding the admixture weights constant. All drift edge
1896 weights are required to be non-negative, which makes this a constrained quadratic programming
1897 problem (hence *qpGraph*). Additional upper and lower bounds can be specified for individual graph
1898 edges.

1899

1900 Step 2 turns the edge weights into fitted f_3 -statistics. To see how edge weights in an admixture
1901 graph translate to f_3 -statistics, it helps to first consider how they translate into f_2 -statistics for a pair
1902 of populations. Without any admixture events, there is exactly one path p connecting any two
1903 populations. The fitted f_2 -statistic ($f_{2,fit}$) is the sum of edge weights w_e along this path p connecting
1904 two populations. The fitted f_2 -statistic is the sum of edge weights w_e along this path:

$$f_{2,fit} = \sum_{e \in p} w_e$$

1905 In the presence of admixture events, two populations may be connected via multiple paths. Each
1906 admixture node that lies between the two populations increases the number of possible paths. The
1907 fitted f_2 -statistic for the two populations now becomes the weighted sum of all these paths, where
1908 the weight of each path is given by the product of all estimated admixture proportions w_a along this
1909 path ($\prod_{a \in p} w_a$):

$$1910 f_{2,fit} = \sum_{p \in P} \prod_{a \in p} w_a \sum_{e \in p} w_e \quad (\text{Eq. 4})$$

1911 The fitted f_2 -statistics are then used to obtain fitted f_3 -statistics using Eq. 1.

1912

1913 Step 3 uses the fitted and observed f_3 -statistics to estimate the likelihood score using Eq. 3.

1914

1915 Prior to these three steps, initial admixture weights are drawn randomly. To ensure that the end
1916 results do not depend on the random initialization, the whole optimization process is repeated
1917 multiple times with different random initial values. The original *ADMIXTOOLS* implementation
1918 retains only the results from the initial values resulting in the lowest absolute likelihood score. The
1919 new *ADMIXTOOLS* implementation provides an option to retrieve the results for all random
1920 initializations. This can be useful, as large fluctuations between different random initializations can
1921 be an indicator of an overparameterized or otherwise poorly fitting model.

1922

1923

1924 *Automated admixture graph inference*

1925

1926 To find graph topologies that could conceivably have given rise to the observed f -statistics, we start
1927 with a randomly generated graph with a fixed number of admixture events, apply a number of
1928 modifications to this graph, and evaluate each of the resulting graphs. We then pick the best-fitting
1929 graph and repeat this procedure until graph modifications no longer lead to improved scores. We
1930 use a number of random graph modifications, as well as targeted modifications which are informed
1931 by parameters obtained during the fitting of the current graph.

1932

1933 For the targeted modifications we change the optimization of a single graph from a constrained
1934 optimization problem, in which drift edges are constrained to be positive and admixture weights are
1935 constrained to be between zero and one, to an unconstrained optimization problem in which both
1936 types of parameters can take any real values. Rearranging the nodes adjacent to edges which were
1937 estimated to be negative results in an improved fit at a much higher rate than random graph
1938 adjustments.

1939
1940 The random modifications include (1) pruning and randomly re-grafting leaf nodes, (2) pruning and
1941 randomly re-grafting a set of connected nodes in the graph, (3) swapping the orientation of
1942 admixture edges, (4) shifting admixture edges, (5) re-rooting the graph, (6) combinations of two or
1943 more of any of these modifications.

1944
1945 The number of admixture events is not affected by the graph modifications described so far. A
1946 significant score improvement can often be achieved by adding a single admixture edge to several
1947 random positions in a graph. This is unsurprising since it increases the degrees of freedom of the
1948 original graph. However, picking the best fitting graph with one admixture edge added, and testing
1949 all graphs that result from removing a single admixture edge from that graph often results in a graph
1950 with the same number of admixture events and a better fit than the original graph. We employ this
1951 strategy whenever the regular graph modifications described above do not lead to any further
1952 improvements.

1953
1954 We keep track of the search tree of all previously evaluated graphs and their scores in order to not
1955 evaluate any graph more than once, and so that backtracking in the search space is possible in cases
1956 where no more local improvements can be identified. Nevertheless, multiple iterations with
1957 different random starting graphs are usually necessary to find graphs with good fits. The number of
1958 iterations needed to approach a global optimum depends on the size of the search space, but the
1959 optimal number of iterations is hard to estimate in practice.

1960
1961 For revisiting published studies, we used the following settings of the *findGraphs* algorithm:
1962 • *mutfuncs = namedList(spr_leaves, spr_all, swap_leaves, move_admixedge_once,*
1963 *flipadmix_random, place_root_random, mutate_n)*, a list of functions used to modify graphs.
1964 • *numgraphs = 10*, number of alternative graphs produced by randomly applying the mutation
1965 functions at the start of each generation.
1966 • *stop_gen = 10000*, total number of generations after which to stop.
1967 • *stop_gen2 = 30*, number of generations without LL score improvement after which to stop.
1968 • *plusminus_generations = 10*. If the best score does not improve after *plusminus_generations*
1969 generations, another approach to improving the score is attempted: A number of graphs
1970 with an additional admixture edge is generated and evaluated. The resulting graph with the
1971 best score is picked, and new graphs are created by removing any one admixture edge
1972 (bringing the number back to what it was originally). The graph with the lowest score is then
1973 selected. This approach often makes it possible to break out of local optima.
1974 • *opt_worst_residual = FALSE*. Optimize for lowest worst residual instead of best score.
1975 “FALSE” by default, because the LL score is generally a better indicator of the quality of the
1976 model fit, and because optimizing for the lowest worst residual is much slower since f_4 -
1977 statistics need to be computed.
1978 • *reject_f4z = 0*. If this is a number greater than zero, all f_4 -statistics with $|Z\text{-score}| > reject_f4z$
1979 will be used to constrain the search space of admixture graphs: Any graphs in which f_4 -
1980 statistics greater than *reject_f4z* are expected to be zero will not be evaluated.
1981 • *diag = 1e-04*. This argument is passed to the *qpgraph* function and determines the
1982 regularization term added to the diagonal elements of the covariance matrix of fitted branch
1983 lengths (after scaling by the matrix trace). Default is 0.0001.

- 1984 • *numstart* = 10. This argument is passed to the *qpgraph* function and determines the number
1985 of random initializations of starting weights (defaults to 10). Increasing this number will
1986 make the optimization slower but reduce the risk of not finding the optimal weights.
1987 • *lsqmode* = FALSE. This argument is passed to the *qpgraph* function. If set to “FALSE”, the
1988 inverse f_3 -statistic covariance matrix is not discarded by the algorithm.

1989 The arguments “*admix_constraints*” (constraints on the number of admixture events in the history
1990 of a given population), “*event_constraints*” (constraints on the branching order of specified
1991 lineages), and “*outpop*” (the population assigned as an outgroup) were set according to **Table S1**.
1992 Each *findGraphs* run was initiated by a random graph with a specified number of admixture events.
1993 Usually, the same topology constraints were applied at the stage of random graph generation and
1994 the topology search stage, for exceptions see **Table S1**.

1995
1996

1997 *Evaluating automated admixture graph inference through simulations*

1998

1999 We evaluated the performance of *findGraphs* by simulating genetic data under a large number of
2000 different admixture graph models, applying *findGraphs* to each simulated data set in three
2001 independent iterations, and comparing the resulting best graph across three iterations to the
2002 simulated graph. We applied *TreeMix* to the same simulated data for comparison. We simulated
2003 between 8 and 16 populations per graph, and between 0 and 10 admixture events. For each
2004 parameter combination, we simulated 20 different admixture graphs generated by the
2005 *random_admixturegraph* function. We counted both the fraction of random simulated graphs where
2006 the best inferred graph was identical to the simulated graph, as well as the fraction of random
2007 simulated graphs where the best inferred graph was either identical to the simulated graph or had a
2008 better score than the simulated graph. For models with a large number of admixture events the
2009 number of possible models is so large that it becomes increasingly likely that there will be some
2010 alternative models which fit the data better than the model under which the data were simulated.

2011

2012 We used *msprime* and the *msprime_sim* wrapper function in *ADMIXTOOLS 2* to simulate data for
2013 100,000 unlinked SNPs and 100 diploid samples per population for each admixture graph. The
2014 simulation parameters we chose were aimed at facilitating fast simulations of large numbers of
2015 informative SNPs rather than at being as realistic as possible. We therefore expect that the
2016 simulation results allow us to make comparisons across groups, but not that they are informative
2017 about the rate at which “true” models can be recovered in empirical data. We simulated under a
2018 constant mutation rate of 0.001 per site per generation, a constant haploid effective population size
2019 of 1000, with neighboring nodes in the graph separated by 1000 or more generations, and all
2020 admixture events occurring in discrete pulses of 50/50 proportions.

2021

2022 To allow for a fair comparison between *findGraphs* and *TreeMix*, we made sure that small
2023 differences in the way admixture graphs are modeled in *findGraphs* and in *TreeMix* were accounted
2024 for before testing graphs for identical topology. For example, *TreeMix* admixture graphs can have
2025 lineages terminating at an admixture node, whereas in *findGraphs* lineages always end at a ‘leaf’
2026 node with a single ancestor.

2027

2028

2029 *Comparing the fits of different admixture graphs*

2030

2031 We are interested in determining whether one admixture graph fits the data significantly better
2032 than another admixture graph, or whether an observed score difference $\Delta = S_1 - S_2$ can be
2033 attributed to variability across independent SNPs.

2034

2035 We first consider two admixture graphs with the same number of admixture events, where we can
2036 ignore the problem of comparing two models with different complexity. As in other bootstrap
2037 standard error calculations, we divide the genome into n blocks indexed by i , and we draw b sets of

2038 n blocks with replacement, indexed by j . We fit both graphs b times - once for each bootstrap set of
2039 SNP blocks. This results in a set of b score differences Δ_j . The bootstrap confidence interval for the
2040 difference in scores is given by the quantiles of the distribution of Δ_j . We also compute an empirical
2041 bootstrap p -value, testing the null hypothesis that two different graphs fit the data equally well. It is
2042 computed as $p = \max(\frac{1}{b}, 2\delta)$ (Boos 2003), where δ is either the fraction of $\Delta_j > 0$, or the
2043 fraction of $\Delta_j < 0$, whichever is smaller. The reason for applying bootstrap resampling, as opposed
2044 to jackknife resampling in this case, is that the distribution of score differences tends to have a high
2045 kurtosis, which can make jackknife estimates inaccurate. Simulating data under the null hypothesis is
2046 not straightforward in this case, because it involves finding two non-identical graphs which in
2047 expectation fit the data equally well. We decided to simulate under one graph and compare two
2048 graphs which are symmetrically related to the simulated graph (Figure 3). This confirmed that the p -
2049 value follows a uniform distribution under the null hypothesis.

2050

2051 Next, we consider comparisons of two graphs of different complexity. The problem here is that more
2052 complex graphs have more degrees of freedom which allow them to overfit the data better, without
2053 necessarily being any closer to the truth. To solve this problem, we introduce an out-of-sample
2054 likelihood score. The regular likelihood score is given by:

2055 $L(g) = -\frac{1}{2} (f_{3,obs} - f_{3,fit})' Q^{-1} (f_{3,obs} - f_{3,fit})$, with $f_{3,obs}$ and $f_{3,fit}$ defined on the same set of
2056 SNPs. The out-of-sample likelihood score is defined in the same way, except that $f_{3,obs}$ and $f_{3,fit}$ are
2057 defined on mutually exclusive sets of SNP blocks, thereby preventing any overfitting. The covariance
2058 matrix Q is defined on the same set of SNP blocks as $f_{3,fit}$. As described earlier, we use block-
2059 bootstrap to fit both graphs multiple times on different SNP blocks. In each bootstrap iteration, we
2060 use all SNP blocks which are not used in fitting the graph for estimating $f_{3,obs}$.

2061

2062

2063 *Admixture graph identifiability*

2064

2065 An edge in an admixture graph is unidentifiable, if small changes to the weight of this edge
2066 (admixture proportions in the case of an admixture edge, drift length in the case of a drift edge) do
2067 not necessarily lead to changes in expected f -statistics. This is the case if the small change in weight
2068 can be offset by small changes in other graph edges, leading to a situation where observed f -
2069 statistics can be explained by more than one weight estimate for that edge. To find unidentifiable
2070 edges, we derive the Jacobi matrix of the graph's system of f_2 equations (Eq. 4 applied to each
2071 population pair). In principle, whether or not a parameter is identifiable can depend on the values of
2072 all other parameters. However, in practice this is rarely the case, and so we draw values for all
2073 parameters from a uniform distribution, which gives us a Jacobi matrix with numeric values. We
2074 then determine the rank of the Jacobi matrix, along with the rank of all matrices that result from
2075 dropping a single column (a parameter corresponding to a graph edge). For identifiable edges, the
2076 rank of the full matrix will be greater than the rank of the reduced matrix, and for unidentifiable
2077 edges, the ranks will be the same.

2078

2079

2080 *Drawing conclusions from a large number of fitting models*

2081

2082 We developed several methods that aim to summarize a collection of graphs which all fit the data
2083 similarly well. By highlighting features which are observed repeatedly across graphs, it becomes
2084 possible to extract interpretable conclusions from an otherwise hard to interpret collection of
2085 possible models. These graph summaries identify features in each graph that can be compared to
2086 different graphs describing the same populations. We summarize each graph in several ways:

2087 (1) Admixture status of each population

2088 For each population, we count the total number of admixture events that is encountered
2089 along all paths to the root.

- 2090 (2) Order of population split events
2091 For each pair of population pairs, we determine if the most recent split of the first pair has
2092 occurred before or after the most recent split of the second pair, or whether the graph does
2093 not specify the order in which those splits occurred.
- 2094 (3) Proxy populations
2095 For each admixed population in a graph, we attempt to identify proxy sources: populations
2096 closest to the admixing populations. In contrast to the other approaches to summarizing
2097 graphs which are based only on the topology of each graph, this can also rely on information
2098 about the estimated graph parameters.
- 2099 (4) Cladality
2100 For each group of four populations, we test whether the graph implies that any f_d -statistic
2101 describing the relationship between the four populations is expected to be zero.
- 2102 (5) Node descendants
2103 Each internal node in an admixture graph is an ancestor to a specific set of leaf populations.
2104 An admixture graph can be characterized by the sets of leaf populations formed by the
2105 internal nodes. Multiple admixture graphs may be compared by counting the number of
2106 overlapping sets. This also makes it possible to quantify for each internal node in a single
2107 graph, how often a matching internal node can be found across a collection of alternative
2108 graphs, which is conceptually similar to bootstrap support values in phylogenetic trees.
- 2109 While these methods provide some help in comparing features across many graphs, they are not
2110 able to reliably answer the question whether the fitting graphs are relatively similar or dissimilar
2111 from each other, and whether they are similar to any particular graph. This is in part due to the fact
2112 that small topological changes involving populations of interest may be more relevant than similar
2113 topological changes involving only populations that are not the focus of the study.

2114
2115

2116 **Acknowledgements**

2117

2118 We thank Anders Bergström, Esther Brielle, Mateja Hajdinjak, Iosif Lazaridis, Pablo Librado, Mark
2119 Lipson, Vagheesh Narasimhan, Ludovic Orlando, Nick Patterson, Mary Prendergast, Jakob Sedig,
2120 Kendra Sirak, Pontus Skoglund, and Chuanchao Wang, for suggestions for how to improve specific
2121 analyses, and for conversations and critical comments. We thank Matthew Mah, Shop Mallick, Adam
2122 Micco, Nadin Rohland, Ron Pinhasi for help in generating additional data from an ancient DNA
2123 library from individual I8726 for which 1.24 million SNP capture data was generated and published in
2124 Narasimhan *et al.* 2019 and for which we report 2.6-fold shotgun data here (**Table S2**). P.F., P.C., and
2125 O.F. were supported by the Czech Ministry of Education, Youth and Sports (program ERC CZ, project
2126 no. LL2103) and by the Czech Science Foundation (project no. 21-27624S). P.F. and P.C. were also
2127 supported by the Czech Ministry of Education, Youth and Sports (Large Infrastructures for Research,
2128 Experimental Development and Innovations project "IT4Innovations National Supercomputing
2129 Center – LM2015070"; Inter-Excellence program, project no. LTAUSA18153). R.M. and D.R. were
2130 supported by grants from the National Institutes of Health (GM100233 and HG012287), the John
2131 Templeton Foundation (grant 61220), and the Allen Discovery Center program, a Paul G. Allen
2132 Frontiers Group advised program of the Paul G. Allen Family Foundation. D.R. was also supported by
2133 a private gift from J.-F. Clin and is an Investigator of the Howard Hughes Medical Institute.

2134

2135

2136 **References**

2137

- 2138 1. Bellwood P. The checkered prehistory of rice movement southwards as a domesticated cereal—
2139 from the Yangzi to the equator. *Rice*. 2011; 4:93–103. doi: 10.1007/s12284-011-9068-9
- 2140 2. Bergström A, Frantz L, Schmidt R, Ersmark E, Lebrasseur O, Girdland-Flink L, Lin AT, Storå J,
2141 Sjögren KG, Anthony D, Antipina E, Amiri S, Bar-Oz G, Bazaliiskii VI, Bulatović J, Brown D,
2142 Carmagnini A, Davy T, Fedorov S, Fiore I, Fulton D, Germonpré M, Haile J, Irving-Pease EK,
2143 Jamieson A, Janssens L, Kirillova I, Horwitz LK, Kuzmanovic-Cvetković J, Kuzmin Y, Losey RJ,

- 2144 Dizdar DL, Mashkour M, Novak M, Onar V, Orton D, Pasarić M, Radivojević M, Rajković D,
2145 Roberts B, Ryan H, Sablin M, Shidlovskiy F, Stojanović I, Tagliacozzo A, Trantalidou K, Ullén I,
2146 Villaluenga A, Wapnish P, Dobney K, Götherström A, Linderholm A, Dalén L, Pinhasi R, Larson G,
2147 Skoglund P. Origins and genetic legacy of prehistoric dogs. *Science*. 2020 Oct 30;370(6516):557-
2148 564. doi: 10.1126/science.aba9572.
- 2149 3. Boos DD. Introduction to the Bootstrap World. *Statist. Sci.* 18(2): 168-174 (May 2003). DOI:
2150 10.1214/ss/1063994971.
- 2151 4. Border R, Johnson EC, Evans LM, Smolen A, Berley N, Sullivan PF, Keller MC (2019) No support
2152 for historical candidate gene or candidate gene-by-interaction hypotheses for major depression
2153 across multiple large samples. *American Journal of Psychiatry* 176: 376-387.
- 2154 5. Chen FH, Dong GH, Zhang DJ, Liu XY, Jia X, An CB, Ma MM, Xie YW, Barton L, Ren XY, Zhao ZJ, Wu
2155 XH, Jones MK. Agriculture facilitated permanent human occupation of the Tibetan Plateau after
2156 3600 B.P. *Science*. 2015 Jan 16;347(6219):248-50. doi: 10.1126/science.1259172.
- 2157 6. Collins AL, Kim Y, Sklar P; International Schizophrenia Consortium, O'Donovan MC, Sullivan PF
2158 (2012) Hypothesis-driven candidate genes for schizophrenia compared to genome-wide
2159 association results. *Psychological Medicine* 42: 607-616.
- 2160 7. Duncan LE, Ostacher M, Ballon J (2019) How genome-wide association studies (GWAS) made
2161 traditional candidate gene studies obsolete. *Neuropsychopharmacology* 44: 1518-
2162 1523. Campbell MC, Tishkoff SA. African genetic diversity: implications for human demographic
2163 history, modern human origins, and complex disease mapping. *Annu Rev Genomics Hum Genet.*
2164 2008;9:403-33. doi: 10.1146/annurev.genom.9.081307.164258.
- 2165 8. Durvasula A, Sankararaman S. Recovering signals of ghost archaic introgression in African
2166 populations. *Sci Adv*. 2020 Feb 12;6(7):eaax5097. doi: 10.1126/sciadv.aax5097.
- 2167 9. Flegontov P, Altınışık NE, Changmai P, Rohland N, Mallick S, Adamski N, Bolnick DA,
2168 Broomandkoshbacht N, Candilio F, Culleton BJ, Flegontova O, Friesen TM, Jeong C, Harper TK,
2169 Keating D, Kennett DJ, Kim AM, Lamnidis TC, Lawson AM, Olalde I, Oppenheimer J, Potter BA,
2170 Raff J, Sattler RA, Skoglund P, Stewardson K, Vajda EJ, Vasilyev S, Veselovskaya E, Hayes MG,
2171 O'Rourke DH, Krause J, Pinhasi R, Reich D, Schiffels S. Palaeo-Eskimo genetic ancestry and the
2172 peopling of Chukotka and North America. *Nature*. 2019 Jun;570(7760):236-240. doi:
2173 10.1038/s41586-019-1251-y.
- 2174 10. Fu Q, Posth C, Hajdinjak M, Petr M, Mallick S, Fernandes D, Furtwängler A, Haak W, Meyer M,
2175 Mittnik A, Nickel B, Peltzer A, Rohland N, Slon V, Talamo S, Lazaridis I, Lipson M, Mathieson I,
2176 Schiffels S, Skoglund P, Derevianko AP, Drozdov N, Slavinsky V, Tsybankov A, Cremonesi RG,
2177 Mallegni F, Gély B, Vacca E, Morales MR, Straus LG, Neugebauer-Maresch C, Teschler-Nicola M,
2178 Constantin S, Moldovan OT, Benazzi S, Peresani M, Coppola D, Lari M, Ricci S, Ronchitelli A,
2179 Valentin F, Thevenet C, Wehrberger K, Grigorescu D, Rougier H, Crevecoeur I, Flas D, Semal P,
2180 Mannino MA, Cupillard C, Bocherens H, Conard NJ, Harvati K, Moiseyev V, Drucker DG, Svoboda
2181 J, Richards MP, Caramelli D, Pinhasi R, Kelso J, Patterson N, Krause J, Pääbo S, Reich D. The
2182 genetic history of Ice Age Europe. *Nature*. 2016 Jun 9;534(7606):200-5. doi:
2183 10.1038/nature17993.
- 2184 11. Gray RD, Jordan FM. Language trees support the express-train sequence of Austronesian
2185 expansion. *Nature*. 2000 Jun 29;405(6790):1052-5. doi: 10.1038/35016575.
- 2186 12. Green RE, Krause J, Briggs AW, Maricic T, Stenzel U, Kircher M, Patterson N, Li H, Zhai W, Fritz
2187 MH, Hansen NF, Durand EY, Malaspinas AS, Jensen JD, Marques-Bonet T, Alkan C, Prüfer K,
2188 Meyer M, Burbano HA, Good JM, Schultz R, Aximu-Petri A, Butthof A, Höber B, Höffner B,
2189 Siegemund M, Weihmann A, Nusbaum C, Lander ES, Russ C, Novod N, Affourtit J, Egholm M,
2190 Verna C, Rudan P, Brajkovic D, Kucan Ž, Gušić I, Doronichev VB, Golovanova LV, Lalueza-Fox C, de
2191 la Rasilla M, Fortea J, Rosas A, Schmitz RW, Johnson PLF, Eichler EE, Falush D, Birney E, Mullikin
2192 JC, Slatkin M, Nielsen R, Kelso J, Lachmann M, Reich D, Pääbo S. A draft sequence of the
2193 Neandertal genome. *Science*. 2010 May 7;328(5979):710-722. doi: 10.1126/science.1188021.
- 2194 13. Gronau I, Hubisz MJ, Gulko B, Danko CG, Siepel A. Bayesian inference of ancient human
2195 demography from individual genome sequences. *Nat Genet*. 2011 Sep 18;43(10):1031-4. doi:
2196 10.1038/ng.937.
- 2197 14. Gutenkunst RN, Hernandez RD, Williamson SH, Bustamante CD. Inferring the joint demographic

- 2198 history of multiple populations from multidimensional SNP frequency data. *PLoS Genet.* 2009
2199 Oct;5(10):e1000695. doi: 10.1371/journal.pgen.1000695.
- 2200 15. Haak W, Lazaridis I, Patterson N, Rohland N, Mallick S, Llamas B, Brandt G, Nordenfelt S, Harney
2201 E, Stewardson K, Fu Q, Mittnik A, Bánffy E, Economou C, Francken M, Friederich S, Pena RG,
2202 Hallgren F, Khartanovich V, Khokhlov A, Kunst M, Kuznetsov P, Meller H, Mochalov O, Moiseyev
2203 V, Nicklisch N, Pichler SL, Risch R, Rojo Guerra MA, Roth C, Szécsényi-Nagy A, Wahl J, Meyer M,
2204 Krause J, Brown D, Anthony D, Cooper A, Alt KW, Reich D. Massive migration from the steppe
2205 was a source for Indo-European languages in Europe. *Nature.* 2015 Jun 11;522(7555):207-11.
2206 doi: 10.1038/nature14317.
- 2207 16. Hajdinjak M, Mafessoni F, Skov L, Vernot B, Hübner A, Fu Q, Essel E, Nagel S, Nickel B, Richter J,
2208 Moldovan OT, Constantin S, Endarova E, Zahariev N, Spasov R, Welker F, Smith GM, Sinet-
2209 Mathiot V, Paskulin L, Fewlass H, Talamo S, Rezek Z, Sirakova S, Sirakov N, McPherron SP,
2210 Tsanova T, Hublin JJ, Peter BM, Meyer M, Skoglund P, Kelso J, Pääbo S. Initial Upper Palaeolithic
2211 humans in Europe had recent Neanderthal ancestry. *Nature.* 2021 Apr;592(7853):253-257. doi:
2212 10.1038/s41586-021-03335-3.
- 2213 17. Hammer MF, Woerner AE, Mendez FL, Watkins JC, Wall JD. Genetic evidence for archaic
2214 admixture in Africa. *Proc Natl Acad Sci U S A.* 2011 Sep 13;108(37):15123-8. doi:
2215 10.1073/pnas.1109300108.
- 2216 18. Hirschhorn JN, Daly MJ. Genome-wide association studies for common diseases and complex
2217 traits. *Nat Rev Genet.* 2005 Feb;6(2):95-108. doi: 10.1038/nrg1521.
- 2218 19. Hubisz MJ, Williams AL, Siepel A. Mapping gene flow between ancient hominins through
2219 demography-aware inference of the ancestral recombination graph. *PLoS Genet.* 2020 Aug
2220 6;16(8):e1008895. doi: 10.1371/journal.pgen.1008895.
- 2221 20. Hubisz M, Siepel A. Inference of Ancestral Recombination Graphs Using ARGweaver. *Methods*
2222 *Mol Biol.* 2020;2090:231-266. doi: 10.1007/978-1-0716-0199-0_10.
- 2223 21. Ioannidis JP (2005) Why most published research findings are false. *PLoS Med* 2: e124.
- 2224 22. Jeong C, Balanovsky O, Lukianova E, Kahbatkzy N, Flegontov P, Zaporozhchenko V, Immel A,
2225 Wang CC, Ixan O, Khussainova E, Bekmanov B, Zaibert V, Lavryashina M, Pocheshkhova E,
2226 Yusupov Y, Agdzhoyan A, Koshel S, Bukin A, Nymadawa P, Turdikulova S, Dalimova D, Churnosov
2227 M, Skhalyakho R, Daragan D, Bogunov Y, Bogunova A, Shtrunov A, Dubova N, Zhabagin M,
2228 Yepiskoposyan L, Churakov V, Pislegin N, Damba L, Saroyants L, Dibirova K, Atramentova L,
2229 Utevska O, Idrisov E, Kamenshchikova E, Evseeva I, Metspalu M, Outram AK, Robbeets M,
2230 Djansugurova L, Balanovska E, Schiffels S, Haak W, Reich D, Krause J. The genetic history of
2231 admixture across inner Eurasia. *Nat Ecol Evol.* 2019 Jun;3(6):966-976. doi: 10.1038/s41559-019-
2232 0878-2.
- 2233 23. Kamm J, Terhorst J, Durbin R, Song YS. Efficiently inferring the demographic history of many
2234 populations with allele count data. *J Am Stat Assoc.* 2020;115(531):1472-1487. doi:
2235 10.1080/01621459.2019.1635482.
- 2236 24. Ko AM, Chen CY, Fu Q, Delfin F, Li M, Chiu HL, Stoneking M, Ko YC. Early Austronesians: into and
2237 out of Taiwan. *Am J Hum Genet.* 2014 Mar 6;94(3):426-36. doi: 10.1016/j.ajhg.2014.02.003.
- 2238 25. Lachance J, Vernot B, Elbers CC, Ferwerda B, Froment A, Bodo JM, Lema G, Fu W, Nyambo TB,
2239 Rebbeck TR, Zhang K, Akey JM, Tishkoff SA. Evolutionary history and adaptation from high-
2240 coverage whole-genome sequences of diverse African hunter-gatherers. *Cell.* 2012 Aug
2241 3;150(3):457-69. doi: 10.1016/j.cell.2012.07.009.
- 2242 26. Lazaridis I, Patterson N, Mittnik A, Renaud G, Mallick S, Kirsanow K, Sudmant PH, Schraiber JG,
2243 Castellano S, Lipson M, Berger B, Economou C, Bollongino R, Fu Q, Bos KI, Nordenfelt S, Li H, de
2244 Filippo C, Prüfer K, Sawyer S, Posth C, Haak W, Hallgren F, Fornander E, Rohland N, Delsate D,
2245 Francken M, Guinet JM, Wahl J, Ayodo G, Babiker HA, Bailliet G, Balanovska E, Balanovsky O,
2246 Barrantes R, Bedoya G, Ben-Ami H, Bene J, Berrada F, Bravi CM, Brisighelli F, Busby GB, Cali F,
2247 Churnosov M, Cole DE, Corach D, Damba L, van Driem G, Dryomov S, Dugoujon JM, Fedorova SA,
2248 Gallego Romero I, Gubina M, Hammer M, Henn BM, Hervig T, Hodoglugil U, Jha AR, Karachanak-
2249 Yankova S, Khusainova R, Khusnutdinova E, Kittles R, Kivisild T, Klitz W, Kučinskas V,
2250 Kushniarevich A, Laredj L, Litvinov S, Loukidis T, Mahley RW, Melegh B, Metspalu E, Molina J,
2251 Mountain J, Näkkäläjärvi K, Nesheva D, Nyambo T, Osipova L, Parik J, Platonov F, Posukh O,

- 2252 Romano V, Rothhammer F, Rudan I, Ruizbakiev R, Sahakyan H, Sajantila A, Salas A, Starikovskaya
2253 EB, Tarekegn A, Toncheva D, Turdikulova S, Uktveryte I, Utevska O, Vasquez R, Villena M,
2254 Voevoda M, Winkler CA, Yepiskoposyan L, Zalloua P, Zemunik T, Cooper A, Capelli C, Thomas
2255 MG, Ruiz-Linares A, Tishkoff SA, Singh L, Thangaraj K, Vilems R, Comas D, Sukernik R, Metspalu
2256 M, Meyer M, Eichler EE, Burger J, Slatkin M, Pääbo S, Kelso J, Reich D, Krause J. Ancient human
2257 genomes suggest three ancestral populations for present-day Europeans. *Nature*. 2014 Sep
2258 18;513(7518):409-13. doi: 10.1038/nature13673.
- 2259 27. Lazaridis I, Nadel D, Rollefson G, Merrett DC, Rohland N, Mallick S, Fernandes D, Novak M,
2260 Gamarra B, Sirak K, Connell S, Stewardson K, Harney E, Fu Q, Gonzalez-Fortes G, Jones ER,
2261 Roodenberg SA, Lengyel G, Bocquentin F, Gasparian B, Monge JM, Gregg M, Eshed V, Mizrahi AS,
2262 Meiklejohn C, Gerritsen F, Bejenaru L, Blüher M, Campbell A, Cavalleri G, Comas D, Froguel P,
2263 Gilbert E, Kerr SM, Kovacs P, Krause J, McGettigan D, Merrigan M, Merriwether DA, O'Reilly S,
2264 Richards MB, Semino O, Shamoony-Pour M, Stefanescu G, Stumvoll M, Tönjes A, Torroni A,
2265 Wilson JF, Yengo L, Hovhannisyan NA, Patterson N, Pinhasi R, Reich D. Genomic insights into the
2266 origin of farming in the ancient Near East. *Nature*. 2016 Aug 25;536(7617):419-24. doi:
2267 10.1038/nature19310.
- 2268 28. Leppälä K, Nielsen SV, Mailund T. admixturegraph: an R package for admixture graph
2269 manipulation and fitting. *Bioinformatics*. 2017 Jun 1;33(11):1738-1740. doi:
2270 10.1093/bioinformatics/btx048.
- 2271 29. Librado P, Khan N, Fages A, Kusliy MA, Suchan T, Tonasso-Calvière L, Schiavinato S, Alioglu D,
2272 Fromentier A, Perdereau A, Aury JM, Gaunitz C, Chauvey L, Seguin-Orlando A, Der Sarkissian C,
2273 Southon J, Shapiro B, Tishkin AA, Kovalev AA, Alquraishi S, Alfarhan AH, Al-Rasheid KAS, Seregely
2274 T, Klassen L, Iversen R, Bignon-Lau O, Bodu P, Olive M, Castel JC, Boudadi-Maligne M, Alvarez N,
2275 Germonpré M, Moskal-Del Hoyo M, Wilczyński J, Pospuła S, Lasota-Kuś A, Tunia K, Nowak M,
2276 Rannamäe E, Saarma U, Boeskorov G, Lõugas L, Kyselý R, Peške L, Bălăşescu A, Dumitraşcu V,
2277 Dobrescu R, Gerber D, Kiss V, Szécsényi-Nagy A, Mende BG, Gallina Z, Somogyi K, Kulcsár G, Gál
2278 E, Bendrey R, Allentoft ME, Sirbu G, Dergachev V, Shephard H, Tomadini N, Grouard S, Kasparov
2279 A, Basilyan AE, Anisimov MA, Nikolskiy PA, Pavlova EY, Pitulko V, Brem G, Wallner B, Schwall C,
2280 Keller M, Kitagawa K, Bessudnov AN, Bessudnov A, Taylor W, Magail J, Gantulga JO,
2281 Bayarsaikhan J, Erdenebaatar D, Tabaldiev K, Mijiddorj E, Boldgiv B, Tsagaan T, Pruvost M, Olsen
2282 S, Makarewicz CA, Valenzuela Lamas S, Albizuri Canadell S, Nieto Espinet A, Iborra MP, Lira
2283 Garrido J, Rodríguez González E, Celestino S, Olària C, Arsuaga JL, Kotova N, Pryor A, Crabtree P,
2284 Zhumatayev R, Toleubaev A, Morgunova NL, Kuznetsova T, Lordkipanize D, Marzullo M, Prato O,
2285 Bagnasco Gianni G, Tecchiati U, Clavel B, Lepetz S, Davoudi H, Mashkour M, Berezina NY,
2286 Stockhammer PW, Krause J, Haak W, Morales-Muñiz A, Benecke N, Hofreiter M, Ludwig A,
2287 Graphodatsky AS, Peters J, Kiryushin KY, Iderkhangai TO, Bokovenko NA, Vasiliev SK, Seregin NN,
2288 Chugunov KV, Plasteeva NA, Baryshnikov GF, Petrova E, Sablin M, Ananyevskaya E, Logvin A,
2289 Shevnina I, Logvin V, Kalieva S, Loman V, Kukushkin I, Merz I, Merz V, Sakenov S, Varfolomeyev
2290 V, Usmanova E, Zaibert V, Arbuckle B, Belinskiy AB, Kalmykov A, Reinhold S, Hansen S, Yudin AI,
2291 Vybornov AA, Epimakhov A, Berezina NS, Roslyakova N, Kosintsev PA, Kuznetsov PF, Anthony D,
2292 Kroonen GJ, Kristiansen K, Wincker P, Outram A, Orlando L. The origins and spread of domestic
2293 horses from the Western Eurasian steppes. *Nature*. 2021 Oct;598(7882):634-640. doi:
2294 10.1038/s41586-021-04018-9.
- 2295 30. Lipson M, Loh PR, Levin A, Reich D, Patterson N, Berger B. Efficient moment-based inference of
2296 admixture parameters and sources of gene flow. *Mol Biol Evol*. 2013 Aug;30(8):1788-802. doi:
2297 10.1093/molbev/mst099.
- 2298 31. Lipson M, Szécsényi-Nagy A, Mallick S, Pósa A, Stégmár B, Keerl V, Rohland N, Stewardson K,
2299 Ferry M, Michel M, Oppenheimer J, Broomandkoshbacht N, Harney E, Nordenfelt S, Llamas B,
2300 Gusztáv Mende B, Köhler K, Oross K, Bondár M, Marton T, Osztás A, Jakucs J, Paluch T, Horváth
2301 F, Csengeri P, Koós J, Sebők K, Anders A, Raczky P, Regenye J, Barna JP, Fábíán S, Serlegi G, Toldi
2302 Z, Gyöngyvér Nagy E, Dani J, Molnár E, Pálfi G, Márk L, Melegh B, Bánfai Z, Domboróczi L,
2303 Fernández-Eraso J, Antonio Mujika-Alustiza J, Alonso Fernández C, Jiménez Echevarría J,
2304 Bollongino R, Orschiedt J, Schierhold K, Meller H, Cooper A, Burger J, Bánffy E, Alt KW, Lalueza-
2305 Fox C, Haak W, Reich D. Parallel palaeogenomic transects reveal complex genetic history of early

- 2306 European farmers. *Nature*. 2017 Nov 16;551(7680):368-372. doi: 10.1038/nature24476.
- 2307 32. Lipson M, Ribot I, Mallick S, Rohland N, Olalde I, Adamski N, Broomandkhoshbacht N, Lawson
2308 AM, López S, Oppenheimer J, Stewardson K, Asombang RN, Bocherens H, Bradman N, Culleton
2309 BJ, Cornelissen E, Crevecoeur I, de Maret P, Fomine FLM, Lavachery P, Mindzie CM, Orban R,
2310 Sawchuk E, Semal P, Thomas MG, Van Neer W, Veeramah KR, Kennett DJ, Patterson N,
2311 Hellenthal G, Lalueza-Fox C, MacEachern S, Prendergast ME, Reich D. Ancient West African
2312 foragers in the context of African population history. *Nature*. 2020 Jan;577(7792):665-670. doi:
2313 10.1038/s41586-020-1929-1.
- 2314 33. Lipson M. Applying f4-statistics and admixture graphs: Theory and examples. *Mol Ecol Resour*.
2315 2020 Nov;20(6):1658-1667. doi: 10.1111/1755-0998.13230.
- 2316 34. Lipson M, Sawchuk EA, Thompson JC, Oppenheimer J, Tryon CA, Ranhorn KL, de Luna KM, Sirak
2317 KA, Olalde I, Ambrose SH, Arthur JW, Arthur KJW, Ayodo G, Bertacchi A, Cerezo-Román JI,
2318 Culleton BJ, Curtis MC, Davis J, Gidna AO, Hanson A, Kaliba P, Katongo M, Kwekason A, Laird MF,
2319 Lewis J, Mabulla AZP, Mapemba F, Morris A, Mudenda G, Mwafulirwa R, Mwangomba D,
2320 Ndiema E, Ogola C, Schilt F, Willoughby PR, Wright DK, Zipkin A, Pinhasi R, Kennett DJ, Manthi
2321 FK, Rohland N, Patterson N, Reich D, Prendergast ME. Ancient DNA and deep population
2322 structure in sub-Saharan African foragers. *Nature*. 2022 Mar;603(7900):290-296. doi:
2323 10.1038/s41586-022-04430-9.
- 2324 35. Lu D, Lou H, Yuan K, Wang X, Wang Y, Zhang C, Lu Y, Yang X, Deng L, Zhou Y, Feng Q, Hu Y, Ding
2325 Q, Yang Y, Li S, Jin L, Guan Y, Su B, Kang L, Xu S. Ancestral Origins and Genetic History of Tibetan
2326 Highlanders. *Am J Hum Genet*. 2016 Sep 1;99(3):580-594. doi: 10.1016/j.ajhg.2016.07.002.
- 2327 36. Mallick S, Li H, Lipson M, Mathieson I, Gymrek M, Racimo F, Zhao M, Chennagiri N, Nordenfelt S,
2328 Tandon A, Skoglund P, Lazaridis I, Sankararaman S, Fu Q, Rohland N, Renaud G, Erlich Y, Willems
2329 T, Gallo C, Spence JP, Song YS, Poletti G, Balloux F, van Driem G, de Knijff P, Romero IG, Jha AR,
2330 Behar DM, Bravi CM, Capelli C, Hervig T, Moreno-Estrada A, Posukh OL, Balanovska E,
2331 Balanovsky O, Karachanak-Yankova S, Sahakyan H, Toncheva D, Yepiskoposyan L, Tyler-Smith C,
2332 Xue Y, Abdullah MS, Ruiz-Linares A, Beall CM, Di Rienzo A, Jeong C, Starikovskaya EB, Metspalu E,
2333 Parik J, Vilems R, Henn BM, Hodoglugil U, Mahley R, Sajantila A, Stamatoyannopoulos G, Wee
2334 JT, Khusai nova R, Khusnutdinova E, Litvinov S, Ayodo G, Comas D, Hammer MF, Kivisild T, Klitz
2335 W, Winkler CA, Labuda D, Bamshad M, Jorde LB, Tishkoff SA, Watkins WS, Metspalu M, Dryomov
2336 S, Sukernik R, Singh L, Thangaraj K, Pääbo S, Kelso J, Patterson N, Reich D. The Simons Genome
2337 Diversity Project: 300 genomes from 142 diverse populations. *Nature*. 2016 Oct
2338 13;538(7624):201-206. Doi: 10.1038/nature18964.
- 2339 37. McColl H, Racimo F, Vinner L, Demeter F, Gakuhari T, Moreno-Mayar JV, van Driem G, Gram
2340 Wilken U, Seguin-Orlando A, de la Fuente Castro C, Wasef S, Shoocongdej R, Souksavatdy V,
2341 Sayavongkhamdy T, Saidin MM, Allentoft ME, Sato T, Malaspinas AS, Aghakhanian FA,
2342 Korneliussen T, Prohaska A, Margaryan A, de Barros Damgaard P, Kaewsutthi S, Lertrit P, Nguyen
2343 TMH, Hung HC, Minh Tran T, Nghia Truong H, Nguyen GH, Shahidan S, Wiradnyana K, Matsumae
2344 H, Shigehara N, Yoneda M, Ishida H, Masuyama T, Yamada Y, Tajima A, Shibata H, Toyoda A,
2345 Hanihara T, Nakagome S, Deviese T, Bacon AM, Durringer P, Ponche JL, Shackelford L, Patole-
2346 Edoumba E, Nguyen AT, Bellina-Pryce B, Galipaud JC, Kinaston R, Buckley H, Pottier C,
2347 Rasmussen S, Higham T, Foley RA, Lahr MM, Orlando L, Sikora M, Phipps ME, Oota H, Higham C,
2348 Lambert DM, Willerslev E. The prehistoric peopling of Southeast Asia. *Science*. 2018 Jul
2349 6;361(6397):88-92. Doi: 10.1126/science.aat3628.
- 2350 38. McVean G. A genealogical interpretation of principal components analysis. *PLoS Genet*. 2009
2351 Oct;5(10):e1000686. Doi: 10.1371/journal.pgen.1000686. Epub 2009 Oct 16.
- 2352 39. Molloy EK, Durvasula A, Sankararaman S. Advancing admixture graph estimation via maximum
2353 likelihood network orientation. *Bioinformatics*. 2021 Jul 12;37(Suppl_1):i142-i150. Doi:
2354 10.1093/bioinformatics/btab267.
- 2355 40. Moreno-Mayar JV, Vinner L, de Barros Damgaard P, de la Fuente C, Chan J, Spence JP, Allentoft
2356 ME, Vimala T, Racimo F, Pinotti T, Rasmussen S, Margaryan A, Iraeta Orbegozo M,
2357 Mylopotamitaki D, Wooller M, Bataille C, Becerra-Valdivia L, Chivall D, Comeskey D, Deviese T,
2358 Grayson DK, George L, Harry H, Alexandersen V, Primeau C, Erlandson J, Rodrigues-Carvalho C,
2359 Reis S, Bastos MQR, Cybulski J, Vullo C, Morello F, Vilar M, Wells S, Gregersen K, Hansen KL,

- 2360 Lynnerup N, Mirazón Lahr M, Kjær K, Strauss A, Alfonso-Durruty M, Salas A, Schroeder H,
2361 Higham T, Malhi RS, Rasic JT, Souza L, Santos FR, Malaspinas AS, Sikora M, Nielsen R, Song YS,
2362 Meltzer DJ, Willerslev E. Early human dispersals within the Americas. *Science*. 2018 Dec
2363 7;362(6419):eaav2621. Doi: 10.1126/science.aav2621.
- 2364 41. Narasimhan VM, Patterson N, Moorjani P, Rohland N, Bernardos R, Mallick S, Lazaridis I,
2365 Nakatsuka N, Olalde I, Lipson M, Kim AM, Olivieri LM, Coppa A, Vidale M, Mallory J, Moiseyev V,
2366 Kitov E, Monge J, Adamski N, Alex N, Broomandkhoshbacht N, Candilio F, Callan K, Cheronet O,
2367 Culleton BJ, Ferry M, Fernandes D, Freilich S, Gamarra B, Gaudio D, Hajdinjak M, Harney É,
2368 Harper TK, Keating D, Lawson AM, Mah M, Mandl K, Michel M, Novak M, Oppenheimer J, Rai N,
2369 Sirak K, Slon V, Stewardson K, Zalzal F, Zhang Z, Akhatov G, Bagashev AN, Bagnera A, Baitanayev
2370 B, Bendezu-Sarmiento J, Bissembaev AA, Bonora GL, Charynov TT, Chikisheva T, Dashkovskiy
2371 PK, Derevianko A, Dobeš M, Douka K, Dubova N, Duisengali MN, Enshin D, Epimakhov A, Fribus
2372 AV, Fuller D, Goryachev A, Gromov A, Grushin SP, Hanks B, Judd M, Kazizov E, Khokhlov A, Krygin
2373 AP, Kupriyanova E, Kuznetsov P, Luiselli D, Maksudov F, Mamedov AM, Mamirov TB, Meiklejohn
2374 C, Merrett DC, Micheli R, Mochalov O, Mustafokulov S, Nayak A, Pettener D, Potts R, Razhev D,
2375 Rykun M, Sarno S, Savenkova TM, Sikhymbaeva K, Slepchenko SM, Soltobaev OA, Stepanova N,
2376 Svyatko S, Tabaldiev K, Teschler-Nicola M, Tishkin AA, Tkachev VV, Vasilyev S, Velemínský P,
2377 Voyakin D, Yermolayeva A, Zahir M, Zubkov VS, Zubova A, Shinde VS, Lalueza-Fox C, Meyer M,
2378 Anthony D, Boivin N, Thangaraj K, Kennett DJ, Frachetti M, Pinhasi R, Reich D. The formation of
2379 human populations in South and Central Asia. *Science*. 2019 Sep 6;365(6457):eaat7487. Doi:
2380 10.1126/science.aat7487.
- 2381 42. Patterson N, Moorjani P, Luo Y, Mallick S, Rohland N, Zhan Y, Genschoreck T, Webster T, Reich D.
2382 Ancient admixture in human history. *Genetics*. 2012 Nov;192(3):1065-93. Doi:
2383 10.1534/genetics.112.145037.
- 2384 43. Pickrell JK, Pritchard JK. Inference of population splits and mixtures from genome-wide allele
2385 frequency data. *PLoS Genet*. 2012;8(11):e1002967. Doi: 10.1371/journal.pgen.1002967.
- 2386 44. Posth C, Nakatsuka N, Lazaridis I, Skoglund P, Mallick S, Lamnidis TC, Rohland N, Nägele K,
2387 Adamski N, Bertolini E, Broomandkhoshbacht N, Cooper A, Culleton BJ, Ferraz T, Ferry M,
2388 Furtwängler A, Haak W, Harkins K, Harper TK, Hünemeier T, Lawson AM, Llamas B, Michel M,
2389 Nelson E, Oppenheimer J, Patterson N, Schiffels S, Sedig J, Stewardson K, Talamo S, Wang CC,
2390 Hublin JJ, Hubbe M, Harvati K, Nuevo Delaunay A, Beier J, Francken M, Kaulicke P, Reyes-
2391 Centeno H, Rademaker K, Trask WR, Robinson M, Gutierrez SM, Prüfer KM, Salazar-García DC,
2392 Chim EN, Müller Plumm Gomes L, Alves ML, Liryo A, Ingles M, Oliveira RE, Bernardo DV, Barioni
2393 A, Wesolowski V, Scheifler NA, Rivera MA, Plens CR, Messineo PG, Figuti L, Corach D, Scabuzzo C,
2394 Eggers S, DeBlasis P, Reindel M, Méndez C, Politis G, Tomasto-Cagigao E, Kennett DJ, Strauss A,
2395 Fehren-Schmitz L, Krause J, Reich D. Reconstructing the Deep Population History of Central and
2396 South America. *Cell*. 2018 Nov 15;175(5):1185-1197.e22. doi: 10.1016/j.cell.2018.10.027.
- 2397 45. Prüfer K, Racimo F, Patterson N, Jay F, Sankararaman S, Sawyer S, Heinze A, Renaud G, Sudmant
2398 PH, de Filippo C, Li H, Mallick S, Dannemann M, Fu Q, Kircher M, Kuhlwilm M, Lachmann M,
2399 Meyer M, Ongyerth M, Siebauer M, Theunert C, Tandon A, Moorjani P, Pickrell J, Mullikin JC,
2400 Vohr SH, Green RE, Hellmann I, Johnson PL, Blanche H, Cann H, Kitzman JO, Shendure J, Eichler
2401 EE, Lein ES, Bakken TE, Golovanova LV, Doronichev VB, Shunkov MV, Derevianko AP, Viola B,
2402 Slatkin M, Reich D, Kelso J, Pääbo S. The complete genome sequence of a Neanderthal from the
2403 Altai Mountains. *Nature*. 2014 Jan 2;505(7481):43-9. Doi: 10.1038/nature12886.
- 2404 46. Raghavan M, Skoglund P, Graf KE, Metspalu M, Albrechtsen A, Moltke I, Rasmussen S, Stafford
2405 TW Jr, Orlando L, Metspalu E, Karmin M, Tambets K, Rootsi S, Mägi R, Campos PF, Balanovska E,
2406 Balanovsky O, Khusnutdinova E, Litvinov S, Osipova LP, Fedorova SA, Voevoda MI, DeGiorgio M,
2407 Sicheritz-Ponten T, Brunak S, Demeshchenko S, Kivisild T, Vilems R, Nielsen R, Jakobsson M,
2408 Willerslev E. Upper Palaeolithic Siberian genome reveals dual ancestry of Native Americans.
2409 *Nature*. 2014 Jan 2;505(7481):87-91. Doi: 10.1038/nature12736.
- 2410 47. Raghavan M, Steinrücken M, Harris K, Schiffels S, Rasmussen S, DeGiorgio M, Albrechtsen A,
2411 Valdiosera C, Ávila-Arcos MC, Malaspinas AS, Eriksson A, Moltke I, Metspalu M, Homburger JR,
2412 Wall J, Cornejo OE, Moreno-Mayar JV, Korneliusson TS, Pierre T, Rasmussen M, Campos PF, de
2413 Barros Damgaard P, Allentoft ME, Lindo J, Metspalu E, Rodríguez-Varela R, Mansilla J,

- 2414 Henrickson C, Seguin-Orlando A, Malmström H, Stafford T Jr, Shringarpure SS, Moreno-Estrada
2415 A, Karmin M, Tambets K, Bergström A, Xue Y, Warmuth V, Friend AD, Singarayer J, Valdes P,
2416 Balloux F, LeBoreiro I, Vera JL, Rangel-Villalobos H, Pettener D, Luiselli D, Davis LG, Heyer E,
2417 Zollikofer CPE, Ponce de León MS, Smith CI, Grimes V, Pike KA, Deal M, Fuller BT, Arriaza B,
2418 Standen V, Luz MF, Ricaut F, Guidon N, Osipova L, Voevoda MI, Posukh OL, Balanovsky O,
2419 Lavryashina M, Bogunov Y, Khusnutdinova E, Gubina M, Balanovska E, Fedorova S, Litvinov S,
2420 Malyarchuk B, Derenko M, Mosher MJ, Archer D, Cybulski J, Petzelt B, Mitchell J, Worl R,
2421 Norman PJ, Parham P, Kemp BM, Kivisild T, Tyler-Smith C, Sandhu MS, Crawford M, Villems R,
2422 Smith DG, Waters MR, Goebel T, Johnson JR, Malhi RS, Jakobsson M, Meltzer DJ, Manica A,
2423 Durbin R, Bustamante CD, Song YS, Nielsen R, Willerslev E. Genomic evidence for the Pleistocene
2424 and recent population history of Native Americans. *Science*. 2015 Aug 21;349(6250):aab3884.
2425 Doi: 10.1126/science.aab3884.
- 2426 48. Reich D, Thangaraj K, Patterson N, Price AL, Singh L. Reconstructing Indian population history.
2427 *Nature*. 2009 Sep 24;461(7263):489-94. Doi: 10.1038/nature08365.
- 2428 49. Reich D, Patterson N, Kircher M, Delfin F, Nandineni MR, Pugach I, Ko AM, Ko YC, Jinam TA,
2429 Phipps ME, Saitou N, Wollstein A, Kayser M, Pääbo S, Stoneking M. Denisova admixture and the
2430 first modern human dispersals into Southeast Asia and Oceania. *Am J Hum Genet*. 2011 Oct
2431 7;89(4):516-28. Doi: 10.1016/j.ajhg.2011.09.005.
- 2432 50. Reich D, Patterson N, Campbell D, Tandon A, Mazieres S, Ray N, Parra MV, Rojas W, Duque C,
2433 Mesa N, García LF, Triana O, Blair S, Maestre A, Dib JC, Bravi CM, Bailliet G, Corach D, Hünemeier
2434 T, Bortolini MC, Salzano FM, Petzl-Erler ML, Acuña-Alonzo V, Aguilar-Salinas C, Canizales-
2435 Quinteros S, Tusié-Luna T, Riba L, Rodríguez-Cruz M, Lopez-Alarcón M, Coral-Vazquez R, Canto-
2436 Cetina T, Silva-Zolezzi I, Fernandez-Lopez JC, Contreras AV, Jimenez-Sanchez G, Gómez-Vázquez
2437 MJ, Molina J, Carracedo A, Salas A, Gallo C, Poletti G, Witonsky DB, Alkorta-Aranburu G, Sukernik
2438 RI, Osipova L, Fedorova SA, Vasquez R, Villena M, Moreau C, Barrantes R, Pauls D, Excoffier L,
2439 Bedoya G, Rothhammer F, Dugoujon JM, Larrouy G, Klitz W, Labuda D, Kidd J, Kidd K, Di Rienzo
2440 A, Freimer NB, Price AL, Ruiz-Linares A. Reconstructing Native American population history.
2441 *Nature*. 2012 Aug 16;488(7411):370-4. Doi: 10.1038/nature11258.
- 2442 51. Rogers AR. Legofit: estimating population history from genetic data. *BMC Bioinformatics*. 2019
2443 Oct 28;20(1):526. Doi: 10.1186/s12859-019-3154-1.
- 2444 52. Schiffels S, Durbin R. Inferring human population size and separation history from multiple
2445 genome sequences. *Nat Genet*. 2014 Aug;46(8):919-25. Doi: 10.1038/ng.3015.
- 2446 53. Schiffels S, Haak W, Paajanen P, Llamas B, Popescu E, Loe L, Clarke R, Lyons A, Mortimer R, Sayer
2447 D, Tyler-Smith C, Cooper A, Durbin R. Iron Age and Anglo-Saxon genomes from East England
2448 reveal British migration history. *Nat Commun*. 2016 Jan 19;7:10408. Doi:
2449 10.1038/ncomms10408.
- 2450 54. Seguin-Orlando A, Korneliussen TS, Sikora M, Malaspina AS, Manica A, Moltke I, Albrechtsen A,
2451 Ko A, Margaryan A, Moiseyev V, Goebel T, Westaway M, Lambert D, Khartanovich V, Wall JD,
2452 Nigst PR, Foley RA, Lahr MM, Nielsen R, Orlando L, Willerslev E. Paleogenomics. Genomic
2453 structure in Europeans dating back at least 36,200 years. *Science*. 2014 Nov 28;346(6213):1113-
2454 8. Doi: 10.1126/science.aaa0114.
- 2455 55. Shinde V, Narasimhan VM, Rohland N, Mallick S, Mah M, Lipson M, Nakatsuka N, Adamski N,
2456 Broomandkoshbacht N, Ferry M, Lawson AM, Michel M, Oppenheimer J, Stewardson K, Jadhav
2457 N, Kim YJ, Chatterjee M, Munshi A, Panyam A, Waghmare P, Yadav Y, Patel H, Kaushik A,
2458 Thangaraj K, Meyer M, Patterson N, Rai N, Reich D. An Ancient Harappan Genome Lacks
2459 Ancestry from Steppe Pastoralists or Iranian Farmers. *Cell*. 2019 Oct 17;179(3):729-735.e10. doi:
2460 10.1016/j.cell.2019.08.048.
- 2461 56. Sikora M, Pitulko VV, Sousa VC, Allentoft ME, Vinner L, Rasmussen S, Margaryan A, de Barros
2462 Damgaard P, de la Fuente C, Renaud G, Yang MA, Fu Q, Dupanloup I, Giampoudakis K, Nogués-
2463 Bravo D, Rahbek C, Kroonen G, Peyrot M, McColl H, Vasilyev SV, Veselovskaya E, Gerasimova M,
2464 Pavlova EY, Chasnyk VG, Nikolskiy PA, Gromov AV, Khartanovich VI, Moiseyev V, Grebenyuk PS,
2465 Fedorchenko AY, Lebedintsev AI, Slobodin SB, Malyarchuk BA, Martiniano R, Meldgaard M,
2466 Arppe L, Palo JU, Sundell T, Mannermaa K, Putkonen M, Alexandersen V, Primeau C,
2467 Baimukhanov N, Malhi RS, Sjögren KG, Kristiansen K, Wessman A, Sajantila A, Lahr MM, Durbin

- 2468 R, Nielsen R, Meltzer DJ, Excoffier L, Willerslev E. The population history of northeastern Siberia
2469 since the Pleistocene. *Nature*. 2019 Jun;570(7760):182-188. Doi: 10.1038/s41586-019-1279-z.
- 2470 57. Skoglund P, Posth C, Sirak K, Spriggs M, Valentin F, Bedford S, Clark GR, Reepmeyer C, Petchey F,
2471 Fernandes D, Fu Q, Harney E, Lipson M, Mallick S, Novak M, Rohland N, Stewardson K, Abdullah
2472 S, Cox MP, Friedlaender FR, Friedlaender JS, Kivisild T, Koki G, Kusuma P, Merriwether DA, Ricaut
2473 FX, Wee JT, Patterson N, Krause J, Pinhasi R, Reich D. Genomic insights into the peopling of the
2474 Southwest Pacific. *Nature*. 2016 Oct 27;538(7626):510-513. Doi: 10.1038/nature19844.
- 2475 58. Speidel L, Forest M, Shi S, Myers SR. A method for genome-wide genealogy estimation for
2476 thousands of samples. *Nat Genet*. 2019 Sep;51(9):1321-1329. Doi: 10.1038/s41588-019-0484-x.
- 2477 59. Tambets K, Yunusbayev B, Hudjashov G, Ilumäe AM, Rootsi S, Honkola T, Vesakoski O, Atkinson
2478 Q, Skoglund P, Kushniarevich A, Litvinov S, Reidla M, Metspalu E, Saag L, Rantanen T, Karmin M,
2479 Parik J, Zhadanov SI, Gubina M, Damba LD, Bermisheva M, Reisberg T, Dibirova K, Evseeva I,
2480 Nelis M, Klovins J, Metspalu A, Esko T, Balanovsky O, Balanovska E, Khusnutdinova EK, Osipova
2481 LP, Voevoda M, VILLEMS R, Kivisild T, Metspalu M. Genes reveal traces of common recent
2482 demographic history for most of the Uralic-speaking populations. *Genome Biol*. 2018 Sep
2483 21;19(1):139. Doi: 10.1186/s13059-018-1522-1.
- 2484 60. Terhorst J, Kamm JA, Song YS. Robust and scalable inference of population history from
2485 hundreds of unphased whole genomes. *Nat Genet*. 2017 Feb;49(2):303-309. Doi:
2486 10.1038/ng.3748.
- 2487 61. Vallini L, Marciani G, Aneli S, Bortolini E, Benazzi S, Pievani T, Pagani L. Genetics and material
2488 culture support repeated expansions into Paleolithic Eurasia from a population hub out of Africa.
2489 *Genome Biol Evol*. 2022 Apr 10;14(4):evac045. Doi: 10.1093/gbe/evac045.
- 2490 62. van de Loosdrecht M, Bouzouggar A, Humphrey L, Posth C, Barton N, Aximu-Petri A, Nickel B,
2491 Nagel S, Talbi EH, El Hajraoui MA, Amzazi S, Hublin JJ, Pääbo S, Schiffels S, Meyer M, Haak W,
2492 Jeong C, Krause J. Pleistocene North African genomes link Near Eastern and sub-Saharan African
2493 human populations. *Science*. 2018 May 4;360(6388):548-552. Doi: 10.1126/science.aar8380.
- 2494 63. Wang CC, Reinhold S, Kalmykov A, Wissgott A, Brandt G, Jeong C, Cheronet O, Ferry M, Harney E,
2495 Keating D, Mallick S, Rohland N, Stewardson K, Kantorovich AR, Maslov VE, Petrenko VG, Erlikh
2496 VR, Atabiev BC, Magomedov RG, Kohl PL, Alt KW, Pichler SL, Gerling C, Meller H, Vardanyan B,
2497 Yeganyan L, Rezepkin AD, Mariaschk D, Berezina N, Gresky J, Fuchs K, Knipper C, Schiffels S,
2498 Balanovska E, Balanovsky O, Mathieson I, Higham T, Berezin YB, Buzhilova A, Trifonov V, Pinhasi
2499 R, Belinskij AB, Reich D, Hansen S, Krause J, Haak W. Ancient human genome-wide data from a
2500 3000-year interval in the Caucasus corresponds with eco-geographic regions. *Nat Commun*. 2019
2501 Feb 4;10(1):590. Doi: 10.1038/s41467-018-08220-8.
- 2502 64. Wang CC, Yeh HY, Popov AN, Zhang HQ, Matsumura H, Sirak K, Cheronet O, Kovalev A, Rohland
2503 N, Kim AM, Mallick S, Bernardos R, Tumen D, Zhao J, Liu YC, Liu JY, Mah M, Wang K, Zhang Z,
2504 Adamski N, Broomandkhoshbacht N, Callan K, Candilio F, Carlson KSD, Culleton BJ, Eccles L,
2505 Freilich S, Keating D, Lawson AM, Mandl K, Michel M, Oppenheimer J, Özdoğan KT, Stewardson
2506 K, Wen S, Yan S, Zalzal F, Chuang R, Huang CJ, Looch H, Shiung CC, Nikitin YG, Tabarev AV,
2507 Tishkin AA, Lin S, Sun ZY, Wu XM, Yang TL, Hu X, Chen L, Du H, Bayarsaikhan J, Mijiddorj E,
2508 Erdenebaatar D, Iderkhangai TO, Myagmar E, Kanzawa-Kiriyama H, Nishino M, Shinoda KI,
2509 Shubina OA, Guo J, Cai W, Deng Q, Kang L, Li D, Li D, Lin R, Nini, Shrestha R, Wang LX, Wei L, Xie
2510 G, Yao H, Zhang M, He G, Yang X, Hu R, Robbeets M, Schiffels S, Kennett DJ, Jin L, Li H, Krause J,
2511 Pinhasi R, Reich D. Genomic insights into the formation of human populations in East Asia.
2512 *Nature*. 2021 Mar;591(7850):413-419. Doi: 10.1038/s41586-021-03336-2.
- 2513 65. Yan J, Patterson N, Narasimhan VM. miqoGraph: fitting admixture graphs using mixed-integer
2514 quadratic optimization. *Bioinformatics*. 2021 Aug 25;37(16):2488-2490. Doi:
2515 10.1093/bioinformatics/btaa988.
- 2516 66. Yang MA, Gao X, Theunert C, Tong H, Aximu-Petri A, Nickel B, Slatkin M, Meyer M, Pääbo S,
2517 Kelso J, Fu Q. 40,000-Year-Old Individual from Asia Provides Insight into Early Population
2518 Structure in Eurasia. *Curr Biol*. 2017 Oct 23;27(20):3202-3208.e9. doi:
2519 10.1016/j.cub.2017.09.030. Epub 2017 Oct 12.
- 2520 67. Yang MA, Fan X, Sun B, Chen C, Lang J, Ko YC, Tsang CH, Chiu H, Wang T, Bao Q, Wu X, Hajdinjak
2521 M, Ko AM, Ding M, Cao P, Yang R, Liu F, Nickel B, Dai Q, Feng X, Zhang L, Sun C, Ning C, Zeng W,

- 2522 Zhao Y, Zhang M, Gao X, Cui Y, Reich D, Stoneking M, Fu Q. Ancient DNA indicates human
2523 population shifts and admixture in northern and southern China. *Science*. 2020 Jul
2524 17;369(6501):282-288. Doi: 10.1126/science.aba0909.
2525 68. Zhang M, Yan S, Pan W, Jin L. Phylogenetic evidence for Sino-Tibetan origin in northern China in
2526 the Late Neolithic. *Nature*. 2019 May;569(7754):112-115. doi: 10.1038/s41586-019-1153-z.

2527
2528
Supplementary Tables

Publication	Flag in the original publication	Groups (populations)					Admix. events	Padj. model I; log likelihood (LL)	Padj. model I; mean of bootstrap dist.	SNPs used	Settings for calculating f2-statistics	Topology search constraints and population modifications	Iterations														p-value best alternative vs. publ.	Used in Table 1																				
		Single on pseudo-haploids	No. of negative f2s as a fraction of total	Pub I	Pub II	Pub III							max miss=0, b_lg_size=4000000, adjust_pseudo_haploid=T	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	minac2=2	max miss=0, b_lg_size=4000000, adjust_pseudo_haploid=T	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	minac2=2	max miss=0, b_lg_size=4000000, adjust_pseudo_haploid=T	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	minac2=2	max miss=0, b_lg_size=4000000, adjust_pseudo_haploid=T	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	minac2=2	max miss=0, b_lg_size=4000000, adjust_pseudo_haploid=T	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T			minac2=2	max miss=0, b_lg_size=4000000, adjust_pseudo_haploid=T	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	minac2=2	max miss=0, b_lg_size=4000000, adjust_pseudo_haploid=T	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	minac2=2	max miss=0, b_lg_size=4000000, adjust_pseudo_haploid=T	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	minac2=2										
Bergström et al. 2020 1e		6	4	0	2	4.7	4.9			312,282	max miss=0, b_lg_size=4000000, adjust_pseudo_haploid=T	Andon Fox OG	100	38	14	0	0	3	11	0	0	0	21.4	78.6	0.060*	10,000	337	221	1	5	37	178	0.5	0.0	16.7	80.5	0.332	yes										
		7	5	0	3	8.4	11.8	2.1	100	32			14	0	0	0	14	0	0	0	0	0	100.0	0.032*	yes																							
		6	1	6	2	3.0	4.2	624,583	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	1,000			0	306	3	37	247	19	1.0	12.1	80.7	6.2	0.032	yes																								
Lazaridis et al. 2014 3		9	3	6	3	22.1	40.9	2.8	19,017	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	Mota OG [4]	4,000	0	398	0	89	266	43	0	0	22.4	66.8	10.8	0.136	1,000	32	14	0	0	0	14	0	0	0	0	0	0	100.0	0.032*	yes								
		7	1	6	4	8.7	11.1	2.2	4,000			0	216	4	61	79	72	1.9	28.2	36.6	33.3	0.072																										
		8	2	2	3	29	40.8	3.2	470,389			max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	4,000	13	143	0	4	5	134	0	2.8	3.5	93.7	0.088*		yes																						
Shinde et al. 2019 3		8	1	10	4	N/A	N/A	N/A	249,009	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T, minac2=2	Mota OG; Indus_Periphery group expanded to 4 Ind., relatives and contaminated Ind. removed	4,000	0	398	0	89	266	43	0	0	22.4	66.8	10.8	0.136	1,000	32	14	0	0	0	14	0	0	0	0	0	0	100.0	0.032*	yes								
		8	1	10	4	N/A	N/A	N/A	4,000			0	216	4	61	79	72	1.9	28.2	36.6	33.3	0.072																										
		3b	3	907.3	942	28.3	max miss=0, b_lg_size=4000000, adjust_pseudo_haploid=T	Donkey OG; group composition was altered: for all populations individuals from only one archaeological site and temporal horizon were included	1,000			25	335											0.002		1,000	10	531												0.016								
		Ext sd	4	363.3	386.2	15.7			1,000			0	747														0.002																					
		Ext se	5	267.7	286.5	10.4			1,000			0	894														0.002																					
		N/A	10	3	3	6			N/A			N/A	N/A	6,343,116	1,000	N/A	986	N/A	N/A	N/A	N/A	N/A	N/A	N/A			N/A	N/A	N/A																			
		3b	3	679.9	711.1	23.9			max miss=0, b_lg_size=4000000, adjust_pseudo_haploid=T, minac2=2			Denison OG [4]	1,000	1	324	22	51	78	173	6.8	15.7	24.1	53.4	0.040			1,000	30	535	0	0	24	511	0	0	0	4.5	95.5	0.788*	yes								
		Ext sd	4	202.8	224.4	14.1							1,000	3	784	0	2	223	559	0	0	3	28.4	71.3				0.720	yes																			
		Ext se	5	121	139.6	6.9							1,000	954																0.002																		
		N/A	10	3	3	6							N/A	N/A	N/A	1,767,419	1,000	N/A	1,000	N/A	N/A	N/A	N/A	N/A				N/A	N/A	N/A	N/A	N/A																
		3b	3	1133	1163	33.8							max miss=0, b_lg_size=4000000, adjust_pseudo_haploid=T	Donkey OG	1,000	0	363															0.002	1,000	3	567													0.002
		Ext sd	4	477.1	500.5	16.3									1,000	0	729															0.002																
Ext se	5	380.5	395	13.2	1,000	906																			0.002																							
N/A	10	1	1	6	N/A	N/A				N/A	7,403,037				1,000	N/A	981	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A			N/A	N/A																			
3b	3	1133	1163	33.8	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	Denison OG [4]	4,000	0		3,996								15.9	80.8	3.2	0.0	0.008	1,000	12	6	0		8	35.8	65.1	2.8	283,890		2,000	0	1,999									0.002			
Ext sd	4	477.1	500.5	16.3			2,000	0		1,995	a fraction of models tested				26.0	63.5	6.6	3.9	0.002																													
Ext se	5	380.5	395	13.2			2,000	0		1,988	a fraction of models tested				56.8	41.4	1.8	0.0	0.002																													
N/A	11	5	0	8			71.6	112.4		4.8	263,698				2,000	0	1,988																															
S3.24	7	1	0	4			5.4	8.7	1.9	801,362	max miss=0, b_lg_size=4000000, adjust_pseudo_haploid=T	2,000			0	778	2	199	489	89	0.3	25.6		62.9	11.4	0.260	1,000	10	1	15	8	33.6		51.2	2.7	838,910	10,000	0	9,927	a fraction of models tested	6.8	83.5	9.1	0.6	0.032	yes		
S3.25	10	1	15	8			33.6	51.2	2.7	838,910	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T, minac2=2	2,000			0	2,000										0.064																						
Ext 4	12	2	22	12			29.5	61.6	2.2	363,131	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T, minac2=2	2,000			0	2,000										0.002																						
S13-9 [2]	9	2	0	4			16	29.6	2.5	543,124	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	2,000			0	906					0.0	28.5		68.5	3.0	0.100																						
S13-10a [3]	10	3	0	5			23.3	41.8	2.5	544,068	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	2,000	0	2,000					0.0	11.9	77.1	10.4		0.176	yes																							
Wang et al. 2021		9	2	0			4	16	29.6	2.5	543,124	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	2,000	0	906					0.0	38.8	60.8		0.4	0.524	1,000		10	3	0	5	23.3	41.8	2.5	544,068	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	2,000	0	1,524					0.0	44.0	56.0	0.0	0.288
		Ext 6	12	3			3	8	62.1	99.2	3.1	496,233	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	1,900	0	1,895												0.004																				
									66.4	106	3.8	203,753	same as above + minac2=2	2,000	0	1,993													0.004																			
Sikora et al. 2019		10	65.7	116.4	3.2	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	Altai, Dinka unadmixed, Denisovan max. 1 admix. event, Chimp OG	1,000	0	1,000					33.2	13.6	50.0	3.2	0.002	1,000	13	4	0	10	85.3	147.8		3.1	613,509	1,000	0	996					28.5	68.7	2.8	0.0	0.002							
		6	76.5	132	3.8			1,000	0	894							0.3	17.1	34.6		48.0	0.016	yes																									
		10	85.3	147.8	3.1			1,000	0	966							0.1	37.0	52.3		10.6	0.024																										
		3f [left]	13	4	0			10	65.7	116.4	3.2	344,903	1,000	0	1,000						13.2	30.0	17.6	39.3	0.003																							
									6	76.5	132	3.8	max miss=0, b_lg_size=0.05, adjust_pseudo_haploid=T	1,000	0	1,000						26.0	44.4	15.6	14.0	0.002																						
		3f [right]	14	2	0			6	102.4	171.9	4.2	613,509	1,000	0	1,000						0.1	0.9	9.8	89.2	0.112	yes																						
					6	102.4	171.9	4.2	613,509	same as above + minac2=2	4,446	0	2,785					0.1	6.7	40.3	52.8	0.056																										
					6	102.4	171.9	4.2	613,509	same as above + minac2=2	3,505	0	2,894					0.1	6.7	40.3	52.8	0.056																										

- [1] Counting at most one topology per iteration, not double counting topologies found in multiplier iterations
 - [2] a 1% gene flow from Loschbour to the Mongolia Neolithic group was dropped from the published model to decrease the complexity of the baseline model. That increased model LL slightly (by 0.5 log-units).
 - [3] a 1% gene flow from Loschbour to the Mongolia Neolithic group was dropped from the published model to decrease the complexity of the baseline model. That increased model LL slightly (by 2.4 log-units).
 - [4] OG was set for generating random starting graphs, but not for the "find_graphs" algorithm itself
- * here the direction of comparison is reversed since the published model is the highest-ranking among all models found

2529
2530

2531 **Table S1: Published graphs in the context of automatically found graphs.** We compared 22 different graphs from 8 publications to alternative graphs inferred on the same
2532 or very similar data; these *findGraphs* runs are highlighted in blue in the “iterations” column. In total, 51 *findGraphs* runs are summarized here since in some cases models more
2533 complex or less complex than the published one were explored and/or different population compositions were tested (see the “Topology search constraints and population
2534 modifications” column and footnotes for details). The columns with names in blue show various information on the published graphs or their modified versions and some
2535 properties of the published population sets. The columns with names in magenta show settings used for calculating f -statistics and for exploring the admixture graph space, and
2536 the number of SNPs used that depends on them. The columns with names in black summarize the outcomes of *findGraphs* runs, i.e., the properties of alternative model sets
2537 found.

2538 **Publication:** Last name of the first author and year of the relevant publication.

2539 **Figure in the original publication:** Figure number in the original paper where the admixture graph is presented.

2540 **Groups (populations):** The number of populations in each graph.

2541 **Singleton pseudo-diploid populations:** The number of populations in the graph composed of a single pseudo-diploid individual. Calculation of negative “admixture” f_3 -statistics is
2542 impossible for such populations since their heterozygosity cannot be estimated (see the text for details).

2543 **No. of negative f_3 -stats (allsnps: YES):** The number of negative f_3 -statistics among all possible f_3 -statistics for a given set of populations when all available sites are used for each
2544 statistic. If no negative f_3 -statistics exist for a set of populations, admixture graph fits are not affected by the “minac2=2” setting intended for accurate calculation of f -statistics for
2545 non-singleton pseudo-diploid groups.

2546 **Admixture events:** The number of admixture events in each graph.

2547 **Publ. model: log-likelihood (LL):** Log-likelihood score of the published graph fitted to the SNP set shown in the “SNPs used” column.

2548 **Publ. model: LL, median of bootstrap distr.:** Median of the log-likelihood scores of 100 or 500 fits of the published graph using bootstrap resampled SNPs.

2549 **Publ. model: worst residual (WR), SE:** The worst f -statistic residual of the published graph fitted to the SNP set shown in the “SNPs used” column, measured in standard errors
2550 (SE).

2551 **SNPs used:** The number of SNPs (with no missing data at the group level) used for fitting the admixture graph. For all case studies, we tested the original data (SNPs, population
2552 composition, and the published graph topology) and obtained model fits very similar to the published ones. However, for the purpose of efficient topology search we adjusted
2553 settings for f_3 -statistic calculation, population composition, or graph complexity as shown here and discussed in the text.

2554 **Settings for calculating f_2 -statistics:** Arguments of the *extract_f2* function used for calculating all possible f_2 -statistics for a set of groups, which were then used by *findGraphs* for
2555 calculating f_3 -statistics needed for fitting admixture graph models. See Methods for descriptions of each argument.

2556 **Topology search constraints and population modifications:** Constraints applied when generating random starting graphs and/or when searching the topology space.
2557 Modifications of the original population composition are also described in this column, where applicable.

2558 **Iterations:** The number of *findGraphs* iterations, each started from a random graph of a certain complexity. For each case study, *findGraphs* setups that were considered optimal
2559 are highlighted in blue in this column.

2560 **Iterations confirming published graph:** The number of iterations in which the resulting graph was topologically identical to the published graph. In the cases when the published
2561 model was irrelevant since more complex graphs were explored, “N/A” appears in this and subsequent columns. If less complex models were explored, the published model was
2562 still relevant since its version without selected admixture edges was tested.

2563 **Distinct alternative topologies found:** The number of distinct newly found topologies. If graph complexity was equal to (or less than) that of the published graph, the published
2564 topology (or its simplified version) is not counted here. If graph complexity exceeded that of the published graph, all newly found topologies are counted. If the published
2565 topology was recovered by *findGraphs*, the numbers in this column are shown in bold.

2566 **Significantly better fitting topologies:** The number of distinct topologies that fit significantly better than the published graph according to the bootstrap model comparison test
2567 (two-tailed empirical p -value <0.05). If the number of distinct topologies was very large, a representative sample of models (1/20 to 1/3 of models evenly distributed along the
2568 log-likelihood spectrum) was compared to the published one instead. These cases are marked as “a fraction of models tested” in this column. If model complexity was higher than
2569 that of the published model, model comparison was irrelevant and was not performed.

2570 **Non-significantly better fitting topologies:** The number of distinct topologies that fit non-significantly (nominally) better than the published graph according to the bootstrap
2571 model comparison test (two-tailed empirical p -value ≥ 0.05).

2572 **Non-significantly worse fitting topologies:** The number of distinct topologies that fit non-significantly (nominally) worse than the published graph according to the bootstrap
2573 model comparison test (two-tailed empirical p -value ≥ 0.05).

2574 **Significantly worse fitting topologies:** The number of distinct topologies that fit significantly worse than the published graph according to the bootstrap model comparison test
2575 (two-tailed empirical p -value ≥ 0.05).

2576 **Significantly better fitting topologies, %:** The percentage of distinct topologies that fit significantly better than the published graph according to the bootstrap model comparison
2577 test (two-tailed empirical p -value <0.05). If the number of distinct topologies was very large, a representative sample of models (1/20 to 1/3 of models evenly distributed along
2578 the log-likelihood spectrum) was compared to the published one instead, and the percentages in this and following columns were calculated on this sample.

2579 **Non-significantly better fitting topologies, %:** The percentage of distinct topologies that fit non-significantly (nominally) better than the published graph according to the
2580 bootstrap model comparison test (two-tailed empirical p -value ≥ 0.05).

2581 **Non-significantly worse fitting topologies, %:** The percentage of distinct topologies that fit non-significantly (nominally) worse than the published graph according to the
2582 bootstrap model comparison test (two-tailed empirical p -value ≥ 0.05).

2583 **Significantly worse fitting topologies, %:** The percentage of distinct topologies that fit significantly worse than the published graph according to the bootstrap model comparison
2584 test (two-tailed empirical p -value ≥ 0.05).

2585 **P-value best alternative vs. publ.:** An empirical two-tailed p -value of a test comparing log-likelihood distributions across bootstrap replicates for two topologies, the highest-
2586 ranking newly found topology and the published topology. In some cases, the highest ranking newly found topology (according to LL) has a fit that is not significantly better than
2587 that of the published model, but other newly found models fit significantly better despite having higher LL. P -values below 0.05 are highlighted in green.

2588 **Used in Table 1:** Here the *findGraphs* runs featured in **Table 1** are marked.

2589 **Table S2: Statistics for shotgun sequencing of individual I8726**

2590

Individual ID	I8726
Archaeological IDs	SHAR_201 (Grave 201)
Skeletal element	Petrous bone
Location	Seistan, Shahr-i-Sokhta, Iran
Archaeological context date	3100-3000 BCE
Latitude	30.649857
Longitude	61.400311
First publication of library	Narasimhan, Patterson et al. Science 2019
Library ID	S8726.E1.L1
HiSeqX10 lanes for shotgun sequencing	3
Molecular sex	Male
Mean coverage measured on 1240k autosomal targets	2.60306

2591

2592 **Supplementary Items Available as Separate Files**

2593

2594 **5 supplementary tables are available as separate files (Tables S3-S7)**

2595

2596 **13 supplementary figures are available as separate files (Figures S1-S13)**