

1 Applying f_4 -statistics and admixture graphs: theory 2 and examples 3

4 Mark Lipson^{1,2}

5 ¹Department of Genetics, Harvard Medical School, Boston, MA 02115, USA

6 ²Department of Human Evolutionary Biology, Harvard University, Cambridge, MA 02138, USA

7 Email: mlipson@genetics.med.harvard.edu
8

9 Abstract

10 A popular approach to learning about admixture from population genetic data
11 is by computing the allele-sharing summary statistics known as f -statistics. Com-
12 pared to some methods in population genetics, f -statistics are relatively simple, but
13 interpreting them can still be complicated at times. In addition, f -statistics can
14 be used to build admixture graphs (multi-population trees allowing for admixture
15 events), which provide more explicit and thorough modeling capabilities but are
16 correspondingly more complex to work with. Here, I discuss some of these issues
17 to provide users of these tools with a basic guide for protocols and procedures. My
18 focus is on the kinds of conclusions that can or cannot be drawn from the results of
19 f_4 -statistics and admixture graphs, illustrated with real-world examples involving
20 human populations.

21 Keywords: f -statistics, admixture graphs, admixture, parameter estimation

22 Introduction

23 f -statistics (Reich et al., 2009; Patterson et al., 2012) are a widely used toolkit for making
24 inferences about phylogeny and admixture from population genetic data, particularly in
25 humans. The statistics measure correlations in allele frequencies among sets of two, three,
26 or four populations. Observed values reflect degrees of shared ancestry and can serve as a
27 means for testing hypotheses regarding population split orders and past gene flow events
28 under historical models.

29 As compared to some other common methods in population genetics, f -statistics are
30 quite simple and flexible, but interpreting them is not always straightforward. Addition-
31 ally, one of the primary applications of f -statistics is in building admixture graphs (i.e.,
32 phylogenetic trees augmented with admixture events) with more than four populations,
33 which introduces a greater level of complexity. In this note, I hope to clarify some of these
34 potential difficulties and provide a range of tips for practitioners. Some of the topics have
35 been addressed previously but are covered here as well for the sake of completeness.

36 f -statistics and admixture

37 Basic definitions and properties

More complete introductions to f -statistics have been published elsewhere (Reich et al., 2009; Patterson et al., 2012; Lipson et al., 2013; Peter, 2016; Soraggi and Wiuf, 2019), but the following are some basics that are used in other sections of the paper. The most general definition is that of the f_4 -statistic $f_4(A, B; C, D)$, which measures the average correlation in allele frequency differences between (i) populations A and B and (ii) populations C and D (i.e., $(p_A - p_B) * (p_C - p_D)$), for allele frequencies p , typically

averaged over many biallelic single-nucleotide polymorphisms [SNPs]). This f_4 -statistic is the same as the (perhaps more familiar) D -statistic up to a normalization factor. If the four populations are related by the (unrooted) phylogeny $((A, B), (C, D))$, then the expected value of $f_4(A, B; C, D)$ will be zero, while the expected values of $f_4(A, C; B, D)$ and $f_4(A, D; B, C)$ will be positive. (When I refer to expectations of f -statistics, I mean with respect to the random noise in real data—typically assumed to be normally distributed—caused by sampling finite numbers of independent SNPs and individuals.) Simple algebra shows that

$$\begin{aligned} f_4(A, B; C, D) &= f_4(C, D; A, B), \\ f_4(A, B; C, D) &= -f_4(B, A; C, D) = -f_4(A, B; D, C), \\ f_4(A, B; C, D) &= f_4(A, C; B, D) + f_4(A, D; C, B). \end{aligned}$$

38 The other two basic definitions are of the f_2 - and f_3 -statistics, which can be formulated
39 as $f_2(A, B) = f_4(A, B; A, B)$ and $f_3(A; B, C) = f_4(A, B; A, C)$.

40 The most important usage for f -statistics is in the context of admixture. If a popu-
41 lation C has a mixture of ancestry derived from sources C' and C'' in proportions α and
42 $(1 - \alpha)$, then in expectation,

$$43 \quad f_4(A, B; C, D) = \alpha f_4(A, B; C', D) + (1 - \alpha) f_4(A, B; C'', D).$$

44 Expected values of f -statistics can be visualized in terms of overlapping paths in an
45 admixture graph (Fig. 1; see also Patterson et al. (2012); Peter (2016); Soraggi and Wiuf
46 (2019)). In the case of admixture, the above equation can be used to derive the expectation
47 in terms of a weighted sum of path-overlaps involving each source (Fig. 1C). Thus, if C
48 is admixed, the typical expected value of $f_4(A, B; C, D)$ will be a branch length times a
49 mixture proportion (Fig. 1C).

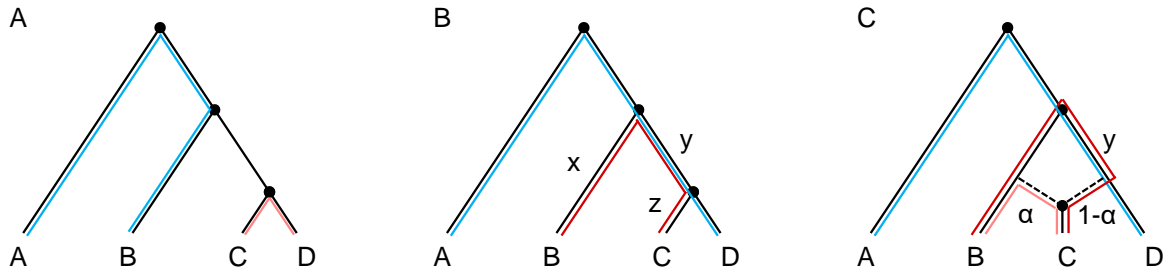


Figure 1. Expected values of f_4 -statistics under specified admixture graph models. (A) The expected value of $f_4(A, B; C, D)$ is given by the intersection between the path from A to B with the path from C to D . Under the model shown, $E[f_4(A, B; C, D)] = 0$. (B) The expected value of $f_4(A, D; B, C)$ is given by the intersection between the path from A to D with the path from B to C . Under the model shown, $E[f_4(A, D; B, C)] = y$. (C) With population C admixed, the path from B to C can be decomposed into two components. Under the model shown, with a proportion of α B -related ancestry and $1 - \alpha$ D -related ancestry, the former yields a path (lighter red) that has a weight of α but does not intersect the path from A to D , while the latter yields a path (darker red) that has a weight of $1 - \alpha$ and intersects the path from A to D over the branch with length y . In total, $E[f_4(A, D; B, C)] = (1 - \alpha)y$.

50 Unlike F_{ST} (and normalized D -statistics, at least approximately), the values of f -
 51 statistics (including branch lengths in admixture graphs that are defined in f -statistic
 52 units, as in Fig. 1) depend on the absolute allele frequencies of the SNPs used to calculate
 53 them (cf. Lipson et al. (2013)). For example, adding fixed sites to the SNP set will shrink
 54 f -statistics toward zero. As a result, when comparing multiple f -statistics, it is important
 55 that each one should be computed on the same set of SNPs (or as similar as possible). In
 56 applications involving ancient DNA, where missing data is common, I typically make the
 57 assumption that the SNPs covered for each individual or population are a random subset
 58 with respect to allele frequency. By contrast, comparisons across different genotyping
 59 arrays are likely to be biased.

60 **Interpreting non-zero f_4 -statistics**

61 If a set of four populations are unadmixed relative to each other, then some permutation
62 of them will yield an f_4 -statistic of zero (in expectation), as in Fig. 1A. Equivalently, if all
63 three permutations of f_4 -statistics for a certain set of four populations are (significantly)
64 non-zero, then at least one of the populations must be admixed; this is one of the most
65 common signals of admixture used in the literature. In this paper, I will use the example
66 of a quartet consisting of four present-day human populations: Mixe (from Mexico), Han
67 Chinese, French, and Baka (hunter-gatherers from Cameroon). The common ancestral
68 population of all Native Americans is known to have been admixed with approximately
69 70% ancestry from an eastern Eurasian lineage and 30% from a western Eurasian lineage
70 (Fig. 2) (Raghavan et al., 2014). Thus, in the context of this quartet, Mixe can be modeled
71 as admixed with ancestry related to Han ($\sim 70\%$) and to French ($\sim 30\%$). I computed the
72 three possible f_4 -statistics for the quartet and obtained significantly non-zero values, with
73 the signs as expected based on the known history (Table 1). (These and all results in
74 the paper are computed from previously published whole-genome sequence data (Mallick
75 et al., 2016; Fan et al., 2019), on a set of ~ 1.1 million autosomal SNPs (Mathieson et al.,
76 2015), using the implementation in ADMIXTOOLS (Patterson et al., 2012), including
77 standard errors estimated by block jackknife.)

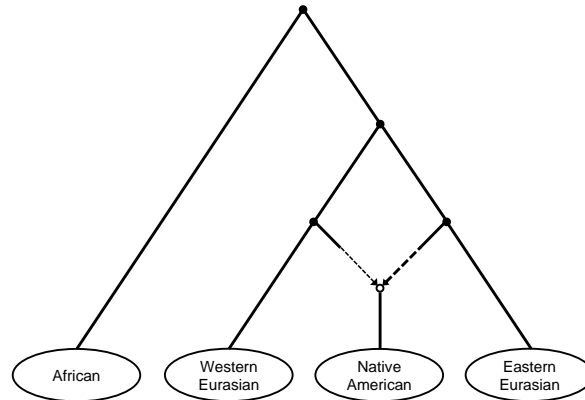


Figure 2. Major human lineages used for examples in the paper, represented by Baka (African), French (western Eurasian), Mixe (Native American), and Han (eastern Eurasian). Setting aside other complexities in the histories of these populations, the admixture event being modeled involves eastern and western Eurasian lineages contributing ancestry to Native Americans (Raghavan et al., 2014). See Figs. 3A and 5A for fitted models using this correct topology.

Table 1. Observed f_4 -statistics (values and Z -scores for difference from zero) for the example populations.

Populations				$f_4(\mathbf{A}, \mathbf{B}; \mathbf{C}, \mathbf{D})$	
A	B	C	D	Value	Z -score
Mixe	Baka	Han	French	0.011	27.1
Mixe	French	Han	Baka	0.013	35.8
Mixe	Han	Baka	French	-0.0025	-8.9

78 In this case, there is prior knowledge available about the admixture in Mixe, but in
 79 general, without additional information, the existence of such a quartet does not identify
 80 which of the four populations is admixed. Here, for example, it could also be that Han is
 81 admixed with most of its ancestry related to Mixe but a small amount related to Baka,
 82 and likewise for the other two (see further discussion in the admixture graph sections
 83 below). In real-world applications, it can also be true that more than one population is
 84 admixed, making the interpretation more complicated. Sometimes, in fact, two admixture

85 events together can cause an f_4 -statistic to be close to zero and thereby mask the signal
86 of admixture (at first glance).

87 Another observation is that as depicted in Fig. 1, f_4 -statistics are not only zero or
88 non-zero but also carry quantitative information about amounts of shared drift between
89 populations. One implication is that populations sharing more drift (i.e., yielding longer
90 intersecting paths in an admixture graph) will have greater-magnitude f_4 -statistics asso-
91 ciated with them. For example, in the trees of Fig. 1B–C, if one replaced population D
92 with a population D' that split halfway between D and the root of the tree, then the
93 expected magnitude of $f_4(A, B; C, D')$ would be smaller, since the length of the shared
94 drift branch would now be less than y . As a result, under the model in Fig. 1C, one could
95 use the fact that $f_4(A, B; C, D) > f_4(A, B; C, D')$ to conclude that D is a better proxy
96 than D' for the ancestry in C (the component with proportion $1 - \alpha$). However, this pro-
97 cedure is complicated by the fact that if the D -related source was in fact itself admixed,
98 with ancestry related to X and Y , then the f_4 -statistic can sometimes be maximized by
99 X or Y instead of by D , even though one would consider D to be a better proxy (Pickrell
100 et al., 2014). It is also good to remember that if a certain signal is weak compared to
101 the noise in the data—for example, if one were testing for admixture in C and the shared
102 drift branch length y was short—then one may not have enough power to identify it.

103 Finally, f -statistics can be subject to certain kinds of biases and batch effects (to
104 varying degrees, as with other methods) arising from SNP ascertainment, sample type
105 and processing (ancient versus present-day, sequencing platform, etc.), and other aspects
106 of the data, so it is important to keep such factors in mind when interpreting results.
107 For ancient DNA data, challenges include C-to-T errors induced by postmortem deami-
108 nation (Hofreiter et al., 2001), as well as short fragment lengths and (often) low coverage,
109 which can exacerbate reference bias (Günther and Nettelblad, 2019). All of these effects

110 can cause ancient individuals to appear artificially closely related to one another and to
 111 certain other populations (e.g., deep outgroups). In general, statistics $f_4(A, B; C, D)$ in
 112 which A and C share a data type and B and D share a different data type are most prone
 113 to this kind of artifact.

114 **Admixture graphs: modeling and inference procedure**

115 **Fitting an admixture graph with *qpGraph***

116 In addition to their stand-alone usage, f -statistics can serve as a means to fit admixture
 117 graphs from allele frequency data. (Other kinds of statistics can also be used to fit
 118 admixture graphs, but I will not discuss such methods in detail here; see Discussion.) In
 119 this context, an admixture graph consists of an ordering of population splits, positions
 120 of admixture events, branch length parameters, and mixture proportions. Given the first
 121 two, the third and fourth can be inferred by solving a system of equations (linear in
 122 terms of the branch lengths) in which observed f -statistic values are matched to their
 123 expectations in terms of the model parameters. For example, one such equation for the
 124 model in Fig. 1B would be $f_2(B, C) = x + y + z$. With n populations, there are $3 \times \binom{n}{4}$
 125 possible f_4 -statistics, $3 \times \binom{n}{3}$ possible f_3 -statistics, and $\binom{n}{2}$ possible f_2 -statistics, but many
 126 of these are linearly dependent; for example, $f_4(A, B; C, D) = f_3(A; B, D) - f_3(A; B,$
 127 $C)$. In fact, there are a total of $\binom{n}{2}$ linearly independent f -statistic equations, or in other
 128 words, f -statistics form a vector space of dimension $\binom{n}{2}$. Possible choices of basis include
 129 (1) the set of all f_2 -statistics, and (2) the set of all f_2 - and f_3 -statistics with a given
 130 population in the first position.

131 The software I typically use to build admixture graphs is *qpGraph* (also referred to as
 132 ADMIXTUREGRAPH) (Patterson et al., 2012). In *qpGraph*, the user manually specifies

133 the topology of the model, and the program then solves for the optimal values of the
134 parameters. In theory, one might wish to search the entire space of all topologies and
135 parameter values (for a given number of admixture events) to find the best-fitting model,
136 but the size of the space (exponential in the number of populations) makes this impractical
137 for larger graphs (Leppälä et al., 2017). The set of basis statistics used for fitting is the
138 set (2) alluded to in the previous paragraph, with the first population listed in the input
139 file as the “base” population.

140 In its standard mode, *qpGraph* attempts to minimize the quantity $S(G) = 1/2(g -$
141 $f)'Q^{-1}(g - f)$, known as the “score” of the model, where f is the vector of observed basis
142 f -statistics (of length $\binom{n}{2}$), g is the vector of predicted f -statistics under the model, and
143 Q is the (estimated) covariance matrix of the statistics. Assuming multivariate normal
144 errors, the score gives the negative log-likelihood of the model; it measures the total
145 amount by which the system of f -statistic equations (one for each basis statistic) fails to
146 be satisfied, taking into account the empirical correlation among the statistics (see also
147 the next section on fit quality). To help insure that Q^{-1} does not become unstable, one
148 can use the “diag” input parameter to add a small number (“diag: 0.0001” works well
149 in my experience, but smaller values may be sufficient as well) to the diagonal entries of
150 Q . The program can also be run using simple least-squares optimization without the Q
151 matrix by specifying “lsqmode: YES,” but in this case highly correlated statistics will be
152 treated as independent for the sake of the fitting, and the score will no longer represent
153 a log-likelihood, both of which make the full objective function preferable. Other input
154 parameters I typically set are “outpop: NULL” (meaning no specified outgroup population
155 in which SNPs are required to be polymorphic) and “lambdascale: 1” (leaving the f -
156 statistics in typical units rather than scaling into approximate F_{ST}). More extensive
157 descriptions of the *qpGraph* software can be found in Patterson et al. (2012) and in the

158 ADMIXTOOLS package repository (<https://github.com/DReichLab/AdmixTools>), and
159 of the f -statistic-based admixture graph inference process more generally in Lipson et al.
160 (2013); Leppälä et al. (2017).

161 By default, *qpGraph* utilizes the set of SNPs that have genotype calls for at least one
162 individual in each population in the model. With low-coverage data (for example, in some
163 ancient DNA applications), this can result in losing the majority of the sites in the initial
164 data set. The program allows an option to use all SNPs instead (“allsnps: YES” or “use-
165 allsnps: YES,” in which case each basis statistic is computed on as many sites as possible
166 for the two or three populations involved), but this mode can give unreliable results, in
167 particular when the base population is highly diverged from the other populations in the
168 model. To the best of my knowledge, this effect is caused by greater absolute noise when
169 estimating larger-magnitude basis statistics, such that the small relative fluctuations in
170 empirical f -statistics caused by modest changes in the SNP set become substantial in
171 the context of the admixture graph. In my own work, my preference has always been
172 to avoid using the all-SNPs option. If this causes an undesirable loss of coverage, then
173 the best approach given the current implementation of *qpGraph* is probably to set as the
174 base a population that (a) is not highly diverged from the others in the model, and (b)
175 preferably has multiple individuals with diploid data (again to reduce the magnitudes of
176 the statistics). Research is currently underway aiming to develop an improved all-SNPs
177 methodology.

178 **Parameters and constraints**

179 An important consideration is whether the system of equations used to infer the param-
180 eters of an admixture graph is over- or under-determined. As mentioned above, a model
181 with n populations has $\binom{n}{2}$ linearly independent constraints (i.e., equations). In the ab-

182 sence of admixture, there are $2n - 3$ parameters, which is the number of branches in an
183 unrooted binary tree with n leaf nodes (with the settings I have described, *qpGraph* results
184 should not depend on where the root of a graph is specified). Converting a population
185 from unadmixed to admixed adds two parameters: one for the mixture proportion and
186 one for the split position of the new source of ancestry. Thus, with a admixture events,
187 the total number of free parameters is $2n + 2a - 3$. One point to note is that in the case
188 of an admixed population with two unsampled sources (which is the typical scenario), the
189 three branch lengths surrounding the admixture event (in Fig. 3A, from the node “East1”
190 to “East2,” from “West1” to “West2,” and from “pAM1” to Mixe) cannot be determined
191 individually but instead form a single compound parameter $\alpha^2x + (1 - \alpha)^2y + z$ (where α is
192 the mixture proportion, x and y are the branch lengths to the two corresponding sources,
193 and z is the terminal branch length). The only exception (to my knowledge) is the case
194 in which at least three populations are included that can be modeled as having different
195 proportions of ancestry from the same two sources, which allows the branch lengths to be
196 solved for individually.

197 Even if the inequality $\binom{n}{2} \geq 2n + 2a - 3$ is satisfied for an admixture graph as a
198 whole, there can be some parameters that are not uniquely determined because of rep-
199 etition across the different equations caused by multiple populations in phylogenetically
200 equivalent positions. Further discussion of this phenomenon can be found in the example
201 sections below. Additionally, having sufficient constraint to estimate parameters is not
202 entirely a yes-or-no proposition. A model can have enough populations in distinct posi-
203 tions to be able to estimate a mixture proportion, but if two of the populations are only
204 slightly separated, then the precision of the estimate will generally be lower. Similarly,
205 if one of the populations providing the constraint is itself admixed, then the power will
206 often be reduced.

207 **Fit quality**

208 To my knowledge, no absolute measure of model fit has been developed for admixture
209 graphs, but there are several ways to evaluate how well a given model fits the data
210 (this is an area of active study; see also Lipson and Reich (2017); Lipson et al. (2017);
211 Leppälä et al. (2017); Flegontov et al. (2019); Shinde et al. (2019); Lipson et al. (2020)).
212 The following discussion is tailored for *qpGraph*, but the ideas also apply more generally.
213 First, the program returns a list of residual poorly-predicted f -statistics and their Z -scores
214 (drawn from the set of all possible f -statistics, not only those in the basis), which can
215 give a good sense for the performance of the model and some idea of which populations
216 are responsible for the greatest inaccuracies. There is no general rule for what threshold
217 constitutes a significantly non-zero residual; the situation is complicated because there
218 are many statistics being tested simultaneously, but many of those are also correlated
219 with each other.

220 Deviations between model predictions and the observed data can be caused either by
221 an incorrectly specified topology or un-modeled admixture. In the first case, assuming that
222 the program does not get stuck at a local optimum, it will try to move the populations as
223 close as possible to their correct positions but will be constrained by the input topology.
224 Thus, an incorrectly specified split order usually manifests as an inferred length-zero
225 internal branch; when such branches (i.e., trifurcations) appear in the results, the order of
226 splits should be adjusted and re-tried. (The default *qpGraph* visualization output rounds
227 branch lengths to the nearest integer, so some non-zero-length but very short branches
228 may initially appear as zero.) As noted in the f -statistics section above, however, one
229 may not have sufficient power to resolve short branches, so some sets of three lineages may
230 be found to be statistically consistent with forming a trifurcation, with all three possible
231 split orders having similar fit quality.

232 In the case of un-modeled admixture, the observed deviations could potentially reflect
233 admixture in one of multiple different populations. Often one can gain information by
234 examining the full list of residuals and noting which populations occur repeatedly. An-
235 other approach is to remove one population from the model and see if the fit improves,
236 although even if it does, that could imply either that the population in question had un-
237 modeled admixture or that it provided a constraint enabling the detection of un-modeled
238 admixture among the other populations.

239 The score of the final graph is also returned as an output from the program, so it can
240 be used to compare the fit quality of different models with the same set of populations,
241 preferring the one with the lower score. (If the equations being fit were independent,
242 then one could apply a chi-squared test for the overall fit, but in practice they are heavily
243 correlated. *qpGraph* returns a naive degrees of freedom count and p -value alongside the
244 score, but they are not well calibrated.) As above, while this approach provides a useful
245 heuristic, evaluating statistical significance is complicated, and I do not have a rigorous
246 set of recommendations. One recent direction that seems promising is using the score to
247 compare alternative models with the same populations and same number of admixture
248 events. In that case, the score difference can be interpreted in an AIC/BIC framework,
249 with the likelihood difference as a Bayes factor (Leppälä et al., 2017; Flegontov et al., 2019;
250 Shinde et al., 2019). The same idea could also be applied in cases with unequal numbers
251 of free parameters—for example, adding one admixture event and testing whether the
252 score improvement is significant. However, defining the change in degrees of freedom is
253 not straightforward in this situation: as noted above, a new admixture event creates two
254 additional parameters in the model, but that does not account for whether the admixture
255 comes from a pre-specified source or from a source that is allowed to be located anywhere
256 in the graph. Finally, the score can additionally be used to compute confidence intervals

257 on parameters (by considering the likelihood as a function of a single branch length or
258 mixture proportion value), although it is worth keeping in mind that the results are
259 model-dependent.

260 **Admixture graphs: examples**

261 One of the strengths of f -statistic-based admixture graphs is that they are computation-
262 ally tractable enough that programs such as *qpGraph* can accommodate a large number
263 of populations and admixture events. Sometimes though it can be difficult to digest all
264 of the information in large admixture graph models and to analyze their behavior. For-
265 tunately, the main principles of admixture graph fitting can be illustrated with simpler
266 examples, which, in particular, carry over directly to larger models by considering subsets
267 of four and five populations.

268 **Four populations**

269 The first examples I will present are four-population admixture graphs containing Mixe,
270 Han, French, and Baka. Given the observed non-zero f_4 -statistics in Table 1, there must
271 be at least one admixture event present in order to fit the data. However, in light of the
272 discussions above about determining which population is admixed and about parameters
273 and constraints in admixture graphs, it would be expected that these models should be
274 insufficiently constrained to determine which population is admixed. Indeed, they have
275 $\binom{4}{2} = 6$ constraints but $2(4) + 2(1) - 3 = 7$ free parameters. Confirming this expectation,
276 perfectly fitting models (i.e., sets of branch length and mixture proportion parameters
277 such that the six basis f -statistics are predicted exactly, yielding $S(G) = 0$) can be
278 obtained with Mixe specified as admixed (Fig. 3A) as well as with any of the other three

279 populations (incorrectly) specified as admixed instead (Fig. 3B–D).

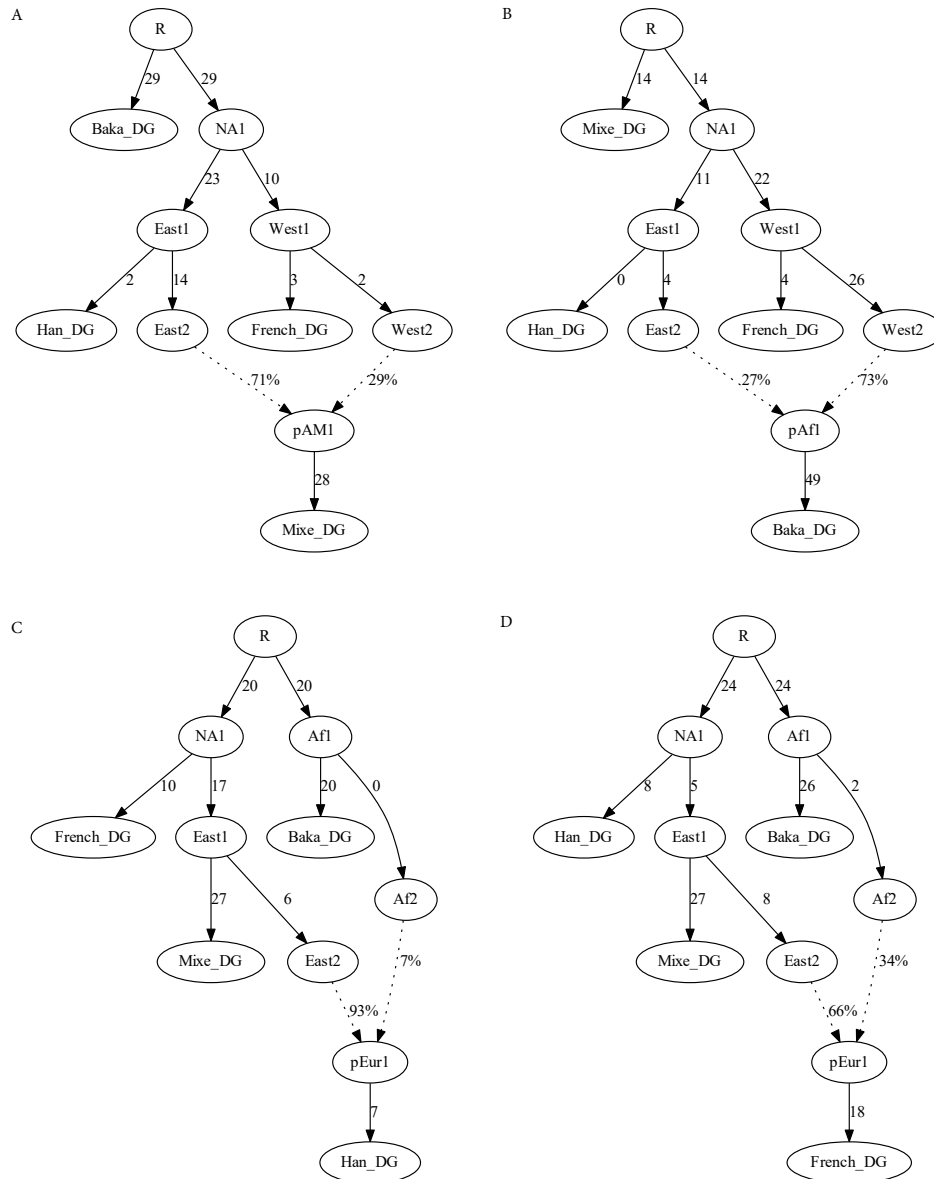


Figure 3. Four-population admixture graphs modeling (A) Mixe, (B) Baka, (C) Han, or (D) French as admixed. All four versions provide perfect fits to the data (exact agreement between observed and predicted f -statistics). In this and all following figures, branch lengths (in f -statistic units, multiplied by 1000) are rounded to the nearest integer.

280 Interestingly, in some scenarios, the admixed population can be determined even with
 281 only four populations in the model: if a negative f_3 -statistic can be formed for some
 282 triple, then the population in the first position of the statistic (i.e., population A if $f_3(A;$
 283 $B, C) < 0$) must be admixed. To give an example, I replaced Mixe with Kyrgyz in the
 284 four-population model. With Kyrgyz modeled as admixed, the fit is perfect as before
 285 (Fig. 4A). With Baka modeled as admixed, however, the fit is very poor, with residuals
 286 up to $Z = 27$ (Fig. 4B). The most extreme residual is the statistic $f_3(\text{Kyrgyz}; \text{Han},$
 287 $\text{French})$, which has an observed value of -0.0064 ($Z = 27$ for difference from zero) and
 288 can only be negative if Kyrgyz is admixed (i.e., in the position of the test population in
 289 a “three-population test” for admixture (Reich et al., 2009; Patterson et al., 2012)).

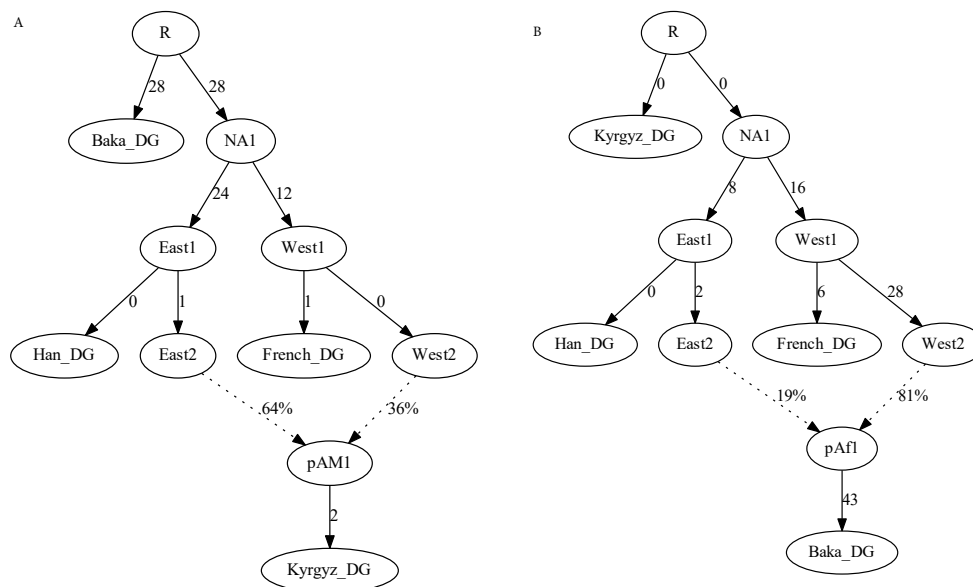


Figure 4. Four-population admixture graphs with Kyrgyz in place of Mixe, modeling either (A) Kyrgyz or (B) Baka as admixed. The first provides a perfect fit to the data, whereas the second has residuals up to $Z = 27$.

290 Another note is that in these examples, I have been focusing on the primary signal
 291 of deep eastern/western Eurasian admixture in Mixe. The other populations are also

292 admixed in their own ways; for example, all of the non-Africans have small proportions
293 of Neanderthal ancestry, and Baka are admixed with ancestry related to nearby Bantu-
294 speaking farmers (Fan et al., 2019). However, the first signal is not evident in the data
295 without deeper outgroups present, and the second without other African populations.
296 Conversely, if the model contained several sub-Saharan African populations plus Mixe as
297 the lone non-Africans, then the primary signal in our examples here would not be visible.
298 In some ways, this inability to detect certain admixture events is beneficial, as it means
299 that models can be constructed so as to focus on events of interest while ignoring some
300 that are outside the desired scope of the work.

301 **Five populations**

302 In general, in order to be able to solve for the parameters of an admixture graph including
303 one admixture event, it is necessary to use at least five populations, providing $\binom{n}{2} = 10$
304 constraints for the $2n + 2a - 3 = 9$ free parameters. Concurrently, in contrast to the four-
305 population examples above, having five populations present allows one to determine which
306 of the populations is admixed, as long as the topological relationships of the populations
307 are all unique relative to the true mixing sources. More detail on this last point can
308 be found elsewhere (Pease and Hahn, 2015; Lipson and Reich, 2017). A simple version
309 of this statement is that, at least in the case of a single admixture event, one four-
310 population subset will be unadmixed, whereas the other four subsets will include the
311 admixed population. Similarly, in order to solve for a given mixture proportion in a larger
312 graph, there must four populations present (aside from the admixed one in question)
313 in distinct positions, yielding a non-redundant five-population subgraph; having three
314 populations in distinct positions allows one to detect the signal of admixture but not to
315 determine the proportion uniquely.

316 As an example, I added Ulchi (from the Amur River Basin of northeastern Asia)
317 as a fifth population alongside the four from above. Ulchi splits closer to the eastern
318 Eurasian source population for Mixe than does Han, which provides the additional degree
319 of constraint. The five-population model is a good fit to the data, but not a perfect one
320 ($Z = 1.9$ for the most significant residual; Fig. 5A). By contrast, if Baka are modeled as
321 admixed instead of Mixe, the fit is poor ($Z = 4.7$; Fig. 5B). I also show an example where
322 the topology is incorrectly specified, with Han closer than Ulchi to the eastern Eurasian
323 source population for Mixe (Fig. 5C); this version fits poorly ($Z = 5.7$), and the branch
324 connecting the split positions of Ulchi and Han collapses to length zero. If I add a second
325 admixture event into the models in Figs. 5A–B, this creates more free parameters (11)
326 than constraints, and indeed there are choices of the parameters that yield perfect fits,
327 even with Mixe modeled as unadmixed (not shown).

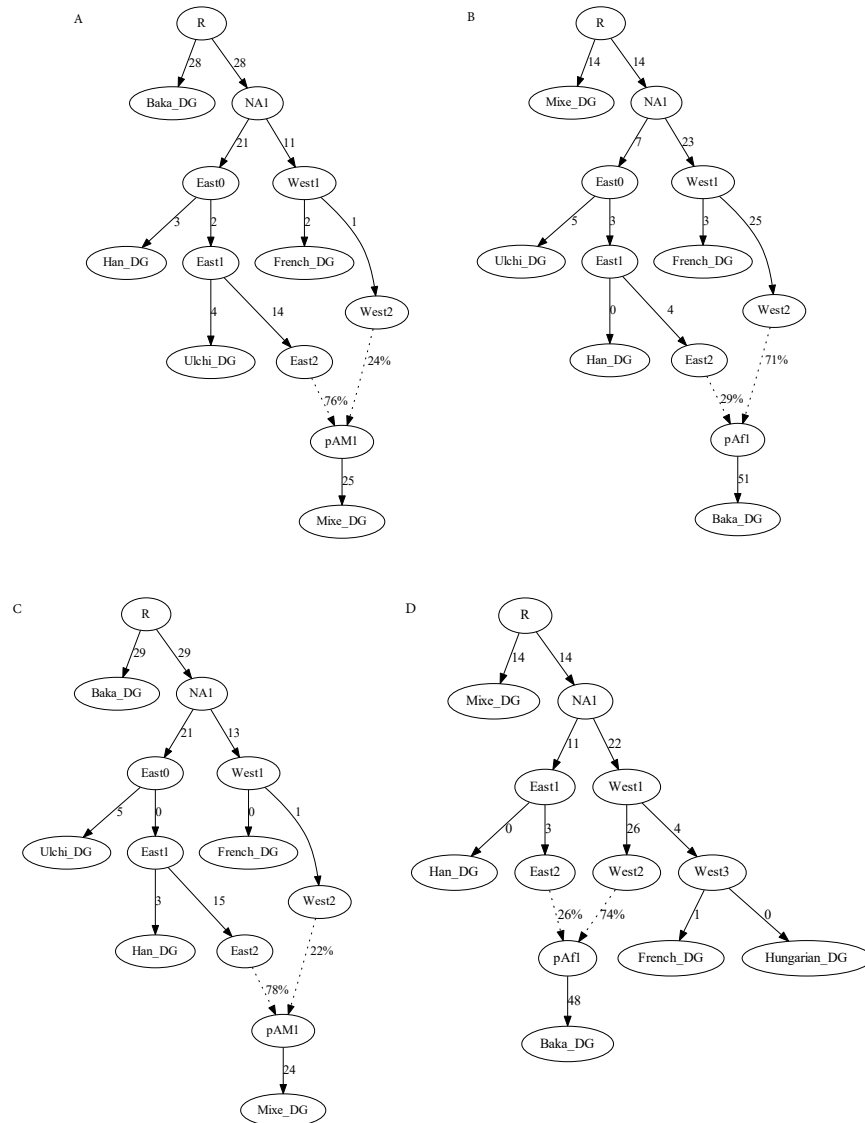


Figure 5. Five-population admixture graphs. (A) Standard four-population example plus Ulchi; all f -statistics are predicted to within 1.9 standard errors of their observed values. (B) Same five populations, but with Baka modeled as admixed; residual statistics are present up to $Z = 4.7$ (C) Same five populations, with Mixe modeled as admixed, but with the positions of Han and Ulchi reversed; residual statistics are present up to $Z = 5.7$. (D) Original four populations plus Hungarian, with Baka modeled as admixed; all f -statistics are predicted to within 1.2 standard errors of their observed values.

329 ability to infer uniquely optimal parameter values. In the four-population example model,
330 the initial estimate of eastern Eurasian ancestry in Mixe was 71%, but with the proportion
331 manually set at 75%, the fit is still perfect (Fig. 6A). Outside of a certain range of mixture
332 proportions (dependent on the values of the branch lengths), the fit will become worse, but
333 within a finite interval, the likelihood is entirely flat. In terms of f_4 -statistics, the observed
334 non-zero value is being fit as equal to a branch length in the admixture graph times the
335 mixture proportion (as in Fig. 1C), but without additional constraint, that product can
336 remain the same while the branch length and mixture proportion covary (where the range
337 is determined by bounds on the individual parameter values, e.g., positivity). With five
338 populations, however, there is a unique optimal solution; for example, if I set the mixture
339 proportion at 70% eastern Eurasian ancestry (as compared to the point estimate of 76%
340 in the five-population model), there are residuals up to $Z = 2.6$ (Fig. 6B), and the score is
341 more than 10 units worse. Even in the example above with Kyrgyz (i.e., a four-population
342 model where the admixed population can be determined because of a negative f_3 -statistic;
343 Fig. 4), the parameters remain not uniquely determined.

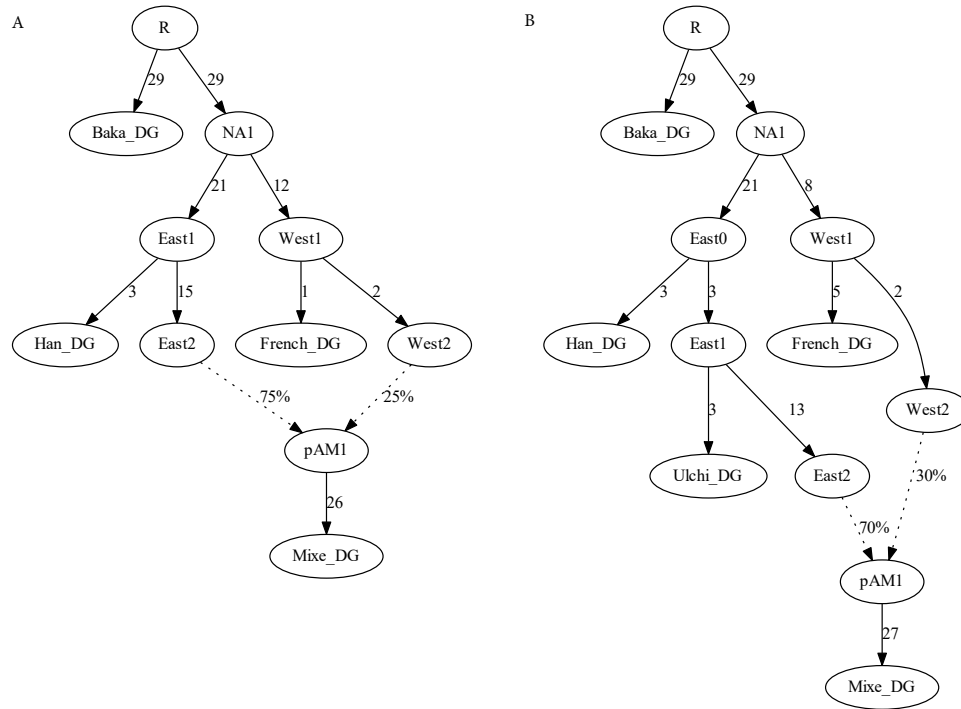


Figure 6. Admixture graphs with pre-specified mixture proportion parameters. (A) Four-population model, with the proportion locked at 75%; the fit is perfect. Note that the branch lengths shift slightly relative to Fig. 3A. (B) Five-population model, with the proportion locked at 70%; residual statistics (indicating a need for more eastern Eurasian ancestry in Mixe) are present up to $Z = 2.6$.

344 Finally, in Fig. 5D, I show a model with the original four populations plus Hungarian
 345 instead of Ulchi. Although there are five populations present, French and Hungarian can
 346 be modeled as sister groups, so equations relating parameters in the graph to statistics
 347 of the form $f_2(\text{French}, X)$ and $f_2(\text{Hungarian}, X)$ are linearly dependent (up to their
 348 terminal branch lengths) and hence do not contribute fully independent constraints. This
 349 can be seen in the results, as Baka can successfully be modeled as the admixed population
 350 (with residuals up to $Z = 1.2$ reflecting small observed asymmetries between French and
 351 Hungarian). This contrasts with Ulchi, which has a distinct phylogenetic position from
 352 Han (relative to the other populations in the model) and thus adds new constraints

353 (although it is worth noting again that a population with only a slightly different position
354 adds constraint but only weakly).

355 Discussion

356 Most of the material in this paper pertaining to admixture graphs has been presented
357 from the perspective of the *qpGraph* software, but other methods are also available, using
358 both different kinds of data and different fitting schemes. At the level of mathematical
359 formulation, the results have assumed that models are fit based on a distance metric
360 (specifically, f -statistics). As an alternative example, the *TreeMix* algorithm (Pickrell
361 and Pritchard, 2012) is based on a maximum-likelihood framework in terms of allele fre-
362 quency covariances, although the information captured is the same; see Peter (2016) for
363 the equivalence and a thorough exploration of alternative interpretations of f -statistics in
364 terms of population genetic models. There are also methods that use richer summaries of
365 the data (for example, the full joint allele frequency spectrum) to infer more complicated
366 demographic models that are similar in form, or in some cases essentially identical, to ad-
367 mixture graphs—for example, $\partial a \partial i$ (Gutenkunst et al., 2009), G-PhoCS (Gronau et al.,
368 2011), fastsimcoal2 (Excoffier et al., 2013), and momi2 (Kamm et al., 2019). The mathe-
369 matical underpinnings of such methods are quite different from those based on f -statistics,
370 and so the results presented here do not pertain to them. The choice of which program to
371 use can depend on aspects of the particular application such as the data set (e.g., number
372 of populations, whole-genome sequencing versus genotyping array, etc.) and the desired
373 level of complexity and parametrization. Even more generally, of course, numerous other
374 approaches exist to model population genetic structure beyond phylogenetic trees with
375 gene flow. While it may sometimes be possible to evaluate empirically the suitability of

376 an admixture graph for a given problem—for example, by exploring whether any graph
377 of a reasonable size provides a good fit to the data—the choice of model is ultimately at
378 the discretion of the analyst.

379 Within the class of f -statistic-based (or equivalent) admixture graph methods, there
380 are different approaches to automation and the selection of which populations to model as
381 admixed. *qpGraph* leaves the choice of how many admixture events to include (and which
382 populations are admixed) up to the user; some guidelines pertaining to this choice have
383 been discussed above. For smaller models, it can also be possible to search some or all of
384 the full graph space (Shinde et al., 2019) to determine best-fitting topologies for a given
385 number of admixture events (for example, using the similar *admixturegraph* R implemen-
386 tation (Leppälä et al., 2017) and *AdmixtureBayes* (Nielsen, 2018); other techniques are
387 the subject of ongoing work). *MixMapper* (Lipson et al., 2013) provides an intermediate
388 level of automation by attempting to infer an unadmixed sub-model and then fitting one
389 or two admixed populations onto this scaffold. With a small set of populations, this can
390 sometimes be a useful approach, but it can largely be recapitulated within *qpGraph*, and
391 the software does not support large models with more admixture events. At the most
392 automated end of the spectrum is *TreeMix* (Pickrell and Pritchard, 2012), which only
393 asks the user to supply the list of populations and the number of admixture events and
394 then returns a single inferred model. The advantage of this strategy is that the program
395 does all of the work of building the graph, which is especially useful if one has limited
396 prior knowledge about the populations. The main drawback, in my view, is that the way
397 the program builds the graph is by starting with an optimal mixture-free tree and then
398 adding admixture events to account for deviations between the predictions of the tree
399 model and the observed data. Depending on the true histories of the populations, this
400 approach can be successful, but it can also increase the chances of falling into local optima

401 imposed by the initial tree (especially if many populations are admixed; see (Lipson et al.,
402 2013)). Additionally—as in other methods—the choice of how many admixture events to
403 include, which can sometimes be difficult, is still left to the user.

404 In my experience, I have found f -statistics and admixture graphs to be very useful
405 tools for learning about phylogeny and admixture. I hope that this guide will help others
406 to get the most out of these tools in a range of real-world applications.

407 Acknowledgments

408 I would like to thank David Reich, Vagheesh Narasimhan, Nick Patterson, Robert Maier,
409 Iosif Lazaridis, and Pavel Flegontov for helpful discussions, and three anonymous reviewers
410 for comments.

411 References

- 412 Excoffier, L., Dupanloup, I., Huerta-Sánchez, E., Sousa, V. C., and Foll, M. (2013).
413 Robust demographic inference from genomic and SNP data. *PLoS Genet.*, 9(10).
- 414 Fan, S., Kelly, D. E., Beltrame, M. H., Hansen, M. E., Mallick, S., Ranciaro, A., Hirbo,
415 J., Thompson, S., Beggs, W., Nyambo, T., et al. (2019). African evolutionary his-
416 tory inferred from whole genome sequence data of 44 indigenous African populations.
417 *Genome Biol.*, 20(1):82.
- 418 Flegontov, P., Altınışık, N. E., Changmai, P., Rohland, N., Mallick, S., Adamski, N.,
419 Bolnick, D. A., Broomandkoshbacht, N., Candilio, F., Culleton, B. J., et al. (2019).
420 Palaeo-Eskimo genetic ancestry and the peopling of Chukotka and North America.
421 *Nature*, 570:236–240.
- 422 Gronau, I., Hubisz, M., Gulko, B., Danko, C., and Siepel, A. (2011). Bayesian inference of
423 ancient human demography from individual genome sequences. *Nat. Genet.*, 43:1031–
424 1034.
- 425 Günther, T. and Nettelblad, C. (2019). The presence and impact of reference bias on popu-
426 lation genomic studies of prehistoric human populations. *PLoS Genet.*, 15(7):e1008302.

- 427 Gutenkunst, R. N., Hernandez, R. D., Williamson, S. H., and Bustamante, C. D. (2009).
428 Inferring the joint demographic history of multiple populations from multidimensional
429 SNP frequency data. *PLoS Genet.*, 5(10):e1000695.
- 430 Hofreiter, M., Jaenicke, V., Serre, D., Haeseler, A. v., and Pääbo, S. (2001). DNA
431 sequences from multiple amplifications reveal artifacts induced by cytosine deamination
432 in ancient DNA. *Nucleic Acids Res.*, 29(23):4793–4799.
- 433 Kamm, J., Terhorst, J., Durbin, R., and Song, Y. S. (2019). Efficiently inferring the
434 demographic history of many populations with allele count data. *J. Am. Statist. Assoc.*,
435 pages 1–16.
- 436 Leppälä, K., Nielsen, S. V., and Mailund, T. (2017). admixturegraph: An R package for
437 admixture graph manipulation and fitting. *Bioinformatics*, 33(11):1738–1740.
- 438 Lipson, M., Loh, P.-R., Levin, A., Reich, D., Patterson, N., and Berger, B. (2013). Effi-
439 cient moment-based inference of admixture parameters and sources of gene flow. *Mol.*
440 *Biol. Evol.*, 30(8):1788–1802.
- 441 Lipson, M. and Reich, D. (2017). A working model of the deep relationships of diverse
442 modern human genetic lineages outside of Africa. *Mol. Biol. Evol.*, 34(4):889–902.
- 443 Lipson, M., Ribot, I., Mallick, S., Rohland, N., Olalde, I., Adamski, N., Broomandkhosh-
444 bacht, N., Lawson, A. M., López, S., Oppenheimer, J., et al. (2020). Ancient West
445 African foragers in the context of African population history. *Nature*, 577:665–670.
- 446 Lipson, M., Szécsényi-Nagy, A., Mallick, S., Pósa, A., Stégmár, B., Keerl, V., Rohland, N.,
447 Stewardson, K., Ferry, M., Michel, M., et al. (2017). Parallel palaeogenomic transects
448 reveal complex genetic history of early European farmers. *Nature*, 551(7680):368–372.
- 449 Mallick, S., Li, H., Lipson, M., Mathieson, I., Gymrek, M., Racimo, F., Zhao, M., Chen-
450 nagiri, N., Nordenfelt, S., Tandon, A., et al. (2016). The Simons Genome Diversity
451 Project: 300 genomes from 142 diverse populations. *Nature*, 538(7624):201–206.
- 452 Mathieson, I., Lazaridis, I., Rohland, N., Mallick, S., Patterson, N., Roodenberg, S. A.,
453 Harney, E., Stewardson, K., Fernandes, D., Novak, M., et al. (2015). Genome-wide
454 patterns of selection in 230 ancient Eurasians. *Nature*, 528(7583):499–503.
- 455 Nielsen, S. V. (2018). *Inferring gene flow between populations with statistical methods.*
456 PhD thesis, Aarhus University.
- 457 Patterson, N., Moorjani, P., Luo, Y., Mallick, S., Rohland, N., Zhan, Y., Genschoreck,
458 T., Webster, T., and Reich, D. (2012). Ancient admixture in human history. *Genetics*,
459 192(3):1065–1093.

- 460 Pease, J. B. and Hahn, M. W. (2015). Detection and polarization of introgression in a
461 five-taxon phylogeny. *Syst. Biol.*, 64(4):651–662.
- 462 Peter, B. M. (2016). Admixture, population structure, and F-statistics. *Genetics*,
463 202(4):1485–1501.
- 464 Pickrell, J. and Pritchard, J. (2012). Inference of population splits and mixtures from
465 genome-wide allele frequency data. *PLoS Genet.*, 8(11):e1002967.
- 466 Pickrell, J. K., Patterson, N., Loh, P.-R., Lipson, M., Berger, B., Stoneking, M., Pak-
467 endorf, B., and Reich, D. (2014). Ancient west Eurasian ancestry in southern and
468 eastern Africa. *Proc. Natl. Acad. Sci. U. S. A.*, 111(7):2632–2637.
- 469 Raghavan, M., Skoglund, P., Graf, K. E., Metspalu, M., Albrechtsen, A., Moltke, I.,
470 Rasmussen, S., Stafford Jr, T. W., Orlando, L., Metspalu, E., et al. (2014). Up-
471 per Palaeolithic Siberian genome reveals dual ancestry of Native Americans. *Nature*,
472 505(7481):87–91.
- 473 Reich, D., Thangaraj, K., Patterson, N., Price, A., and Singh, L. (2009). Reconstructing
474 Indian population history. *Nature*, 461(7263):489–494.
- 475 Shinde, V., Narasimhan, V. M., Rohland, N., Mallick, S., Mah, M., Lipson, M., Nakat-
476 suka, N., Adamski, N., Broomandkshobacht, N., Ferry, M., et al. (2019). An ancient
477 Harappan genome lacks ancestry from Steppe pastoralists or Iranian farmers. *Cell*,
478 179(3):729–735.
- 479 Soraggi, S. and Wiuf, C. (2019). General theory for stochastic admixture graphs and
480 F-statistics. *Theoret. Pop. Biol.*, 125:56–66.

481 **Data Accessibility**

482 The data that support the findings of this study are openly available through the European
483 Nucleotide Archive (ENA), under accession numbers PRJEB9586 and ERP010710, and at
484 the European Genome-phenome Archive (EGA), under accession number EGAS00001001959
485 (Mallick et al., 2016; Fan et al., 2019).