

Peer Review:

Estimating allele frequencies, ancestry proportions and genotype likelihoods in the presence of mapping bias

Mapping bias poses a significant challenge in the analysis of ancient DNA data. This study introduces testable hypotheses that address the impact of mapping bias on allele frequency estimates and admixture proportion estimation, particularly in ancient DNA research. By testing the effect of mapping bias, the study clearly demonstrates its influence on allele frequency estimation in empirical data. The corrected genotype likelihood approach shows the best correlation with “true” allele frequencies. The research further shows that while mapping bias can substantially affect ancestry proportion estimates, the adjusted genotype likelihoods can mitigate this issue. It also emphasizes the critical role of method selection, with some methods exhibiting considerable variability in results. These findings help refine methodologies in the field, making it possible to obtain more reliable results from low-coverage ancient DNA data and thus moving the field forward.

Global impression

The article makes a valuable contribution to the field by introducing a novel method for reducing mapping bias in ancient DNA analysis. It effectively outlines the problem and current challenges, with the proposed approach appearing both innovative and promising. The use of high-quality SNP array data adds value, as it provides a reliable control. Although it would have been interesting to see the effect of mapping bias on real data, the decision to simulate admixture data seems like a good choice to address this. However, while the impact of the corrected genotype likelihood on allele frequency and admixture estimation is significant, it looks very minor when compared to the standard genotype likelihood method. A more detailed biological interpretation of these results would be helpful to clarify why the modified genotype likelihood only has such a modest effect on mapping bias. With this in mind, a discussion of other potential sources of biases is still lacking.

Overall, the study provides valuable results to address the initial research question, but further investigation is needed to fully explain and contextualize these findings.

Major comment

Introduction

The introduction is well-constructed, providing a clear understanding of the challenges associated with mapping bias, the current strategies proposed to address these issues, and the new approach for mitigating mapping bias and assessing its impact. However, given that this project focuses on ancient DNA, it would be beneficial to dedicate more time to introducing ancient DNA and explaining the specific challenges involved in mapping this type of DNA. Additionally, the detailed description of algorithms for estimating admixture proportions is more appropriate for the methodology section, specifically under "2.4 Estimating Admixture Proportions."

Methodology

Regarding the methodology part, the four sections are relevant and well described. However, there are instances where the choice of certain parameters or values could benefit from more detailed justification or references such as bwa and ANGSD parameters. Additionally, I have significant concerns regarding the reproducibility of the simulations, as I encountered difficulties running your code for simulating genomic data, because some of the required packages and modules seem to be internal to the author's system without detailed information about their contents. To improve reproducibility, it is crucial to make these packages and modules available to the community and provide clear instructions on the specific commands and procedures used for the simulations.

Furthermore, I am concerned that selecting only SNPs with matching alleles in both pseudohaploid and SNP chip data might introduce a selection bias. This filtering approach excludes SNPs whose genotype is different due to methods, which might overlook important differences caused by different genotyping methods. Comparing allele frequencies between pseudohaploid data with all SNPs and pseudohaploid data filtered to match the SNP chip could reveal if the filtering process introduces significant biases and demonstrate that the SNP filtering does not significantly alter the results.

Results

The three sections are relevant, but the analyses seem shallow. For example, it would have been interesting to investigate whether certain genomic regions are more susceptible to mapping bias. Is mapping bias more frequent in GC-rich regions, repetitive sequences, or complex genomic areas? Visualizing the locations of these potential differences and correlating them with specific genomic features would provide deeper insights into the sources of mapping bias.

In general, the results lack proper biological interpretation and discussion, especially regarding the admixture simulations, on aspects such as LD pruning and the choice of reference genomes in function of each case. As it stands it is mainly descriptive. For example, it would be valuable to discuss and possibly investigate why allele frequencies from SNP arrays show lower correlation with those derived from pseudohaploid genotype calls, while admixture proportion estimation with qpAdm, which uses pseudohaploid genotype calls as input, appear to perform best.

Additionally, the paper would benefit from clearer conclusions at the end of each result section to highlight the really important information to take home.

Some figures are not readable, particularly those comparing simulated and estimated admixture proportions, as the use of white text on a gray background makes it difficult to read the details. Adding tables with the actual values of estimated admixture proportions would be helpful, as the small differences are hard to discern from the graphs. Including a table with correlations between allele frequencies values of SNP array and the different methods could also be valuable. Additionally, it would be useful to include a figure showing the distribution of read balance values (r) as supplementary material. This could help illustrate the types of mapping bias (reference or alternate) and the ratio between them.

Discussion

The limits of this study are discussed, but the authors should clarify some practical points such as whether it is better to use a reference genome that is closer or more distant genetically in order to compute allele frequencies or compute admixture proportions, since these results seem to be contradictory between computation of allele frequencies and admixture proportion inference.

The authors should consider adding some perspectives to this work, particularly in relation to the limitations of the study. While the issue of mapping bias has been reduced, it has not been completely resolved, as seen in the simulated data. Moreover, the impact on admixture proportion inference with read data remains unknown. As mentioned by the authors, mapping bias might have a greater effect on real datasets due to higher genomic variability, so the genotype likelihood correction could potentially reduce this impact more significantly. It would be valuable to evaluate the effect of the corrected genotype likelihood on non-simulated data.

Additional and minor comments

All other comments can be found in the PDF file as 'comment' elements.

Estimating allele frequencies, ancestry proportions and genotype likelihoods in the presence of mapping bias

Torsten Günther^{1,2,*} & Joshua G. Schraiber³

¹Human Evolution, Department of Organismal Biology, Uppsala University, Uppsala, Sweden

² Science for Life Laboratory, Ancient DNA Unit, Uppsala University, Uppsala, Sweden

³Department of Quantitative and Computational Biology, University of Southern California, Los Angeles, USA

*Corresponding author: torsten.gunther@ebc.uu.se

Abstract

Population genomic analyses rely on an accurate and unbiased characterization of the genetic setup of the studied population. For short-read, high-throughput sequencing data, mapping sequencing reads to a linear reference genome can bias population genetic inference due to mismatches in reads carrying non-reference alleles. In this study, we investigate the impact of mapping bias on allele frequency estimates from pseudohaploid data, commonly used in ultra-low coverage ancient DNA sequencing. To mitigate mapping bias, we propose an empirical adjustment to genotype likelihoods. Simulating ancient DNA data with realistic post-mortem damage, we compare widely used methods for estimating ancestry proportions under different scenarios, including reference genome selection, population divergence, and sequencing depth. Our findings reveal that mapping bias can lead to differences in estimated admixture proportion of up to 4% depending on the reference population. However, the choice of method has a much stronger impact, with some methods showing differences of 10%. qpAdm appears to perform best at estimating simulated ancestry proportions, but it is sensitive to mapping bias and its applicability may vary across species due to its requirement for additional populations beyond the sources and target population. Our adjusted genotype likelihood approach largely mitigates the effect of mapping bias on genome-wide ancestry estimates from genotype likelihood-based tools. However, it cannot account for the bias introduced by the method itself or the noise in individual site allele frequency estimates due to low sequencing depth. Overall, our study provides valuable insights for obtaining **precise** estimates of allele frequencies and ancestry proportions in empirical studies.

1 Introduction

1 A phenomenon gaining an increasing degree of attention in population genomics is mapping bias in
2 re-sequencing studies employing short sequencing reads (Orlando et al., 2013; Gopalakrishnan et al.,
3 2017; Günther and Nettelblad, 2019; Martiniano et al., 2020; Chen et al., 2021; Oliva et al., 2021;
4 Prasad et al., 2022; Gopalakrishnan et al., 2022; Thorburn et al., 2023; Koptekin et al., 2023). As
5 most mapping approaches employ linear reference genomes, reads carrying the same allele as the
6 reference will have fewer mismatches and higher mapping scores than reads carrying an alternative

7 allele leading to some alternative reads being rejected. As a consequence, sequenced individuals may
8 seem more similar to the reference genome (and hence, the individual/population/species it originates
9 from) than it is in reality, biasing variant calling and downstream analysis. **The effect of mapping bias**
10 **is exacerbated in ancient DNA studies due to post-mortem DNA damage such as fragmentation and**
11 **cytosine deamination to uracil (which is sequenced as thymine) (Orlando et al., 2021).** The human
12 reference genome is a mosaic sequence of multiple individuals from different continental ancestries
13 (Green et al., 2010; Church et al., 2015). In most other species with an existing reference genome
14 sequence, this genome represents a single individual from a certain population while for studies in
15 species without a reference genome, researchers are limited to the genomes of related species. One
16 consequence is that the sequence at a locus in the reference genome may either represent an ingroup
17 or an outgroup relative to the other sequences studied in a population genomic analysis. It has been
18 shown that this can bias estimates of heterozygosity, phylogenetic placement, assessment of gene flow,
19 and population affinity (see e.g. Orlando et al., 2013; Heintzman et al., 2017; Gopalakrishnan et al.,
20 2017; Günther and Nettelblad, 2019; van der Valk et al., 2020; Mathieson et al., 2020; Prasad et al.,
21 2022). Notably, while mapping bias mostly manifests as reference bias, it also exists as alternative
22 bias depending on the studied individual and the particular position in the genome (Günther and
23 Nettelblad, 2019).

24 Different strategies have been proposed to mitigate or remove the effect of mapping bias. These
25 include mapping to an outgroup species (Orlando et al., 2013), mapping to multiple genomes simul-
26 taneously (Huang et al., 2013; Chen et al., 2021), mapping to variation graphs (Martiniano et al.,
27 2020), the use of an IUPAC reference genome (Oliva et al., 2021), masking variable sites (Koptekin
28 et al., 2023) or filtering of “biased reads” (Günther and Nettelblad, 2019). All of these strategies
29 have significant limitations, such as exclusion of some precious sequencing reads (outgroup mapping
30 or filtering) or requiring additional data that may not be available for all species prior to the particular
31 study (variation graphs, IUPAC reference genomes, or mapping to multiple genomes). Therefore, it
32 would be preferable to develop a strategy that uses the available sequencing reads and accounts for
33 potential biases in downstream analyses. Genotype likelihoods (Nielsen et al., 2011) represent one
34 promising approach that can be used with low- and medium-depth sequencing data (Lou et al., 2021).
35 Instead of working with hard genotype calls at each position one can use $P(D|G)$, the probability
36 of observing a set of sequencing reads D conditional on a true genotype G . Different approaches
37 exist for calculating genotype likelihoods with the main aim to account for uncertainty due to random
38 sampling of sequencing reads and sequencing error. Genotype likelihoods can be used in a wide range
39 of potential applications for downstream analysis which include imputation (Rubinacci et al., 2021),
40 estimation of admixture proportions (Skotte et al., 2013; Jørsboe et al., 2017; Meisner and Albrecht-
41 sen, 2018), principal component analysis (PCA, Meisner and Albrechtsen, 2018), relatedness analysis
42 (Korneliussen and Moltke, 2015; Hanghøj et al., 2019; Nøhr et al., 2021), or to search for signals of
43 selection (Korneliussen et al., 2013; Fumagalli et al., 2013). Many of these are available as part of the
44 popular software package ANGSD (Korneliussen et al., 2014). **However, some downstream results can**
45 **depend on the specific genotype likelihood model selected (Lou et al., 2021).**

46 To render genotype likelihoods and their downstream applications more robust to the presence of
47 mapping bias, we introduce a modified genotype likelihood, building off of the approach in Günther
48 and Nettelblad (2019). We use modified reads carrying the other allele at biallelic SNP positions to
49 assess the distribution of mapping bias and to obtain an empirical quantification of the locus- and
50 individual-specific mapping bias. We then calculate a modified genotype likelihood to account for
51 mapping bias. The approach is similar to `snpAD` (Prüfer, 2018), with the contrast that our aim is not
52 to call **genotypes all sites** and we are using a set of ascertained biallelic SNPs allowing us to estimate
53 mapping bias locus-specific rather than using one estimate across the full genome of the particular
54 individual.

55 We examine two downstream applications of genetic data to determine the impact of mapping bias,
56 and assess the ability of our corrected genotype likelihood to ameliorate issues with mapping bias.
57 First, we look at a very high-level summary of genetic variation: allele frequencies. Because allele

58 frequencies can be estimated from high-quality SNP array data, we can use them as a control and
59 assess the impact of mapping bias and our corrected genotype likelihood in real short-read data.

60 Next, we examine the assignment of ancestry proportions. Most currently used methods trace
61 their roots back to the software STRUCTURE (Pritchard et al., 2000; Falush et al., 2003, 2007; Hubisz
62 et al., 2009), a model-based clustering approach modeling each individual’s ancestry from K source
63 populations (PSD model). These source populations can be inferred from multi-individual data (unsu-
64 pervised) or groups of individuals can be designated as sources (supervised). Popular implementations
65 of this model differ in terms of input data (e.g. genotype calls or genotype likelihoods), optimization
66 procedure and whether they implement a supervised and/or unsupervised approach (Table 1). In
67 the ancient DNA field, f statistics (Patterson et al., 2012) and their derivatives are fundamental to
68 many studies due to their versatility, efficiency and their ability to work with pseudohaploid data.
69 Consequently, methods based on f statistics are also often used for estimating ancestry proportions in
70 ancient DNA studies. One method that uses f statistics for supervised estimation of ancestry propor-
71 tions is qpAdm (Haak et al., 2015; Harney et al., 2021). In addition to the source populations (“left”
72 populations), a set of more distantly related “right” populations is needed for this approach. Ancestry
73 proportions are then estimated from a set of f_4 statistics calculated between the target population
74 and the “left” and “right” populations. We simulate data sequencing data with realistic ancient DNA
75 damage under a demographic model with recent gene flow (Figure 1) and then compare the different
76 methods in their ability to recover the estimated admixture proportion and how sensitive they are to
77 mapping bias.

78 2 Materials and Methods

79 2.1 Correcting genotype-likelihoods for mapping bias

80 Two versions of genotype likelihoods (Nielsen et al., 2011) were calculated for this study. First, we
81 use the direct method as included in the original version of GATK (McKenna et al., 2010) and also
82 implemented in ANGSD (Korneliussen et al., 2014). For a position ℓ covered by n reads, the genotype
83 likelihood is defined as the probability for observing the bases $D_\ell = \{b_{\ell_1}, b_{\ell_2}, \dots, b_{\ell_n}\}$ if the true
84 genotype is A_1A_2 :

$$P(D_\ell | G_\ell = A_1, A_2) = \prod_{i=1}^n P(b_{\ell_i} | G_\ell = A_1, A_2) = \prod_{i=1}^n \frac{P(b_{\ell_i} | A_1) + P(b_{\ell_i} | A_2)}{2} \quad (1)$$

85 with

$$P(b_{\ell_i} | A) = \begin{cases} \frac{e_{\ell_i}}{3} & \text{if } b = A \\ 1 - e_{\ell_i} & \text{if } b \neq A \end{cases}$$

86 where e_{ℓ_i} is the probability of a sequencing error of read i at position ℓ , calculated from the phred scaled
87 base quality score Q_{ℓ_i} , i.e. $e_{\ell_i} = 10^{-Q_{\ell_i}/10}$. The calculation of genotype likelihoods was implemented
88 in Python 3 using the pysam library (<https://github.com/pysam-developers/pysam>), a wrapper
89 around htslib and the samtools package (Li et al., 2009) or by calling samtools mpileup and parsing
90 the output in the Python script.

91 To quantify the impact of mapping bias, we restrict the following analysis to ascertained biallelic
92 SNPs and modify each original read to carry the other allele at the SNP position, as in Günther
93 and Nettelblad (2019). The modified reads are then remapped to the reference genome using the
94 same mapping parameters. If there were no mapping bias, all modified reads would map to the same
95 position as the unmodified original read. Consequently, when counting both original and modified
96 reads together, we should observe half of our reads carrying the reference allele and the other half
97 carrying the alternative allele at the SNP position. We can summarize the read balance at position ℓ as
98 r_ℓ , which measures the proportion of reference alleles among all original and modified reads mapping

99 to the position. Without mapping bias, we would observe $r_\ell = 0.5$. Under reference bias, we would
100 observe $r_\ell > 0.5$ and under alternative bias $r_\ell < 0.5$. We can see r_ℓ as an empirical quantification
101 of the locus- and individual-specific mapping bias. Similar to [Prüfer \(2018\)](#), we can then modify
102 equation 1 for heterozygous sites to

$$P(D_\ell | G_\ell = R_\ell, A_\ell) = \prod_{i=1}^n r_\ell P(b_{\ell i} | R_\ell) + (1 - r_\ell) P(b_{\ell i} | A_\ell) \quad (2)$$

103 where R_ℓ is the reference allele at position ℓ and A_ℓ is the alternative allele. Genotype likelihood-
104 based methods are tested with both genotype likelihood versions. All code used in this study can be
105 found under https://github.com/tgue/refbias_GL

106 2.2 Empirical Data

107 To estimate the effect of mapping bias in empirical data we obtained low coverage BAM files for **ten**
108 **FIN individuals and 10 YRI individuals from the 1000 Genomes project (Table S1)** ([Auton et al.,](#)
109 [2015](#)). We also downloaded Illumina Omni2.5M chip genotype calls for the same individuals. The
110 SNP data was filtered to restrict to sites without missing data in the 20 selected individuals, a minor
111 **allele frequency of at least 0.2** in the reduced dataset (considering individuals from both populations
112 together), and excluding A/T and C/G SNPs to avoid strand misidentification. Reads mapping
113 to these positions were extracted from the BAM files using **samtools** ([Li et al., 2009](#)). **To make the**
114 **sequence data more similar to fragmented ancient DNA, each read was split into two halves at its mid-**
115 **point and each sub-read was re-mapped separately.** For mapping, we used **bwa aln** ([Li and Durbin,](#)
116 [2009](#)) and the non-default parameters **-l 16500** (to avoid seeding), **-n 0.01** and **-o 2**. Only reads with
117 mapping qualities of 30 or higher were kept for further analysis. **Pseudohaploid genotypes were called**
118 **with ANGSD v0.933** ([Korneliussen et al., 2014](#)) **by randomly drawing one read per SNP as described for**
119 **the simulations below and only SNPs with the same two alleles in pseudohaploid and SNP chip data**
120 **were included in all comparisons.** **Remapping of modified reads and genotype likelihood calculation**
121 **were performed as described above. Allele frequencies were calculated from genotype likelihoods with**
122 **ANGSD v0.933** ([Korneliussen et al., 2014](#)) **using -doMaf 4 and the human reference as “ancestral” allele**
123 **in order to calculate the allele frequency of the reference alleles. SNP calls from the genotyping array**
124 **and pseudohaploid calls were converted to genotype likelihood files assuming no genotyping errors, so**
125 **the allele frequency estimation for this data could be based on ANGSD as well.**

126 2.3 Simulation of genomic data

127 **Population histories are simulated** using **msprime v0.6.2** ([Kelleher et al., 2016](#)). We simulate a demo-
128 graphic history where a target population T receives a single pulse of admixture with proportion f
129 from source $S3$ 50 generations ago. Furthermore, we simulate population $S1$ which forms an outgroup
130 and population $S2$ which is closer to T than $S3$ to serve as second source for estimating ancestry pro-
131 portions (Figure 1). Finally, we simulate populations $O1$, $O2$, $O3$, and $O4$ as populations not involved
132 in the admixture events which split off internal branches of the tree to serve as “right” populations
133 for **qpAdm** ([Haak et al., 2015](#); [Harney et al., 2021](#)). Split times are scaled relative to the deepest split
134 t_{123} : the split between $(S2, T)$ and $S3$, t_{23} , is set to $0.5 \times t_{123}$ while the split between T and $S2$ is set
135 to $0.2 \times t_{123}$. **Different values of 20,000 and 50,000 generations are tested for t_{123} approximately corre-**
136 **sponding to divergence times within and between (sub-)species.** **Mutation rate was set to 2.5×10^{-8}**
137 **and recombination rate was set to 2×10^{-8} .** **The effective population size along all branches is 10,000.**
138 **For each population, 21 diploid individuals (i.e. 42 haploid chromosomes) with 5 chromosome pairs**
139 **of 20,000,000 bp each were simulated.**

140 For each chromosome, a random ancestral sequence was generated with a GC content of 41% corre-
141 sponding to the GC content of the human genome ([Lander et al., 2001](#)). Transversion polymorphisms
142 were then placed along the sequence according to the **msprime** simulations. The first sequences from
143 populations $S1$, $S2$ and $S3$ were used as reference genomes. Pairs of sequences were then considered as

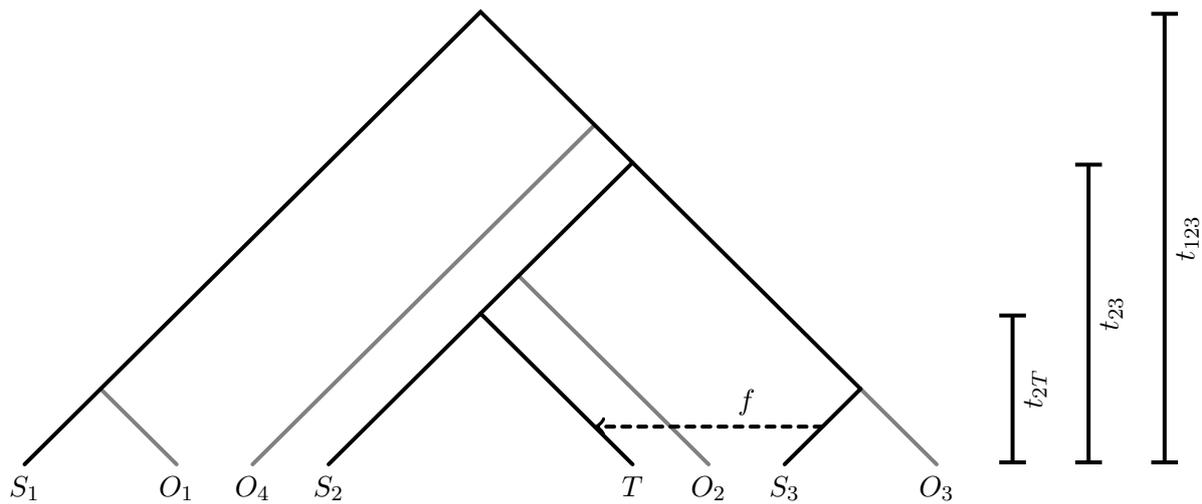


Figure 1: Illustration of the population relationships used in the simulations. Branch lengths are not to scale

144 diploid individuals and `gargamel` (Renaud et al., 2017) was used to simulate next-generation sequencing
 145 ing data with ancient DNA damage. Data were simulated to mimic data generated with an Illumina
 146 HiSeq 2500 sequencing machine assuming the post-mortem damage pattern observed when sequencing
 147 Neandertals in Briggs et al. (2007). For each individual, fragment sizes followed a log-normal distribu-
 148 tion with a location between 3.3 and 3.8 (randomly drawn per individual from a uniform distribution)
 149 and a scale of 0.2, corresponding to an average fragment length per individual between 27 and 46bp.
 150 Fragments shorter than 20bp were excluded. No contaminating sequences were simulated. Sequencing
 151 reads were then trimmed and merged with `AdapterRemoval` (Schubert et al., 2016). Reads were then
 152 mapped to the different reference genomes using `bwa aln v0.7.17` (Li and Durbin, 2009) together with
 153 the commonly used non-default parameters `-l 16500` (to avoid seeding), `-n 0.01` and `-o 2` (Schubert
 154 et al., 2012; Oliva et al., 2021). BAM files were handled using `samtools v1.5` (Li et al., 2009).

155 Genotype calling and downstream analysis were performed separately for the three reference genomes
 156 originating from populations $S1$, $S2$ and $S3$. To avoid ascertainment bias, polymorphic SNPs were as-
 157 certained from the simulated true genotypes and restricted to SNPs with a minimum allele frequency
 158 of 10% in the outgroup population $S1$. 100,000 SNPs were selected at random using `Plink v1.90`
 159 (Chang et al., 2015) `-thin-count`. Pseudohaploid calls were then generated for all individuals at these
 160 sites using `ANGSD v0.917` (Korneliussen et al., 2014) by randomly sampling a single read per position
 161 with minimum base and mapping quality of at least 30. This step was performed using `ANGSD` with
 162 the parameters `-checkBamHeaders 0 -doHaploCall 1 -doCounts 1 -doGeno -4 -doPost 2 -doPlink 2`
 163 `-minMapQ 30 -minQ 30 -doMajorMinor 1 -GL 1 -domaf 1`. Files were then converted to `Plink` format
 164 using `haploToPlink` distributed with `ANGSD` (Korneliussen et al., 2014). For downstream analyses,
 165 the set of SNPs was further restricted to sites with less than 50 % missing data and a minor allele
 166 frequency of at least 10% in $S1$, $S2$, $S3$ and T together. Binary and transposed `Plink` files were
 167 handled using `Plink v1.90` (Chang et al., 2015). `convertf` (Patterson et al., 2006; Price et al., 2006)
 168 was used to convert between `Plink` and `EIGENSTRAT` file formats. `Plink` was also used for linkage
 169 disequilibrium (LD) pruning with parameters `-indep-pairwise 200 25 0.7`.

2.4 Estimating admixture proportions

171 We used five different approaches to estimate ancestry proportions in our target population T . In
 172 addition to differences in the underlying model and implementations, for users the tools differ in the
 173 type of their input data (genotype calls or genotype likelihoods) and whether their approaches are
 174 unsupervised and/or supervised (Table 1).

175 All software was set to estimate ancestry assuming two source populations. Unless stated otherwise,

Table 1: Overview of the different tools used for ancestry estimation.

Method	Genotype calls	Genotype-likelihoods	Unsupervised	Supervised	Citation
ADMIXTURE	X	-	X	X	Alexander et al. (2009); Alexander and Lange (2011)
qpAdm	X	-	-	X	Haak et al. (2015); Harney et al. (2021)
NGSadmix	-	X	X	-	Skotte et al. (2013)
fastNGSadmix	-*	X	-	X	Jørsboe et al. (2017)

* source populations for fastNGSadmix can be either genotype calls or genotype likelihoods

176 $S2$ and $S3$ were set as sources and T as the target population while no other individuals were included
 177 in when running the software. ADMIXTURE (Alexander et al., 2009; Alexander and Lange, 2011) is the
 178 only included method that has both a supervised (i.e. with pre-defined source populations) and an
 179 unsupervised mode. Both options were tested using the `-haploid` option without multithreading as the
 180 genotype calls were pseudo-haploid. For qpAdm (Haak et al., 2015; Harney et al., 2021), populations
 181 $O1$, $O2$, $O3$ and $O4$ served as “right” populations. qpAdm was run with the options `allsnps: YES` and
 182 details: YES. For fastNGSadmix (Jørsboe et al., 2017), allele frequencies in the source populations
 183 were estimated using NGSadmixmap (Skotte et al., 2013) with the option `-printInfo 1`. fastNGSadmix
 184 was then run to estimate ancestry per individual without bootstrapping. NGSadmixmap (Skotte et al.,
 185 2013) was run in default setting. The mean ancestry proportions across all individuals in the target
 186 population was used as an ancestry estimate for the entire population. In the case of unsupervised
 187 approaches, the clusters belonging to the source populations were identified as those where individuals
 188 from $S2$ or $S3$ showed more than 90 % estimated ancestry.

189 3 Results

190 3.1 Mapping bias in empirical data

191 We first tested the effect of mapping bias on allele frequency estimates in empirical data. We selected
 192 low to medium coverage (mostly between 2 and 4X depth except for one individual at 14X, Table S1)
 193 for ten individuals from each of two 1000 Genomes populations (FIN and YRI). We used ANGSD to
 194 estimate allele frequencies and compare them to allele frequencies estimated from the same individuals
 195 genotyped using a SNP array and pseudohaploid genotype data. As the genotyping array should be
 196 less affected by mapping bias, we consider these estimates as “true” allele frequencies.

197 Overall, genotype likelihood-based point estimates of the allele frequencies tend towards more inter-
 198 mediate allele frequencies while pseudohaploid genotypes and “true” genotypes result in more alleles
 199 estimated to have low and high alternative allele frequency (Figure 2A and B). In FIN, the pseu-
 200 dohaploid genotypes lead to a slight underestimation of the reference allele frequencies (Figure 2A),
 201 while this signal is reversed in YRI (Figure 2B), a pattern which could be related to the fact that
 202 most of the human reference genome has European ancestry (Green et al., 2010; Church et al., 2015;
 203 Günther and Nettelblad, 2019). In both tested populations, the default version of genotype likelihood
 204 calculation produced an allele frequency distribution slightly shifted towards lower non-reference allele
 205 frequency estimates (Paired Wilcoxon test $p < 2.2 \times 10^{-22}$ in both populations). The allele frequen-
 206 cies estimated from the corrected genotype likelihoods exhibit a slightly better correlation with the
 207 “true” frequencies in both FIN (Pearson’s correlation coefficient 0.9297 [0.9294, 0.9301] vs. 0.9310
 208 [0.9307, 0.9313] for uncorrected and corrected, respectively; $p = 2.14 \times 10^{-7}$) and YRI (Pearson’s cor-
 209 relation coefficient 0.9444 [0.9442, 0.9447] vs. 0.9459 [0.9457, 0.9462] for uncorrected and corrected,
 210 respectively; $p = 1.8 \times 10^{-14}$). Notably, allele frequency estimates from pseudohaploid data display
 211 the lowest correlation with the “true” frequencies in both FIN ($r = 0.8571$) and YRI ($r = 0.8344$)
 212 indicating that while the distribution of allele frequencies seems close to the true spectrum (Figure

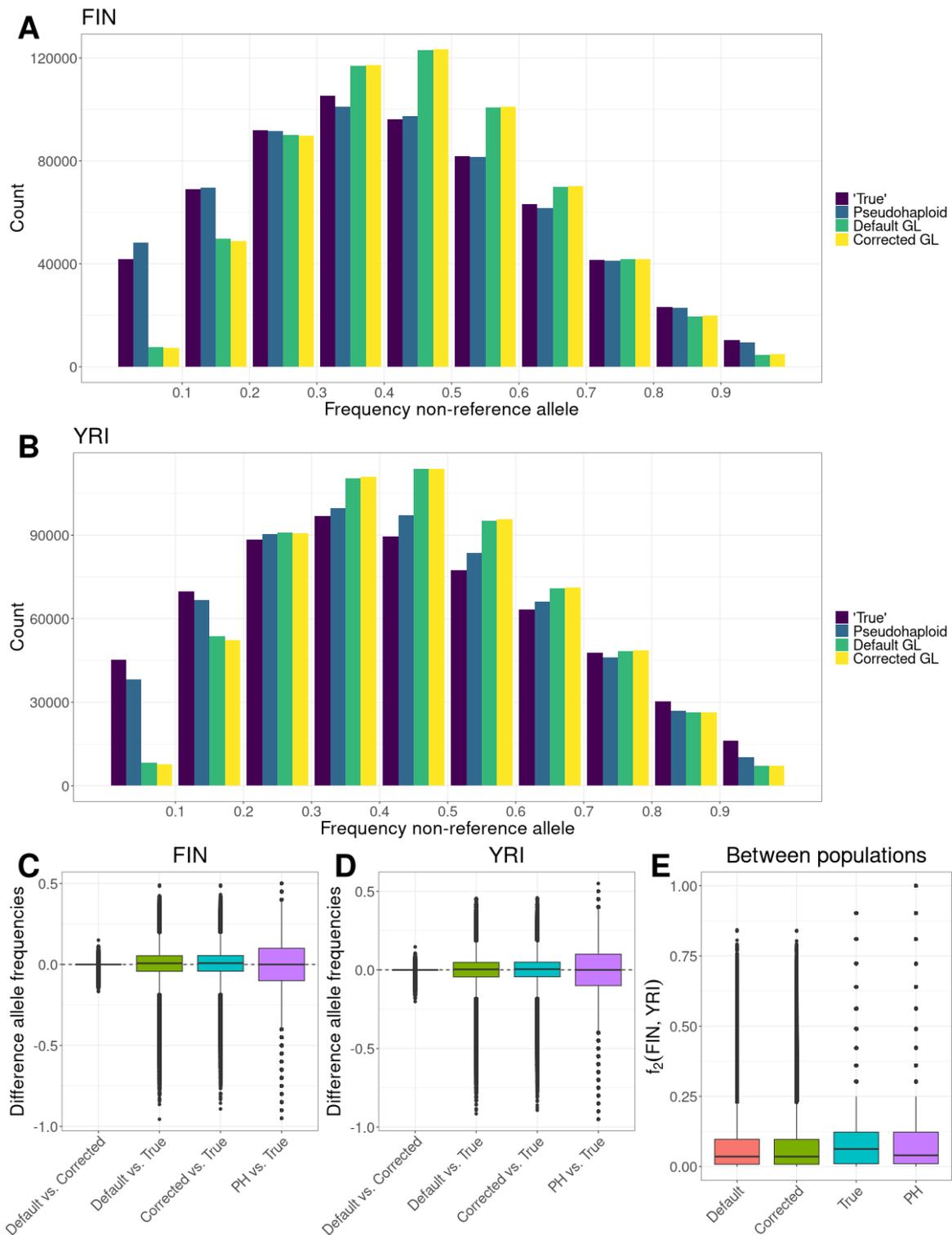


Figure 2: Differences in allele frequency estimates. Binned spectrum of non-reference alleles in FIN (A) and YRI (B) for the four different estimation methods. Note that the specific ascertainment of common SNPs in the joint genotyping data contributes to the enrichment of variants with intermediate frequencies. Boxplots for the differences between default genotype likelihood-based estimates and corrected genotype likelihood-based estimates, default genotype likelihood-based estimates and SNP array-based estimates, corrected genotype likelihood-based estimates, pseudohaploid (PH) genotype-based and SNP array-based estimates (C) in the FIN population and (D) in the YRI population. (E) is showing boxplots of the per-site population differentiation (measured as f_2 statistic) for the four allele frequency estimates.

213 2A and B), the estimates at individual loci are rather noisy.

214 Differences at individual sites, however, display some extreme outliers with $\sim 0.1\%$ of the SNPs
215 showing more than 50% difference between estimates from SNP chips and sequencing data, which could
216 hint at systematic technological differences between the two data types at those sites. This pattern of
217 outliers is slightly less pronounced when using the corrected genotype likelihoods (Table S2). Interest-
218 ingly, despite the overall closer concordance between the pseudohaploid allele frequency spectrum and
219 the SNP array allele frequency spectrum, there is significantly higher variation between pseudohaploid
220 and true frequencies at any particular hint, suggesting that this is a general difference between NGS
221 and SNP chip data. In Günther and Nettelblad (2019), we found that different parts of the human
222 reference genome exhibit different types of mapping bias. We find a similar result here: the parts of
223 the reference genome that can be attributed to African ancestry (Green et al., 2010) display a mean
224 and median difference of nearly 0 in FIN but allele frequencies remain higher than array estimates
225 in YRI (Figure S1). In contrast, the European and East Asian parts of the reference genome show a
226 distribution of differences around 0 in YRI but positive means and median in FIN (Figures S2 and
227 S3). This confirms the utility of reducing the effect of mapping bias by mapping against a reference
228 genome from an outgroup. A consequence of the systematic over-estimation of the allele frequencies
229 when using genotype likelihoods is that the population differentiation (here measured as f_2 statistic)
230 is reduced compared to estimates from the SNP array or pseudohaploid genotype calls (Figure 2E).

231 3.2 Estimation of admixture proportions based on genotype calls

232 We compare the accuracy of the different methods for estimating admixture proportion under a set
233 of different population divergence times, sequencing depths, and with or without LD pruning of the
234 SNP panel. For most parts of this results section, we will focus on the scenario with an average
235 sequencing depth of 0.5X where the deepest population split (t_{123}) was 50,000 generations ago and
236 the split (t_{23}) between the relevant sources dating to 25,000 generations ago. Consequently, mapping
237 the reads against a reference genome sequence from one or the other source would be equivalent to a
238 study comparing (sub-)species where the reference genome originated from one of those populations.
239 Results for other population divergences and sequencing depths are shown in Figures S4-S9.

240 We begin by assessing methods that require hard genotype calls, ADMIXTURE and qpAdm. For these
241 approaches, we used single randomly drawn reads per individual and site to generate pseudo-haploid
242 data in the target population. The popular implementation of the PSD (Pritchard et al., 2000) model
243 working with SNP genotype calls, ADMIXTURE (Alexander et al., 2009; Alexander and Lange, 2011),
244 has both supervised and unsupervised modes. Both modes show similar general patterns: low (10%)
245 admixture proportions are estimated well while medium to high ($\geq 50\%$) admixture proportions are
246 over-estimated (Figure 3). On the full SNP panel, the median estimated admixture proportion differs
247 up to $\sim 4\%$ when mapping to reference genomes representing either of the two sources (S2 or S3)
248 while mapping to the outgroup reference genome (S1) results in estimates intermediate between the
249 two. LD pruning slightly reduces mapping bias and reduces the overestimation, at least for high (90%)
250 admixture proportions. qpAdm (Haak et al., 2015; Harney et al., 2021), on the other hand, estimated all
251 admixture proportions accurately when the outgroup (S1) was used for the reference genome sequence
252 and when the full SNP panel was used. The median estimates of admixture differed up to 3% between
253 mapping to reference genomes from one of the source populations (S2 or S3). Notably, LD pruning
254 increased the noise of the qpAdm estimates (probably due to the reduced number of SNPs) and led
255 to all admixture proportions being slightly underestimated (Figure 3). The extent of mapping bias
256 decreases with lower population divergence across all methods (Figure S4), as mapping bias should
257 correlate with distance to the reference genome sequence. Conversely, increasing sequencing depth
258 mostly reduced noise but not mapping bias (Figures S5 and S8) as the genotype-based methods
259 benefit from the increased number of SNPs but the genotype calls do not increase certainty when
260 multiple reads are mapping to the same position.

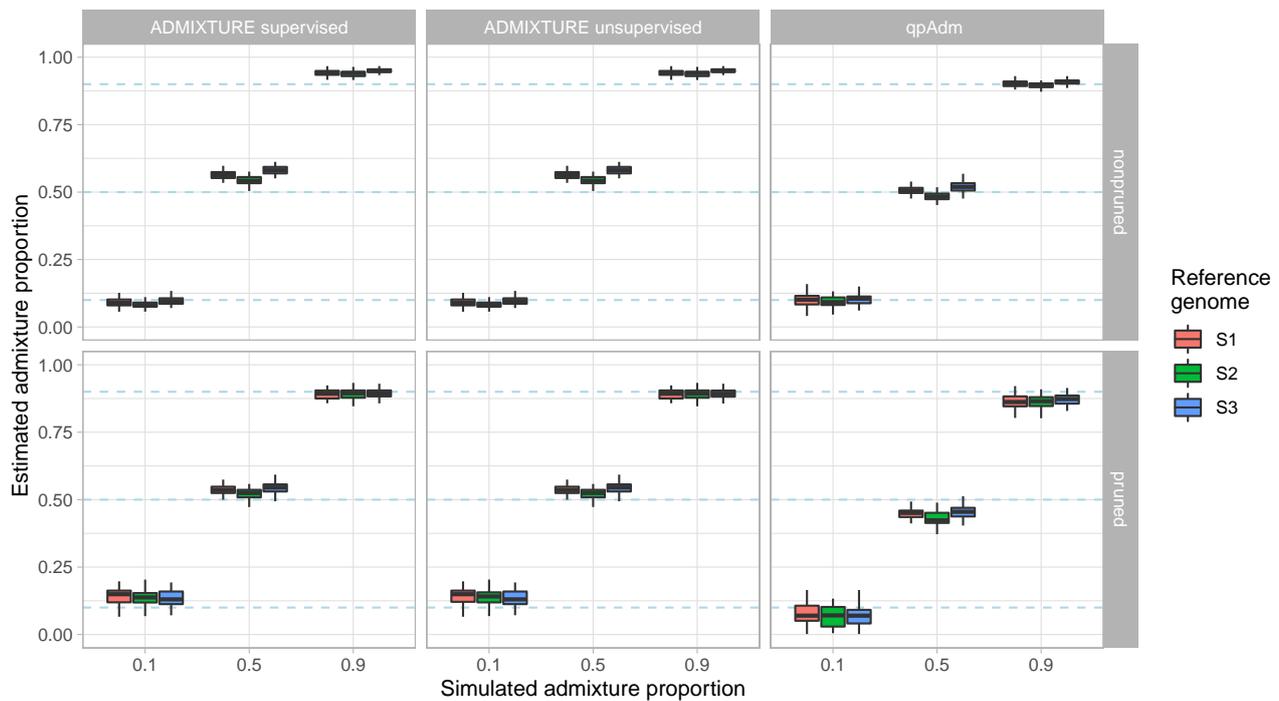


Figure 3: Simulation results for genotype call based methods using $t_{123} = 50000$ generations and a sequencing depth of 0.5X. Dashed blue lines represent the simulated admixture proportions.

3.3 Estimation of admixture proportions based on genotype likelihoods

261

262 We next examined the performance of genotype-likelihood-based approaches to estimate admixture
 263 proportions. In principle, genotype likelihoods should be able to make better use of all of the data in
 264 ancient DNA, because more than a single random read can be used per site. Moreover, we are able
 265 to explicitly incorporate our mapping bias correction into the genotype likelihood. We compared the
 266 supervised *fastNGSadmix* (Jørsboe et al., 2017) to the unsupervised *NGSadmix* (Skotte et al., 2013).
 267 *fastNGSadmix* shows the highest level of overestimation of low to medium admixture proportions
 268 ($\leq 50\%$) among all tested approaches while high admixture proportions (90%) are estimated well
 269 (Figure 4). Mapping bias caused differences of up to $\sim 3\%$ in the admixture estimates when mapping to
 270 the different reference genomes. LD pruning enhances the overestimation of low admixture proportions
 271 while leading to an underestimation of high admixture proportions. Notably, when employing the
 272 corrected genotype-likelihood the estimated admixture proportions when mapping to *S2* or *S3* are
 273 slightly more similar than with the default formula without correction, showing that the correction
 274 makes the genome-wide estimates less dependent on the reference sequence used for mapping while
 275 not fully removing the effect. The estimates when using the outgroup *S1* as reference are slightly
 276 higher for high admixture proportions (90%). The results for *NGSadmix* show similar patterns to
 277 *ADMIXTURE* with a moderate overestimation of admixture proportions $\geq 50\%$ (Figure 4). Mapping
 278 bias caused differences of up to $\sim 4\%$ in the admixture estimates when mapping to the different
 279 reference genomes. After LD pruning, estimated admixture proportions for higher simulated values
 280 were closer to the simulated values. Furthermore, employing the mapping bias corrected genotype-
 281 likelihoods made the estimated admixture proportions less dependent on the reference genome used
 282 during mapping. Notably, the extent of over-estimation for both methods seems to be somewhat
 283 negatively correlated with population divergence (Figures S6 and 4), i.e. increased distances between
 284 the source populations reduces the method bias. Further patterns are as expected: the extent of
 285 mapping bias is correlated with population divergence and increased sequencing depth reduces noise
 286 (Figures S6, 4, S7 and S9).

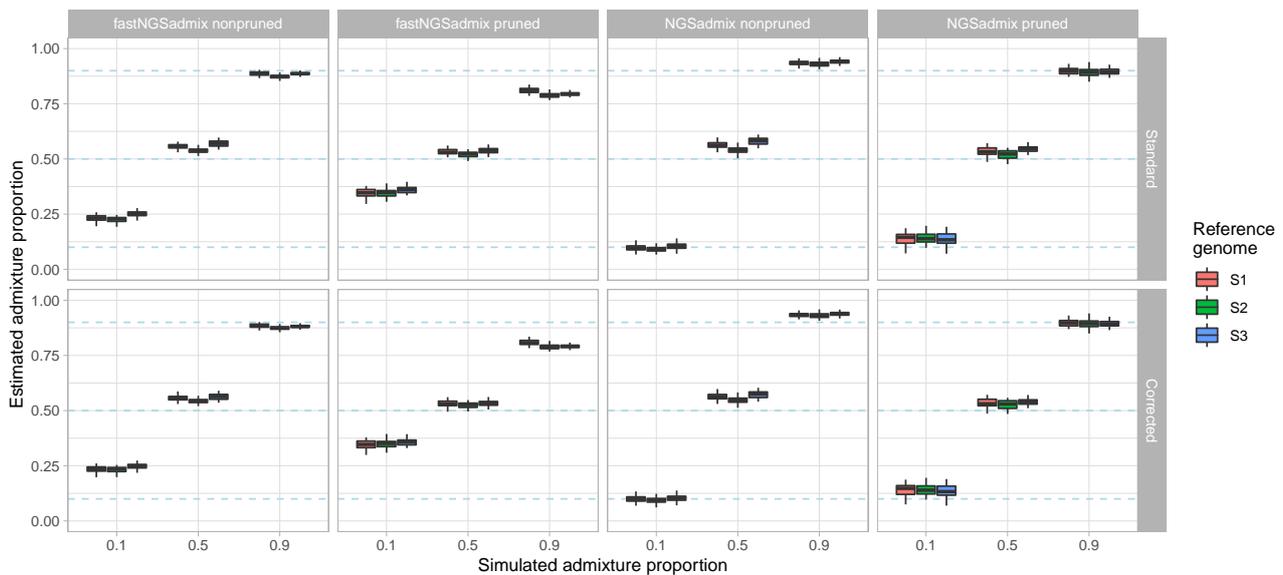


Figure 4: Simulation results for genotype likelihood based methods using $t_{123} = 50000$ generations and a sequencing depth of 0.5X. Dashed blue lines represent the simulated admixture proportions.

287

4 Discussion

288 We illustrate the impacts of mapping bias on downstream applications, such as allele frequency esti-
289 mation and ancestry proportion estimation, and we introduced a new approach to recalibrate genotype
290 likelihoods in the presence of mapping bias to alleviate its effects. The impact of mapping bias in
291 our comparisons is small but pervasive suggesting that it can have an effect on the results of different
292 types of analysis in empirical studies.

293 Increasing sample sizes in ancient DNA studies have motivated a number of studies aiming to detect
294 selection in genome-wide scans or to investigate phenotypes in ancient populations (e.g. Mathieson
295 et al., 2015; Cox et al., 2022; Klunk et al., 2022; Gopalakrishnan et al., 2022; Mathieson and Terhorst,
296 2022; Davy et al., 2023; Barton et al., 2023; Hui et al., 2024). Such investigations are potentially very
297 sensitive to biases and uncertainties in genotype calls or allele frequencies at individual sites while
298 certain effects will average out for genome-wide estimates such as ancestry proportions. Concerns
299 about certain biases and how to estimate allele frequencies have even reduced confidence in the results
300 of some studies (Gopalakrishnan et al., 2022; Barton et al., 2023). Our results indicate that such con-
301 cerns are valid as individual sites can show very strong deviations in their allele frequencies estimated
302 from low-coverage sequencing data. This is due to a combination of effects, including mapping bias
303 and **sampling artifacts**. Allele frequency point estimates from genotype likelihoods tend to be higher
304 than true frequencies because most alleles segregate at low frequencies, and thus appear most com-
305 monly in heterozygotes. However, genotype likelihood approaches without an allele frequency prior
306 will naturally put some weight on individuals being homozygous for the allele, ultimately collectively
307 driving up allele frequency estimates. The risk is then that most downstream analyses will treat the
308 allele frequency point estimates as face values potentially leading to both false positives and negatives.
309 While our new approach to recalibrate genotype likelihoods reduces the number of outlier loci, there
310 is still uncertainty in allele frequency estimates from low coverage data. Therefore, results heavily
311 relying on allele frequency estimates or genotype calls at single loci from low-coverage sequencing data
312 or even ancient DNA data need to be taken with a grain of salt.

313 The simulations in this study revealed a modest but significant effect of mapping bias on ancestry
314 estimates as the difference between reference genomes never exceeded 5 percent. **The differences seen**
315 **in our simulations are likely underestimates of what might occur in empirical studies as real genomes**
316 **are larger and more complex than what was used in the simulations.** For instance, we simulated five

317 20 megabase long chromosomes for a 100 megabase genome, while mammalian genomes are one order
318 of magnitude larger; the human genome is roughly 3 gigabases and the shortest human chromosome
319 alone is ~ 45 megabases long. Furthermore, the only added complexity when generating the random
320 sequences was a GC content of 41%. Real genomes also experience more complex mutation events
321 involving translocations and duplications, which, together with the increased length and the presence
322 of repetitive elements, should increase mapping bias in empirical studies. Finally, the range of possible
323 demographic histories including the relationships of targets and sources, drift as well as the timing
324 and number of gene flow is impossible to explore in a simulation study. The restricted scenarios tested
325 in this study should affect the quantitative results but the qualitative interpretation of mapping bias
326 impacting ancestry estimates should extend beyond the specific model used in the simulations.

327 While the ancestry estimates depended slightly on the reference genome the reads were mapped to,
328 they seemed more influenced by the choice of method or software. Methods easily differed by more
329 than 10% in their ancestry estimates from the same source data. This highlights that **other factors**
330 **and biases play major roles in the performance of these methods**. Depending on the method, the type
331 of input data and the implementation, they showed different sensitivities to e.g. the amount of missing
332 data or linkage. For non-pruned data, **qpAdm** performed best across all scenarios and did not show
333 any method-specific bias in certain ranges of simulated admixture proportions. This supports the
334 common practice of using **qpAdm** in most human ancient DNA studies. However, the requirement of
335 data from additional, “right” populations, might not make it applicable to many non-human species.
336 Furthermore, **qpAdm** only works with genotype calls, so it is influenced by mapping bias in similar
337 ways as **ADMIXTURE** and these methods cannot benefit from the newly introduced genotype likelihood
338 estimation. We also need to note that we tested **qpAdm** under almost ideal settings in our simulations
339 with left and right populations clearly separated and without gene flow between them. More thorough
340 assessments of the performance of **qpAdm** can be found elsewhere ([Harney et al., 2021](#); [Yüncü et al., 2023](#)).
341 In our simulations, unsupervised PSD-model approaches (**ADMIXTURE**, **NGSadmix**) work as well
342 as or even better than supervised PSD-model approaches (**ADMIXTURE**, **fastNGSadmix**) in estimating
343 the ancestry proportions in the target population. **ADMIXTURE** and **NGSadmix** benefit from LD pruning
344 while LD pruning increases the method bias for **fastNGSadmix** and introduces method bias for **qpAdm**.

345 Genotype likelihood-based methods for estimating ancestry proportions are not commonly used in
346 human ancient DNA studies (but they are popular as input for imputation pipelines). This may be
347 surprising, because genotype-likelihood-based approaches are targeted at low coverage data, exactly as
348 seen in ancient DNA studies. However, the definition of “low coverage” differs between fields. While
349 most working with modern DNA would understand 2-4X as “low depth”, the standards for ancient
350 DNA researchers are usually a lot lower due to limited DNA preservation. Genotype likelihood meth-
351 ods perform much better with $>1X$ coverage, an amount of data that is not within reach for most
352 ancient DNA samples investigated so far ([Mallick et al., 2023](#)). The large body of known, common
353 polymorphic sites in human populations allows the use of pseudohaploid calls at those positions in-
354 stead. Nonetheless, this study highlights that unsupervised methods employing genotype-likelihoods
355 (**NGSadmix**) can reach similar accuracies as methods such as **ADMIXTURE** that require (pseudo-haploid)
356 genotype calls. Moreover, methods that incorporate genotype likelihoods have the added benefit that
357 the modified genotype likelihood estimation approach can be used to reduce the effect of mapping bias.
358 Furthermore, if some samples in the dataset have $>1X$ depth, genotype likelihood-based approaches
359 will benefit from the additional data and provide more precise estimates of ancestry proportions while
360 pseudo-haploid data will not gain any information from more than one read at a position. Finally,
361 genotype likelihoods are very flexible and can be adjusted for many other aspects of the data. For
362 example, variations of genotype likelihood estimators exist that incorporate the effect of post-mortem
363 damage ([Hofmanová et al., 2016](#); [Link et al., 2017](#); [Kousathanas et al., 2017](#)) allowing to use of all
364 sequence data without filtering for potentially damaged sites or enzymatic repair of the damages in
365 the wet lab.

366 As the main aim of this study was to show the general impact of mapping bias and introduce a
367 modified genotype likelihood, we opted for a comparison of some of the most popular methods with a

368 limited set of settings. This was done in part to limit the computational load of this study. We also
369 decided to not set this up as a systematic assessment of different factors influencing mapping bias. The
370 effects of fragmentation (Günther and Nettelblad, 2019) and deamination damage (Martiniano et al.,
371 2020) on mapping bias have been explored elsewhere. Our results reiterate that mapping bias can
372 skew results in studies using low-coverage data as is the case in most ancient DNA studies. Different
373 strategies exist for mitigating these effects and we added a modified genotype likelihood approach
374 to the population genomic toolkit. Nevertheless, none of these methods will be the ideal solution in
375 all cases and they will not always fully remove the potential effect of mapping bias, making proper
376 verification and critical presentation of all results crucial.

377 Acknowledgements

378 We are extremely grateful to Amy Goldberg for numerous discussions during the initial phase of this
379 project. We thank Gabriel Renaud for making code for connecting `msprime` and `gargamel` available
380 on Github. The computations were enabled by resources in projects SNIC 2017/7-259, SNIC 2018/8-6,
381 SNIC 2021/2-17, SNIC 2022/22-874, NAISS 2023/22-883, sllstore2017087, UPPMAX 2023/2-30 and
382 NAISS 2023/2-19 provided by the National Academic Infrastructure for Supercomputing in Sweden
383 (NAISS) and the Swedish National Infrastructure for Computing (SNIC) at Uppmax, partially funded
384 by Uppsala University and the Swedish Research Council through grant agreements no. 2022-06725
385 and no. 2018-05973.

386 Funding

387 TG was supported by grants from the Swedish Research Council Vetenskapsrådet (2017-05267) and
388 Svenska Forskningsrådet Formas (2023-01381).

389 References

- 390 D. H. Alexander and K. Lange. Enhancements to the ADMIXTURE algorithm for individ-
391 ual ancestry estimation. *BMC Bioinformatics*, 12(1):246, June 2011. ISSN 1471-2105. doi:
392 10.1186/1471-2105-12-246. URL <https://doi.org/10.1186/1471-2105-12-246>.
- 393 D. H. Alexander, J. Novembre, and K. Lange. Fast model-based estimation of ancestry in unrelated
394 individuals. *Genome research*, 19(9):1655–1664, 2009. ISSN 1088-9051. Number: 9 Publisher: Cold
395 Spring Harbor Lab.
- 396 A. Auton, G. R. Abecasis, D. M. Altshuler, R. M. Durbin, G. R. Abecasis, D. R. Bentley,
397 A. Chakravarti, A. G. Clark, P. Donnelly, E. E. Eichler, P. Flicek, S. B. Gabriel, R. A. Gibbs, E. D.
398 Green, M. E. Hurles, B. M. Knoppers, J. O. Korb, E. S. Lander, C. Lee, H. Lehrach, E. R. Mardis,
399 G. T. Marth, G. A. McVean, D. A. Nickerson, J. P. Schmidt, S. T. Sherry, J. Wang, R. K. Wilson,
400 R. A. Gibbs, E. Boerwinkle, H. Doddapaneni, Y. Han, V. Korchina, C. Kovar, S. Lee, D. Muzny,
401 J. G. Reid, Y. Zhu, J. Wang, Y. Chang, Q. Feng, X. Fang, X. Guo, M. Jian, H. Jiang, X. Jin, T. Lan,
402 G. Li, J. Li, Y. Li, S. Liu, X. Liu, Y. Lu, X. Ma, M. Tang, B. Wang, G. Wang, H. Wu, R. Wu,
403 X. Xu, Y. Yin, D. Zhang, W. Zhang, J. Zhao, M. Zhao, X. Zheng, E. S. Lander, D. M. Altshuler,
404 S. B. Gabriel, N. Gupta, N. Gharani, L. H. Toji, N. P. Gerry, A. M. Resch, P. Flicek, J. Barker,
405 L. Clarke, L. Gil, S. E. Hunt, G. Kelman, E. Kulesha, R. Leinonen, W. M. McLaren, R. Rad-
406 hakrishnan, A. Roa, D. Smirnov, R. E. Smith, I. Streeter, A. Thormann, I. Toneva, B. Vaughan,
407 X. Zheng-Bradley, D. R. Bentley, R. Grocock, S. Humphray, T. James, Z. Kingsbury, H. Lehrach,
408 R. Sudbrak, M. W. Albrecht, V. S. Amstislavskiy, T. A. Borodina, M. Lienhard, F. Mertes, M. Sul-
409 tan, B. Timmermann, M.-L. Yaspo, E. R. Mardis, R. K. Wilson, L. Fulton, R. Fulton, S. T. Sherry,
410 V. Ananiev, Z. Belaia, D. Beloslyudtsev, N. Bouk, C. Chen, D. Church, R. Cohen, C. Cook, J. Gar-
411 ner, T. Hefferon, M. Kimelman, C. Liu, J. Lopez, P. Meric, C. O’Sullivan, Y. Ostapchuk, L. Phan,
412 S. Ponomarov, V. Schneider, E. Shekhtman, K. Sirotkin, D. Slotta, H. Zhang, G. A. McVean, R. M.

413 Durbin, S. Balasubramaniam, J. Burton, P. Danecek, T. M. Keane, A. Kolb-Kokocinski, S. Mc-
414 Carthy, J. Stalker, M. Quail, J. P. Schmidt, C. J. Davies, J. Gollub, T. Webster, B. Wong, Y. Zhan,
415 A. Auton, C. L. Campbell, Y. Kong, A. Marcketta, R. A. Gibbs, F. Yu, L. Antunes, M. Bainbridge,
416 D. Muzny, A. Sabo, Z. Huang, J. Wang, L. J. M. Coin, L. Fang, X. Guo, X. Jin, G. Li, Q. Li,
417 Y. Li, Z. Li, H. Lin, B. Liu, R. Luo, H. Shao, Y. Xie, C. Ye, C. Yu, F. Zhang, H. Zheng, H. Zhu,
418 C. Alkan, E. Dal, F. Kahveci, G. T. Marth, E. P. Garrison, D. Kural, W.-P. Lee, W. Fung Leong,
419 M. Stromberg, A. N. Ward, J. Wu, M. Zhang, M. J. Daly, M. A. DePristo, R. E. Handsaker, D. M.
420 Altshuler, E. Banks, G. Bhatia, G. del Angel, S. B. Gabriel, G. Genovese, N. Gupta, H. Li, S. Kashin,
421 E. S. Lander, S. A. McCarroll, J. C. Nemes, R. E. Poplin, S. C. Yoon, J. Lihm, V. Makarov, A. G.
422 Clark, S. Gottipati, A. Keinan, J. L. Rodriguez-Flores, J. O. Korb, T. Rausch, M. H. Fritz, A. M.
423 Stutz, P. Flicek, K. Beal, L. Clarke, A. Datta, J. Herrero, W. M. McLaren, G. R. S. Ritchie, R. E.
424 Smith, D. Zerbino, X. Zheng-Bradley, P. C. Sabeti, I. Shlyakhter, S. F. Schaffner, J. Vitti, D. N.
425 Cooper, E. V. Ball, P. D. Stenson, D. R. Bentley, B. Barnes, M. Bauer, R. Keira Cheetham, A. Cox,
426 M. Eberle, S. Humphray, S. Kahn, L. Murray, J. Peden, R. Shaw, E. E. Kenny, M. A. Batzer, M. K.
427 Konkel, J. A. Walker, D. G. MacArthur, M. Lek, R. Sudbrak, V. S. Amstislavskiy, R. Herwig, E. R.
428 Mardis, L. Ding, D. C. Koboldt, D. Larson, K. Ye, S. Gravel, A. Swaroop, E. Chew, T. Lappalainen,
429 Y. Erlich, M. Gymrek, T. Frederick Willems, J. T. Simpson, M. D. Shriver, J. A. Rosenfeld, C. D.
430 Bustamante, S. B. Montgomery, F. M. De La Vega, J. K. Byrnes, A. W. Carroll, M. K. DeGorter,
431 P. Lacroute, B. K. Maples, A. R. Martin, A. Moreno-Estrada, S. S. Shringarpure, F. Zakharia,
432 E. Halperin, Y. Baran, C. Lee, E. Cerveira, J. Hwang, A. Malhotra, D. Plewczynski, K. Radew,
433 M. Romanovitch, C. Zhang, F. C. L. Hyland, D. W. Craig, A. Christoforides, N. Homer, T. Izatt,
434 A. A. Kurdoglu, S. A. Sinari, K. Squire, S. T. Sherry, C. Xiao, J. Sebat, D. Antaki, M. Gujral,
435 A. Noor, K. Ye, E. G. Burchard, R. D. Hernandez, C. R. Gignoux, D. Haussler, S. J. Katzman,
436 W. James Kent, B. Howie, A. Ruiz-Linares, E. T. Dermitzakis, S. E. Devine, G. R. Abecasis,
437 H. Min Kang, J. M. Kidd, T. Blackwell, S. Caron, W. Chen, S. Emery, L. Fritsche, C. Fuchsberger,
438 G. Jun, B. Li, R. Lyons, C. Scheller, C. Sidore, S. Song, E. Sliwerska, D. Taliun, A. Tan, R. Welch,
439 M. Kate Wing, X. Zhan, P. Awadalla, A. Hodgkinson, Y. Li, X. Shi, A. Quitadamo, G. Lunter,
440 G. A. McVean, J. L. Marchini, S. Myers, C. Churchhouse, O. Delaneau, A. Gupta-Hinch, W. Kretz-
441 zschmar, Z. Iqbal, I. Mathieson, A. Menelaou, A. Rimmer, D. K. Xifara, T. K. Oleksyk, Y. Fu,
442 X. Liu, M. Xiong, L. Jorde, D. Witherspoon, J. Xing, E. E. Eichler, B. L. Browning, S. R. Brown-
443 ing, F. Hormozdiari, P. H. Sudmant, E. Khurana, R. M. Durbin, M. E. Hurles, C. Tyler-Smith,
444 C. A. Albers, Q. Ayub, S. Balasubramaniam, Y. Chen, V. Colonna, P. Danecek, L. Jostins, T. M.
445 Keane, S. McCarthy, K. Walter, Y. Xue, M. B. Gerstein, A. Abyzov, S. Balasubramaniam, J. Chen,
446 D. Clarke, Y. Fu, A. O. Harmanci, M. Jin, D. Lee, J. Liu, X. Jasmine Mu, J. Zhang, Y. Zhang,
447 Y. Li, R. Luo, H. Zhu, C. Alkan, E. Dal, F. Kahveci, G. T. Marth, E. P. Garrison, D. Kural, W.-P.
448 Lee, A. N. Ward, J. Wu, M. Zhang, S. A. McCarroll, R. E. Handsaker, D. M. Altshuler, E. Banks,
449 G. del Angel, G. Genovese, C. Hartl, H. Li, S. Kashin, J. C. Nemes, K. Shakir, S. C. Yoon,
450 J. Lihm, V. Makarov, J. Degenhardt, J. O. Korb, M. H. Fritz, S. Meiers, B. Raeder, T. Rausch,
451 A. M. Stutz, P. Flicek, F. Paolo Casale, L. Clarke, R. E. Smith, O. Stegle, X. Zheng-Bradley, D. R.
452 Bentley, B. Barnes, R. Keira Cheetham, M. Eberle, S. Humphray, S. Kahn, L. Murray, R. Shaw,
453 E.-W. Lameijer, M. A. Batzer, M. K. Konkel, J. A. Walker, L. Ding, I. Hall, K. Ye, P. Lacroute,
454 C. Lee, E. Cerveira, A. Malhotra, J. Hwang, D. Plewczynski, K. Radew, M. Romanovitch, C. Zhang,
455 D. W. Craig, N. Homer, D. Church, C. Xiao, J. Sebat, D. Antaki, V. Bafna, J. Michaelson, K. Ye,
456 S. E. Devine, E. J. Gardner, G. R. Abecasis, J. M. Kidd, R. E. Mills, G. Dayama, S. Emery,
457 G. Jun, X. Shi, A. Quitadamo, G. Lunter, G. A. McVean, K. Chen, X. Fan, Z. Chong, T. Chen,
458 D. Witherspoon, J. Xing, E. E. Eichler, M. J. Chaisson, F. Hormozdiari, J. Huddleston, M. Ma-
459 lig, B. J. Nelson, P. H. Sudmant, N. F. Parrish, E. Khurana, M. E. Hurles, B. Blackburne, S. J.
460 Lindsay, Z. Ning, K. Walter, Y. Zhang, M. B. Gerstein, A. Abyzov, J. Chen, D. Clarke, H. Lam,
461 X. Jasmine Mu, C. Sis, J. Zhang, Y. Zhang, R. A. Gibbs, F. Yu, M. Bainbridge, D. Challis, U. S.
462 Evani, C. Kovar, J. Lu, D. Muzny, U. Nagaswamy, J. G. Reid, A. Sabo, J. Yu, X. Guo, W. Li,
463 Y. Li, R. Wu, G. T. Marth, E. P. Garrison, W. Fung Leong, A. N. Ward, G. del Angel, M. A.

- 464 DePristo, S. B. Gabriel, N. Gupta, C. Hartl, R. E. Poplin, A. G. Clark, J. L. Rodriguez-Flores,
465 P. Flicek, L. Clarke, R. E. Smith, X. Zheng-Bradley, D. G. MacArthur, E. R. Mardis, R. Fulton,
466 D. C. Koboldt, S. Gravel, C. D. Bustamante, D. W. Craig, A. Christoforides, N. Homer, T. Izatt,
467 S. T. Sherry, C. Xiao, E. T. Dermitzakis, G. R. Abecasis, H. Min Kang, G. A. McVean, M. B.
468 Gerstein, S. Balasubramanian, L. Habegger, H. Yu, P. Flicek, L. Clarke, F. Cunningham, I. Dun-
469 ham, D. Zerbino, X. Zheng-Bradley, K. Lage, J. Berg Jaspersen, H. Horn, S. B. Montgomery, M. K.
470 DeGortor, E. Khurana, C. Tyler-Smith, Y. Chen, V. Colonna, Y. Xue, M. B. Gerstein, S. Balasubra-
471 manian, Y. Fu, D. Kim, A. Auton, A. Marnett, R. Desalle, A. Narechania, M. A. Wilson Sayres,
472 E. P. Garrison, R. E. Handsaker, S. Kashin, S. A. McCarroll, J. L. Rodriguez-Flores, P. Flicek,
473 L. Clarke, X. Zheng-Bradley, Y. Erlich, M. Gymrek, T. Frederick Willems, C. D. Bustamante, F. L.
474 Mendez, G. David Poznik, P. A. Underhill, C. Lee, E. Cervera, A. Malhotra, M. Romanovitch,
475 C. Zhang, G. R. Abecasis, L. Coin, H. Shao, D. Mittelman, C. Tyler-Smith, Q. Ayub, R. Banerjee,
476 M. Cerezo, Y. Chen, T. W. Fitzgerald, S. Louzada, A. Massaia, S. McCarthy, G. R. Ritchie, Y. Xue,
477 F. Yang, R. A. Gibbs, C. Kovar, D. Kalra, W. Hale, D. Muzny, J. G. Reid, J. Wang, X. Dan, X. Guo,
478 G. Li, Y. Li, C. Ye, X. Zheng, D. M. Altshuler, P. Flicek, L. Clarke, X. Zheng-Bradley, D. R. Bent-
479 ley, A. Cox, S. Humphray, S. Kahn, R. Sudbrak, M. W. Albrecht, M. Lienhard, D. Larson, D. W.
480 Craig, T. Izatt, A. A. Kurdoglu, S. T. Sherry, C. Xiao, D. Haussler, G. R. Abecasis, G. A. McVean,
481 R. M. Durbin, S. Balasubramanian, T. M. Keane, S. McCarthy, J. Stalker, A. Chakravarti, B. M.
482 Knoppers, G. R. Abecasis, K. C. Barnes, C. Beiswanger, E. G. Burchard, C. D. Bustamante, H. Cai,
483 H. Cao, R. M. Durbin, N. P. Gerry, N. Gharani, R. A. Gibbs, C. R. Gignoux, S. Gravel, B. Henn,
484 D. Jones, L. Jorde, J. S. Kaye, A. Keinan, A. Kent, A. Kerasidou, Y. Li, R. Mathias, G. A. McVean,
485 A. Moreno-Estrada, P. N. Ossorio, M. Parker, A. M. Resch, C. N. Rotimi, C. D. Royal, K. Sandoval,
486 Y. Su, R. Sudbrak, Z. Tian, S. Tishkoff, L. H. Toji, C. Tyler-Smith, M. Via, Y. Wang, H. Yang,
487 L. Yang, J. Zhu, W. Bodmer, G. Bedoya, A. Ruiz-Linares, Z. Cai, Y. Gao, J. Chu, L. Peltonen,
488 A. Garcia-Montero, A. Orfao, J. Dutil, J. C. Martinez-Cruzado, T. K. Oleksyk, K. C. Barnes,
489 R. A. Mathias, A. Hennis, H. Watson, C. McKenzie, F. Qadri, R. LaRocque, P. C. Sabeti, J. Zhu,
490 X. Deng, P. C. Sabeti, D. Asogun, O. Folarin, C. Happi, O. Omoniwa, M. Stremlau, R. Tariyal,
491 M. Jallow, F. Sisay Joof, T. Corrah, K. Rockett, D. Kwiatkowski, J. Kooner, T. T?nh Hi?n, S. J.
492 Dunstan, N. Thuy Hang, R. Fonnies, R. Garry, L. Kanneh, L. Moses, P. C. Sabeti, J. Schieffelin,
493 D. S. Grant, C. Gallo, G. Poletti, D. Saleheen, A. Rasheed, L. D. Brooks, A. L. Felsenfeld, J. E.
494 McEwen, Y. Vaydylevich, E. D. Green, A. Duncanson, M. Dunn, J. A. Schloss, J. Wang, H. Yang,
495 A. Auton, L. D. Brooks, R. M. Durbin, E. P. Garrison, H. Min Kang, J. O. Korbel, J. L. Marchini,
496 S. McCarthy, G. A. McVean, and G. R. Abecasis. A global reference for human genetic variation.
497 *Nature*, 526(7571):68–74, Sept. 2015. ISSN 0028-0836, 1476-4687. doi: 10.1038/nature15393. URL
498 <http://www.nature.com/doi/10.1038/nature15393>.
- 499 A. R. Barton, C. G. Santander, P. Skoglund, I. Moltke, D. Reich, and I. Mathieson. Insuffi-
500 cient evidence for natural selection associated with the Black Death, Mar. 2023. URL <https://www.biorxiv.org/content/10.1101/2023.03.14.532615v1>. Pages: 2023.03.14.532615 Section:
501 Contradictory Results.
502
- 503 A. W. Briggs, U. Stenzel, P. L. Johnson, R. E. Green, J. Kelso, K. Prüfer, M. Meyer, J. Krause,
504 M. T. Ronan, M. Lachmann, and others. Patterns of damage in genomic DNA sequences from a
505 Neandertal. *Proceedings of the National Academy of Sciences*, 104(37):14616–14621, 2007.
- 506 C. C. Chang, C. C. Chow, L. C. Tellier, S. Vattikuti, S. M. Purcell, and J. J. Lee. Second-generation
507 PLINK: rising to the challenge of larger and richer datasets. *Gigascience*, 4(1):s13742–015, 2015.
508 ISSN 2047-217X. Number: 1 Publisher: Oxford University Press.
- 509 N.-C. Chen, B. Solomon, T. Mun, S. Iyer, and B. Langmead. Reference flow: reducing reference bias
510 using multiple population genomes. *Genome Biology*, 22(1):8, Jan. 2021. ISSN 1474-760X. doi:
511 [10.1186/s13059-020-02229-3](https://doi.org/10.1186/s13059-020-02229-3). URL <https://doi.org/10.1186/s13059-020-02229-3>.

- 512 D. M. Church, V. A. Schneider, K. M. Steinberg, M. C. Schatz, A. R. Quinlan, C.-S. Chin, P. A. Kitts,
513 B. Aken, G. T. Marth, M. M. Hoffman, J. Herrero, M. L. Z. Mendoza, R. Durbin, and P. Flicek.
514 Extending reference assembly models. *Genome Biology*, 16(1):13, Jan. 2015. ISSN 1465-6906. doi:
515 10.1186/s13059-015-0587-3. URL <https://doi.org/10.1186/s13059-015-0587-3>.
- 516 S. L. Cox, H. M. Moots, J. T. Stock, A. Shbat, B. D. Bitarello, N. Nicklisch, K. W. Alt,
517 W. Haak, E. Rosenstock, C. B. Ruff, and I. Mathieson. Predicting skeletal stature using ancient
518 DNA. *American Journal of Biological Anthropology*, 177(1):162–174, 2022. ISSN 2692-7691. doi:
519 10.1002/ajpa.24426. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/ajpa.24426>.
520 eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/ajpa.24426>.
- 521 T. Davy, D. Ju, I. Mathieson, and P. Skoglund. Hunter-gatherer admixture facilitated natural selection
522 in Neolithic European farmers. *Current Biology*, 33(7):1365–1371.e3, Apr. 2023. ISSN 0960-9822.
523 doi: 10.1016/j.cub.2023.02.049. URL <https://www.sciencedirect.com/science/article/pii/S0960982223001896>.
524 [S0960982223001896](https://www.sciencedirect.com/science/article/pii/S0960982223001896).
- 525 D. Falush, M. Stephens, and J. K. Pritchard. Inference of population structure using multilocus
526 genotype data: linked loci and correlated allele frequencies. *Genetics*, 164(4):1567–1587, 2003.
527 ISSN 0016-6731. Number: 4.
- 528 D. Falush, M. Stephens, and J. K. Pritchard. Inference of population structure using multilocus
529 genotype data: dominant markers and null alleles. *Molecular ecology notes*, 7(4):574–578, 2007.
530 ISSN 1471-8278. Number: 4.
- 531 M. Fumagalli, F. G. Vieira, T. S. Korneliussen, T. Linderoth, E. Huerta-Sánchez, A. Albrechtsen,
532 and R. Nielsen. Quantifying Population Genetic Differentiation from Next-Generation Sequencing
533 Data. *Genetics*, 195(3):979–992, Nov. 2013. ISSN 1943-2631. doi: 10.1534/genetics.113.154740.
534 URL <https://doi.org/10.1534/genetics.113.154740>.
- 535 S. Gopalakrishnan, J. A. Samaniego Castruita, M.-H. S. Sinding, L. F. K. Kuderna, J. Räikkönen,
536 B. Petersen, T. Sicheritz-Ponten, G. Larson, L. Orlando, T. Marques-Bonet, A. J. Hansen, L. Dalén,
537 and M. T. P. Gilbert. The wolf reference genome sequence (*Canis lupus lupus*) and its implications
538 for *Canis* spp. population genomics. *BMC Genomics*, 18:495, June 2017. ISSN 1471-2164. doi:
539 10.1186/s12864-017-3883-3. URL <https://doi.org/10.1186/s12864-017-3883-3>.
- 540 S. Gopalakrishnan, S. S. Ebenesersdóttir, I. K. C. Lundstrøm, G. Turner-Walker, K. H. S. Moore,
541 P. Luisi, A. Margaryan, M. D. Martin, M. R. Ellegaard, Magnússon, Sigursson, S. Snorradóttir,
542 D. N. Magnúsdóttir, J. E. Laffoon, L. van Dorp, X. Liu, I. Moltke, M. C. Ávila Arcos, J. G.
543 Schraiber, S. Rasmussen, D. Juan, P. Gelabert, T. de Dios, A. K. Fotakis, M. Iraeta-Orbegozo,
544 J. Vågane, S. D. Denham, A. Christophersen, H. K. Stenøien, F. G. Vieira, S. Liu, T. Günther,
545 T. Kivisild, O. G. Moseng, B. Skar, C. Cheung, M. Sandoval-Velasco, N. Wales, H. Schroeder, P. F.
546 Campos, V. B. Gumundsdóttir, T. Sicheritz-Ponten, B. Petersen, J. Halgunset, E. Gilbert, G. L.
547 Cavalleri, E. Hovig, I. Kockum, T. Olsson, L. Alfredsson, T. F. Hansen, T. Werge, E. Willerslev,
548 F. Balloux, T. Marques-Bonet, C. Lalueza-Fox, R. Nielsen, K. Stefánsson, A. Helgason, and M. T. P.
549 Gilbert. The population genomic legacy of the second plague pandemic. *Current Biology*, 32
550 (21):4743–4751.e6, Nov. 2022. ISSN 0960-9822. doi: 10.1016/j.cub.2022.09.023. URL <https://www.sciencedirect.com/science/article/pii/S0960982222014671>.
551 [S0960982222014671](https://www.sciencedirect.com/science/article/pii/S0960982222014671).
- 552 R. E. Green, J. Krause, A. W. Briggs, T. Maricic, U. Stenzel, M. Kircher, N. Patterson, H. Li, W. Zhai,
553 and M. H.-Y. Fritz. A draft sequence of the Neandertal genome. *science*, 328(5979):710–722, 2010.
554 ISSN 0036-8075. Number: 5979 Publisher: American Association for the Advancement of Science.
- 555 T. Günther and C. Nettelblad. The presence and impact of reference bias on population genomic stud-
556 ies of prehistoric human populations. *PLOS Genetics*, 15(7):e1008302, July 2019. ISSN 1553-7404.
557 doi: 10.1371/journal.pgen.1008302. URL [https://journals.plos.org/plosgenetics/article?
558 id=10.1371/journal.pgen.1008302](https://journals.plos.org/plosgenetics/article?id=10.1371/journal.pgen.1008302).

- 559 W. Haak, I. Lazaridis, N. Patterson, N. Rohland, S. Mallick, B. Llamas, G. Brandt, S. Nordenfelt,
560 E. Harney, and K. Stewardson. Massive migration from the steppe was a source for Indo-European
561 languages in Europe. *Nature*, 522(7555):207–211, 2015. ISSN 1476-4687. Number: 7555 Publisher:
562 Nature Publishing Group.
- 563 K. Hanghøj, I. Moltke, P. A. Andersen, A. Manica, and T. S. Korneliussen. Fast and accu-
564 rate relatedness estimation from high-throughput sequencing data in the presence of inbreed-
565 ing. *GigaScience*, 8(5), May 2019. ISSN 2047-217X. doi: 10.1093/gigascience/giz034. URL
566 <https://doi.org/10.1093/gigascience/giz034>.
- 567 E. Harney, N. Patterson, D. Reich, and J. Wakeley. Assessing the performance of qpAdm: a statistical
568 tool for studying population admixture. *Genetics*, 217(4), Apr. 2021. ISSN 1943-2631. doi: 10.
569 1093/genetics/iyaa045. URL <https://doi.org/10.1093/genetics/iyaa045>.
- 570 P. D. Heintzman, G. D. Zazula, R. D. MacPhee, E. Scott, J. A. Cahill, B. K. McHorse, J. D. Kapp,
571 M. Stiller, M. J. Wooller, L. Orlando, J. Southon, D. G. Froese, and B. Shapiro. A new genus of
572 horse from Pleistocene North America. *eLife*, 6, 2017. ISSN 2050-084X. doi: 10.7554/eLife.29944.
- 573 Z. Hofmanová, S. Kreutzer, G. Hellenthal, C. Sell, Y. Diekmann, D. Díez-del Molino, L. van Dorp,
574 S. López, A. Kousathanas, V. Link, and others. Early farmers from across Europe directly descended
575 from Neolithic Aegeans. *Proceedings of the National Academy of Sciences*, page 201523951, 2016.
- 576 L. Huang, V. Popic, and S. Batzoglou. Short read alignment with populations of genomes. *Bioin-*
577 *formatics*, 29(13):i361–i370, July 2013. ISSN 1367-4803. doi: 10.1093/bioinformatics/btt215. URL
578 <https://doi.org/10.1093/bioinformatics/btt215>.
- 579 M. J. Hubisz, D. Falush, M. Stephens, and J. K. Pritchard. Inferring weak population structure with
580 the assistance of sample group information. *Molecular ecology resources*, 9(5):1322–1332, 2009. ISSN
581 1755-098X. Number: 5.
- 582 R. Hui, C. L. Scheib, E. D’Atanasio, S. A. Inskip, C. Cessford, S. A. Biagini, A. W. Wohms, M. Q.
583 Ali, S. J. Griffith, A. Solnik, H. Niinemäe, X. J. Ge, A. K. Rose, O. Beneker, T. C. O’Connell, J. E.
584 Robb, and T. Kivisild. Genetic history of Cambridgeshire before and after the Black Death. *Science*
585 *Advances*, 10(3):eadi5903, Jan. 2024. doi: 10.1126/sciadv.adi5903. URL <https://www.science.org/doi/10.1126/sciadv.adi5903>. Publisher: American Association for the Advancement of
586 Science.
587
- 588 E. Jørsboe, K. Hanghøj, and A. Albrechtsen. fastNGSadmix: admixture proportions and principal
589 component analysis of a single NGS sample. *Bioinformatics*, 33(19):3148–3150, 2017.
- 590 J. Kelleher, A. M. Etheridge, and G. McVean. Efficient coalescent simulation and genealogical analysis
591 for large sample sizes. *PLoS computational biology*, 12(5):e1004842, 2016.
- 592 J. Klunk, T. P. Vilgalys, C. E. Demeure, X. Cheng, M. Shiratori, J. Madej, R. Beau, D. Elli, M. I.
593 Patino, R. Redfern, S. N. DeWitte, J. A. Gamble, J. L. Boldsen, A. Carmichael, N. Varlik, K. Eaton,
594 J.-C. Grenier, G. B. Golding, A. Devault, J.-M. Rouillard, V. Yotova, R. Sindeaux, C. J. Ye,
595 M. Bikaran, A. Dumaine, J. F. Brinkworth, D. Missiakas, G. A. Rouleau, M. Steinrücken, J. Pizarro-
596 Cerdá, H. N. Poinar, and L. B. Barreiro. Evolution of immune genes is associated with the Black
597 Death. *Nature*, 611(7935):312–319, Nov. 2022. ISSN 1476-4687. doi: 10.1038/s41586-022-05349-x.
598 URL <https://www.nature.com/articles/s41586-022-05349-x>. Number: 7935 Publisher: Na-
599 ture Publishing Group.
- 600 D. Koptekin, E. Yapar, K. B. Vural, E. Sağlıcan, N. E. Altınışık, A.-S. Malaspinas, C. Alkan, and
601 M. Somel. Pre-processing of paleogenomes: Mitigating reference bias and postmortem damage in
602 ancient genome data, Nov. 2023. URL [https://www.biorxiv.org/content/10.1101/2023.11.
603 11.566695v1](https://www.biorxiv.org/content/10.1101/2023.11.11.566695v1). Pages: 2023.11.11.566695 Section: New Results.

- 604 T. S. Korneliussen and I. Moltke. NgsRelate: a software tool for estimating pairwise relatedness
605 from next-generation sequencing data. *Bioinformatics*, 31(24):4009–4011, 2015. ISSN 1460-2059.
606 Number: 24 Publisher: Oxford University Press.
- 607 T. S. Korneliussen, I. Moltke, A. Albrechtsen, and R. Nielsen. Calculation of Tajima’s D and other
608 neutrality test statistics from low depth next-generation sequencing data. *BMC bioinformatics*, 14:
609 289, Oct. 2013. ISSN 1471-2105. doi: 10.1186/1471-2105-14-289.
- 610 T. S. Korneliussen, A. Albrechtsen, and R. Nielsen. ANGSD: Analysis of Next Generation Sequencing
611 Data. *BMC bioinformatics*, 15(1):356, 2014. ISSN 1471-2105. doi: 10.1186/s12859-014-0356-4.
- 612 A. Kousathanas, C. Leuenberger, V. Link, C. Sell, J. Burger, and D. Wegmann. Inferring Heterozygos-
613 ity from Ancient and Low Coverage Genomes. *Genetics*, 205(1):317–332, Jan. 2017. ISSN 0016-6731,
614 1943-2631. doi: 10.1534/genetics.116.189985. URL [http://www.genetics.org/content/205/1/
615 317](http://www.genetics.org/content/205/1/317).
- 616 E. S. Lander, L. M. Linton, B. Birren, C. Nusbaum, M. C. Zody, J. Baldwin, K. Devon, K. De-
617 war, M. Doyle, W. FitzHugh, R. Funke, D. Gage, K. Harris, A. Heaford, J. Howland, L. Kann,
618 J. Lehoczy, R. LeVine, P. McEwan, K. McKernan, J. Meldrim, J. P. Mesirov, C. Miranda, W. Mor-
619 ris, J. Naylor, C. Raymond, M. Rosetti, R. Santos, A. Sheridan, C. Sougnez, N. Stange-Thomann,
620 N. Stojanovic, A. Subramanian, D. Wyman, J. Rogers, J. Sulston, R. Ainscough, S. Beck, D. Bent-
621 ley, J. Burton, C. Clee, N. Carter, A. Coulson, R. Deadman, P. Deloukas, A. Dunham, I. Dunham,
622 R. Durbin, L. French, D. Grafham, S. Gregory, T. Hubbard, S. Humphray, A. Hunt, M. Jones,
623 C. Lloyd, A. McMurray, L. Matthews, S. Mercer, S. Milne, J. C. Mullikin, A. Mungall, R. Plumb,
624 M. Ross, R. Shownkeen, S. Sims, R. H. Waterston, R. K. Wilson, L. W. Hillier, J. D. McPherson,
625 M. A. Marra, E. R. Mardis, L. A. Fulton, A. T. Chinwalla, K. H. Pepin, W. R. Gish, S. L. Chissoe,
626 M. C. Wendl, K. D. Delehaunty, T. L. Miner, A. Delehaunty, J. B. Kramer, L. L. Cook, R. S. Fulton,
627 D. L. Johnson, P. J. Minx, S. W. Clifton, T. Hawkins, E. Branscomb, P. Predki, P. Richardson,
628 S. Wenning, T. Slezak, N. Doggett, J.-F. Cheng, A. Olsen, S. Lucas, C. Elkin, E. Uberbacher, M. Fra-
629 zier, R. A. Gibbs, D. M. Muzny, S. E. Scherer, J. B. Bouck, E. J. Sodergren, K. C. Worley, C. M.
630 Rives, J. H. Gorrell, M. L. Metzker, S. L. Naylor, R. S. Kucherlapati, D. L. Nelson, G. M. Wein-
631 stock, Y. Sakaki, A. Fujiyama, M. Hattori, T. Yada, A. Toyoda, T. Itoh, C. Kawagoe, H. Watanabe,
632 Y. Totoki, T. Taylor, J. Weissenbach, R. Heilig, W. Saurin, F. Artiguenave, P. Brottier, T. Bruls,
633 E. Pelletier, C. Robert, P. Wincker, A. Rosenthal, M. Platzer, G. Nyakatura, S. Taudien, A. Rump,
634 D. R. Smith, L. Doucette-Stamm, M. Rubenfield, K. Weinstock, H. M. Lee, J. Dubois, H. Yang,
635 J. Yu, J. Wang, G. Huang, J. Gu, L. Hood, L. Rowen, A. Madan, S. Qin, R. W. Davis, N. A. Feder-
636 spiel, A. P. Abola, M. J. Proctor, B. A. Roe, F. Chen, H. Pan, J. Ramser, H. Lehrach, R. Reinhardt,
637 W. R. McCombie, M. de la Bastide, N. Dedhia, H. Blöcker, K. Hornischer, G. Nordsiek, R. Agarwala,
638 L. Aravind, J. A. Bailey, A. Bateman, S. Batzoglou, E. Birney, P. Bork, D. G. Brown, C. B. Burge,
639 L. Cerutti, H.-C. Chen, D. Church, M. Clamp, R. R. Copley, T. Doerks, S. R. Eddy, E. E. Eichler,
640 T. S. Furey, J. Galagan, J. G. R. Gilbert, C. Harmon, Y. Hayashizaki, D. Haussler, H. Hermjakob,
641 K. Hokamp, W. Jang, L. S. Johnson, T. A. Jones, S. Kasif, A. Kasprzyk, S. Kennedy, W. J. Kent,
642 P. Kitts, E. V. Koonin, I. Korf, D. Kulp, D. Lancet, T. M. Lowe, A. McLysaght, T. Mikkelsen, J. V.
643 Moran, N. Mulder, V. J. Pollara, C. P. Ponting, G. Schuler, J. Schultz, G. Slater, A. F. A. Smit,
644 E. Stupka, J. Szustakowki, D. Thierry-Mieg, J. Thierry-Mieg, L. Wagner, J. Wallis, R. Wheeler,
645 A. Williams, Y. I. Wolf, K. H. Wolfe, S.-P. Yang, R.-F. Yeh, F. Collins, M. S. Guyer, J. Peterson,
646 A. Felsenfeld, K. A. Wetterstrand, R. M. Myers, J. Schmutz, M. Dickson, J. Grimwood, D. R.
647 Cox, M. V. Olson, R. Kaul, C. Raymond, N. Shimizu, K. Kawasaki, S. Minoshima, G. A. Evans,
648 M. Athanasiou, R. Schultz, A. Patrinos, M. J. Morgan, International Human Genome Sequencing
649 Consortium, C. f. G. R. Whitehead Institute for Biomedical Research, The Sanger Centre:, Wash-
650 ington University Genome Sequencing Center, US DOE Joint Genome Institute:, Baylor College
651 of Medicine Human Genome Sequencing Center:, RIKEN Genomic Sciences Center:, Genoscope
652 and CNRS UMR-8030:, I. o. M. B. Department of Genome Analysis, GTC Sequencing Center:,

- 653 Beijing Genomics Institute/Human Genome Center:, T. I. f. S. B. Multimegabase Sequencing Cen-
654 ter, Stanford Genome Technology Center:, University of Oklahoma’s Advanced Center for Genome
655 Technology:, Max Planck Institute for Molecular Genetics:, L. A. H. G. C. Cold Spring Harbor Lab-
656 oratory, GBF—German Research Centre for Biotechnology:, a. i. i. l. u. o. h. *Genome Analysis
657 Group (listed in alphabetical order, U. N. I. o. H. Scientific management: National Human Genome
658 Research Institute, Stanford Human Genome Center:, University of Washington Genome Center:,
659 K. U. S. o. M. Department of Molecular Biology, University of Texas Southwestern Medical Center
660 at Dallas:, U. D. o. E. Office of Science, and The Wellcome Trust:. Initial sequencing and analysis of
661 the human genome. *Nature*, 409(6822):860–921, Feb. 2001. ISSN 1476-4687. doi: 10.1038/35057062.
662 URL <https://www.nature.com/articles/35057062>. Number: 6822 Publisher: Nature Publishing
663 Group.
- 664 H. Li and R. Durbin. Fast and accurate short read alignment with Burrows–Wheeler transform.
665 *bioinformatics*, 25(14):1754–1760, 2009. ISSN 1367-4803. Number: 14 Publisher: Oxford University
666 Press.
- 667 H. Li, B. Handsaker, A. Wysoker, T. Fennell, J. Ruan, N. Homer, G. Marth, G. Abecasis, R. Durbin,
668 and 1000 Genome Project Data Processing Subgroup. The Sequence Alignment/Map format and
669 SAMtools. *Bioinformatics (Oxford, England)*, 25(16):2078–2079, Aug. 2009. ISSN 1367-4811. doi:
670 10.1093/bioinformatics/btp352.
- 671 V. Link, A. Kousathanas, K. Veeramah, C. Sell, A. Scheu, and D. Wegmann. ATLAS: analysis tools
672 for low-depth and ancient samples. *bioRxiv*, page 105346, 2017.
- 673 R. N. Lou, A. Jacobs, A. P. Wilder, and N. O. Therkildsen. A beginner’s guide to low-coverage
674 whole genome sequencing for population genomics. *Molecular Ecology*, 30(23):5966–5993, 2021.
675 ISSN 1365-294X. doi: 10.1111/mec.16077. URL [https://onlinelibrary.wiley.com/doi/abs/](https://onlinelibrary.wiley.com/doi/abs/10.1111/mec.16077)
676 [10.1111/mec.16077](https://onlinelibrary.wiley.com/doi/abs/10.1111/mec.16077). eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/mec.16077>.
- 677 S. Mallick, A. Micco, M. Mah, H. Ringbauer, I. Lazaridis, I. Olalde, N. Patterson, and D. Reich. The
678 Allen Ancient DNA Resource (AADR): A curated compendium of ancient human genomes, Apr.
679 2023. URL <https://www.biorxiv.org/content/10.1101/2023.04.06.535797v1>.
- 680 R. Martiniano, E. Garrison, E. R. Jones, A. Manica, and R. Durbin. Removing reference bias and
681 improving indel calling in ancient DNA data analysis by mapping to a sequence variation graph.
682 *Genome Biology*, 21(1):250, Sept. 2020. ISSN 1474-760X. doi: 10.1186/s13059-020-02160-7. URL
683 <https://doi.org/10.1186/s13059-020-02160-7>.
- 684 I. Mathieson and J. Terhorst. Direct detection of natural selection in Bronze Age Britain. *Genome*
685 *Research*, 32(11-12):2057–2067, Nov. 2022. ISSN 1088-9051, 1549-5469. doi: 10.1101/gr.276862.122.
686 URL <https://genome.cshlp.org/content/32/11-12/2057>. Company: Cold Spring Harbor Lab-
687 oratory Press Distributor: Cold Spring Harbor Laboratory Press Institution: Cold Spring Harbor
688 Laboratory Press Label: Cold Spring Harbor Laboratory Press Publisher: Cold Spring Harbor Lab.
- 689 I. Mathieson, I. Lazaridis, N. Rohland, S. Mallick, N. Patterson, S. A. Roodenberg, E. Harney, K. Stew-
690 ardson, D. Fernandes, M. Novak, and others. Genome-wide patterns of selection in 230 ancient
691 Eurasians. *Nature*, 528(7583):499–503, 2015.
- 692 I. Mathieson, F. Abascal, L. Vinner, P. Skoglund, C. Pomilla, P. Mitchell, C. Arthur, D. Gurdasani,
693 E. Willerslev, M. S. Sandhu, and G. Dewar. An Ancient Baboon Genome Demonstrates Long-Term
694 Population Continuity in Southern Africa. *Genome Biology and Evolution*, 12(4):407–412, Apr.
695 2020. ISSN 1759-6653. doi: 10.1093/gbe/evaa019. URL <https://doi.org/10.1093/gbe/evaa019>.
- 696 A. McKenna, M. Hanna, E. Banks, A. Sivachenko, K. Cibulskis, A. Kernytsky, K. Garimella, D. Alt-
697 shuler, S. Gabriel, M. Daly, and M. A. DePristo. The Genome Analysis Toolkit: A MapReduce

- 698 framework for analyzing next-generation DNA sequencing data. *Genome Research*, 20(9):1297–1303,
699 Sept. 2010. ISSN 1088-9051. doi: 10.1101/gr.107524.110. URL [https://www.ncbi.nlm.nih.gov/
700 pmc/articles/PMC2928508/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2928508/).
- 701 J. Meisner and A. Albrechtsen. Inferring population structure and admixture proportions in low-
702 depth NGS data. *Genetics*, 210(2):719–731, 2018. ISSN 1943-2631. Number: 2 Publisher: Oxford
703 University Press.
- 704 R. Nielsen, J. S. Paul, A. Albrechtsen, and Y. S. Song. Genotype and SNP calling from next-generation
705 sequencing data. *Nature Reviews Genetics*, 12(6):443, 2011.
- 706 A. K. Nøhr, K. Hanghøj, G. Garcia-Erill, Z. Li, I. Moltke, and A. Albrechtsen. NGSremix: a soft-
707 ware tool for estimating pairwise relatedness between admixed individuals from next-generation
708 sequencing data. *G3*, (jkab174), May 2021. ISSN 2160-1836. doi: 10.1093/g3journal/jkab174. URL
709 <https://doi.org/10.1093/g3journal/jkab174>.
- 710 A. Oliva, R. Tobler, A. Cooper, B. Llamas, and Y. Souilmi. Systematic benchmark of ancient DNA
711 read mapping. *Briefings in Bioinformatics*, (bbab076), Apr. 2021. ISSN 1477-4054. doi: 10.1093/
712 bib/bbab076. URL <https://doi.org/10.1093/bib/bbab076>.
- 713 L. Orlando, A. Ginolhac, G. Zhang, D. Froese, A. Albrechtsen, M. Stiller, M. Schubert, E. Cap-
714 pellini, B. Petersen, I. Moltke, P. L. F. Johnson, M. Fumagalli, J. T. Vilstrup, M. Raghavan,
715 T. Korneliussen, A.-S. Malaspinas, J. Vogt, D. Szklarczyk, C. D. Kelstrup, J. Vinther, A. Dolocan,
716 J. Stenderup, A. M. V. Velazquez, J. Cahill, M. Rasmussen, X. Wang, J. Min, G. D. Zazula,
717 A. Seguin-Orlando, C. Mortensen, K. Magnussen, J. F. Thompson, J. Weinstock, K. Gregersen,
718 K. H. Røed, V. Eisenmann, C. J. Rubin, D. C. Miller, D. F. Antczak, M. F. Bertelsen, S. Brunak,
719 K. A. S. Al-Rasheid, O. Ryder, L. Andersson, J. Mundy, A. Krogh, M. T. P. Gilbert, K. Kjær,
720 T. Sicheritz-Ponten, L. J. Jensen, J. V. Olsen, M. Hofreiter, R. Nielsen, B. Shapiro, J. Wang, and
721 E. Willerslev. Recalibrating Equus evolution using the genome sequence of an early Middle Pleis-
722 tocene horse. *Nature*, 499(7456):74–78, July 2013. ISSN 1476-4687. doi: 10.1038/nature12323.
723 URL <https://www.nature.com/articles/nature12323>. Bandiera_abtest: a Cg_type: Nature Re-
724 search Journals Number: 7456 Primary_atype: Research Publisher: Nature Publishing Group Sub-
725 ject_term: Evolutionary genetics Subject_term_id: evolutionary-genetics.
- 726 L. Orlando, R. Allaby, P. Skoglund, C. Der Sarkissian, P. W. Stockhammer, M. C. Ávila Arcos,
727 Q. Fu, J. Krause, E. Willerslev, A. C. Stone, and C. Warinner. Ancient DNA analysis. *Nature*
728 *Reviews Methods Primers*, 1(1):1–26, Feb. 2021. ISSN 2662-8449. doi: 10.1038/s43586-020-00011-0.
729 URL <https://www.nature.com/articles/s43586-020-00011-0>. Number: 1 Publisher: Nature
730 Publishing Group.
- 731 N. Patterson, A. L. Price, and D. Reich. Population structure and eigenanalysis. *PLoS genetics*, 2
732 (12):e190, 2006. ISSN 1553-7390. Number: 12 Publisher: Public Library of Science San Francisco,
733 USA.
- 734 N. Patterson, P. Moorjani, Y. Luo, S. Mallick, N. Rohland, Y. Zhan, T. Genschoreck, T. Webster, and
735 D. Reich. Ancient admixture in human history. *Genetics*, 192(3):1065–1093, 2012. ISSN 1943-2631.
736 Number: 3 Publisher: Oxford University Press.
- 737 A. Prasad, E. D. Lorenzen, and M. V. Westbury. Evaluating the role of reference-genome phy-
738 logenetic distance on evolutionary inference. *Molecular Ecology Resources*, 22(1):45–55, 2022.
739 ISSN 1755-0998. doi: 10.1111/1755-0998.13457. URL [https://onlinelibrary.wiley.com/doi/
740 abs/10.1111/1755-0998.13457](https://onlinelibrary.wiley.com/doi/abs/10.1111/1755-0998.13457). eprint: [https://onlinelibrary.wiley.com/doi/pdf/10.1111/1755-
741 0998.13457](https://onlinelibrary.wiley.com/doi/pdf/10.1111/1755-0998.13457).

- 742 A. L. Price, N. J. Patterson, R. M. Plenge, M. E. Weinblatt, N. A. Shadick, and D. Reich. Principal
743 components analysis corrects for stratification in genome-wide association studies. *Nature genetics*,
744 38(8):904–909, 2006. ISSN 1546-1718. Number: 8 Publisher: Nature Publishing Group.
- 745 J. K. Pritchard, M. Stephens, and P. Donnelly. Inference of population structure using multilocus
746 genotype data. *Genetics*, 155(2):945–959, 2000. ISSN 0016-6731. Number: 2.
- 747 K. Prüfer. snpAD: An ancient DNA genotype caller. *Bioinformatics*, 2018. doi: 10.1093/
748 bioinformatics/bty507. URL [https://academic.oup.com/bioinformatics/advance-article/
749 doi/10.1093/bioinformatics/bty507/5042170](https://academic.oup.com/bioinformatics/advance-article/doi/10.1093/bioinformatics/bty507/5042170).
- 750 G. Renaud, K. Hanghøj, E. Willerslev, and L. Orlando. gargammel: a sequence simulator for ancient
751 DNA. *Bioinformatics*, 33(4):577–579, Feb. 2017. ISSN 1367-4803. doi: 10.1093/bioinformatics/
752 btw670. URL <https://academic.oup.com/bioinformatics/article/33/4/577/2608651>.
- 753 S. Rubinacci, D. M. Ribeiro, R. J. Hofmeister, and O. Delaneau. Efficient phasing and imputa-
754 tion of low-coverage sequencing data using large reference panels. *Nature Genetics*, 53(1):120–126,
755 Jan. 2021. ISSN 1546-1718. doi: 10.1038/s41588-020-00756-0. URL [https://www.nature.com/
756 articles/s41588-020-00756-0](https://www.nature.com/articles/s41588-020-00756-0). Number: 1 Publisher: Nature Publishing Group.
- 757 M. Schubert, A. Ginolhac, S. Lindgreen, J. F. Thompson, K. A. AL-Rasheid, E. Willerslev, A. Krogh,
758 and L. Orlando. Improving ancient DNA read mapping against modern reference genomes. *BMC*
759 *Genomics*, 13:178, May 2012. ISSN 1471-2164. doi: 10.1186/1471-2164-13-178. URL [https:
760 //doi.org/10.1186/1471-2164-13-178](https://doi.org/10.1186/1471-2164-13-178).
- 761 M. Schubert, S. Lindgreen, and L. Orlando. AdapterRemoval v2: rapid adapter trimming, identifica-
762 tion, and read merging. *BMC research notes*, 9(1):1–7, 2016. ISSN 1756-0500. Number: 1 Publisher:
763 BioMed Central.
- 764 L. Skotte, T. S. Korneliussen, and A. Albrechtsen. Estimating individual admixture proportions from
765 next generation sequencing data. *Genetics*, 195(3):693–702, 2013. ISSN 1943-2631. Number: 3
766 Publisher: Oxford University Press.
- 767 D.-M. J. Thorburn, K. Sagonas, M. Binzer-Panchal, F. J. J. Chain, P. G. D. Feulner, E. Bornberg-
768 Bauer, T. B. H. Reusch, I. E. Samonte-Padilla, M. Milinski, T. L. Lenz, and C. Eizaguirre.
769 Origin matters: Using a local reference genome improves measures in population genomics.
770 *Molecular Ecology Resources*, 23(7):1706–1723, 2023. ISSN 1755-0998. doi: 10.1111/1755-0998.
771 13838. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/1755-0998.13838>. eprint:
772 <https://onlinelibrary.wiley.com/doi/pdf/10.1111/1755-0998.13838>.
- 773 T. van der Valk, C. M. Gonda, H. Silegowa, S. Almanza, I. Sifuentes-Romero, T. B. Hart, J. A. Hart,
774 K. M. Detwiler, and K. Guschanski. The Genome of the Endangered Dryas Monkey Provides New
775 Insights into the Evolutionary History of the Vervets. *Molecular Biology and Evolution*, 37(1):183–
776 194, Jan. 2020. ISSN 0737-4038. doi: 10.1093/molbev/msz213. URL [https://doi.org/10.1093/
777 molbev/msz213](https://doi.org/10.1093/molbev/msz213).
- 778 E. Yüncü, U. Işıldak, M. P. Williams, C. D. Huber, L. A. Vyazov, P. Changmai, and P. Flegontov. False
779 discovery rates of qpAdm-based screens for genetic admixture. *bioRxiv*, Apr. 2023. doi: 10.1101/
780 2023.04.25.538339. URL <https://www.biorxiv.org/content/10.1101/2023.04.25.538339v1>.
781 Pages: 2023.04.25.538339 Section: New Results.

782

Supplementary Figures

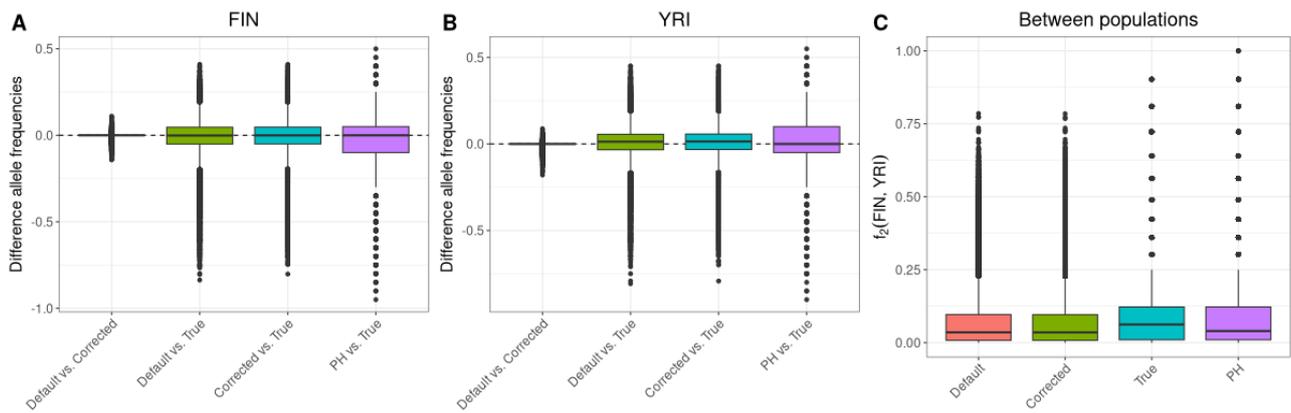


Figure S1: Differences in allele frequency estimates in the parts of the reference genome attributed to African ancestry. Boxplots for the differences between default genotype likelihood-based estimates and corrected genotype likelihood-based estimates, default genotype likelihood-based estimates and SNP array-based estimates, corrected genotype likelihood-based estimates, pseudohaploid (PH) genotype-based and SNP array-based estimates (A) in the FIN population and (B) in the YRI population. (C) is showing boxplots of the per-site population differentiation (measured as f_2 statistic) for the four allele frequency estimates.

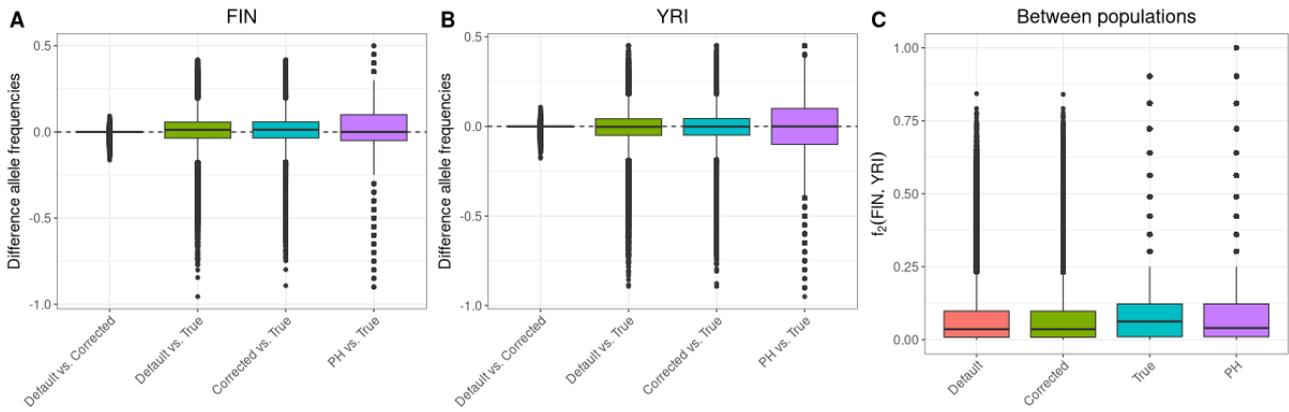


Figure S2: Differences in allele frequency estimates in the parts of the reference genome attributed to European ancestry. Boxplots for the differences between default genotype likelihood-based estimates and corrected genotype likelihood-based estimates, default genotype likelihood-based estimates and SNP array-based estimates, corrected genotype likelihood-based estimates, pseudohaploid (PH) genotype-based and SNP array-based estimates (A) in the FIN population and (B) in the YRI population. (C) is showing boxplots of the per-site population differentiation (measured as f_2 statistic) for the four allele frequency estimates.

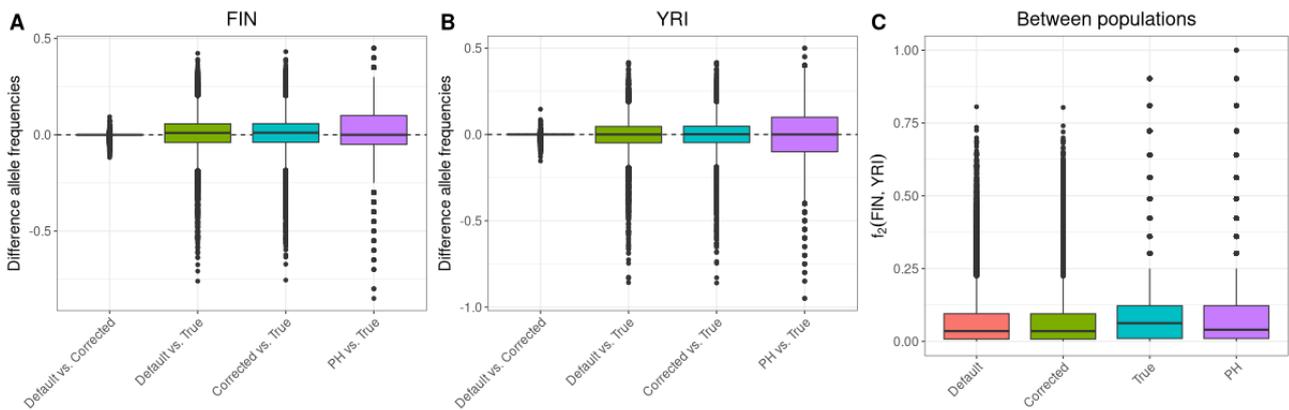


Figure S3: Differences in allele frequency estimates in the parts of the reference genome attributed to East Asian ancestry. Boxplots for the differences between default genotype likelihood-based estimates and corrected genotype likelihood-based estimates, default genotype likelihood-based estimates and SNP array-based estimates, corrected genotype likelihood-based estimates, pseudohaploid (PH) genotype-based and SNP array-based estimates (A) in the FIN population and (B) in the YRI population. (C) is showing boxplots of the per-site population differentiation (measured as f_2 statistic) for the four allele frequency estimates.

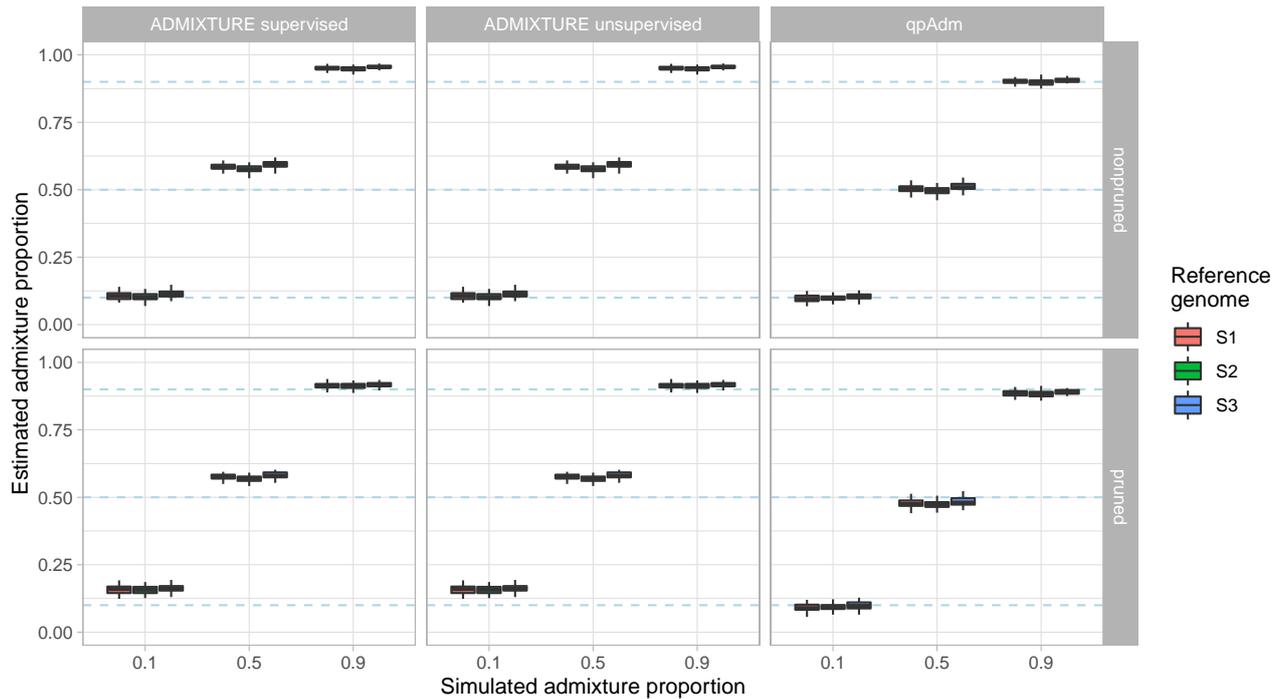


Figure S4: Simulation results for genotype call based methods using $t_{123} = 20000$ generations and a sequencing depth of 0.5X. Dashed blue lines represent the simulated admixture proportions.

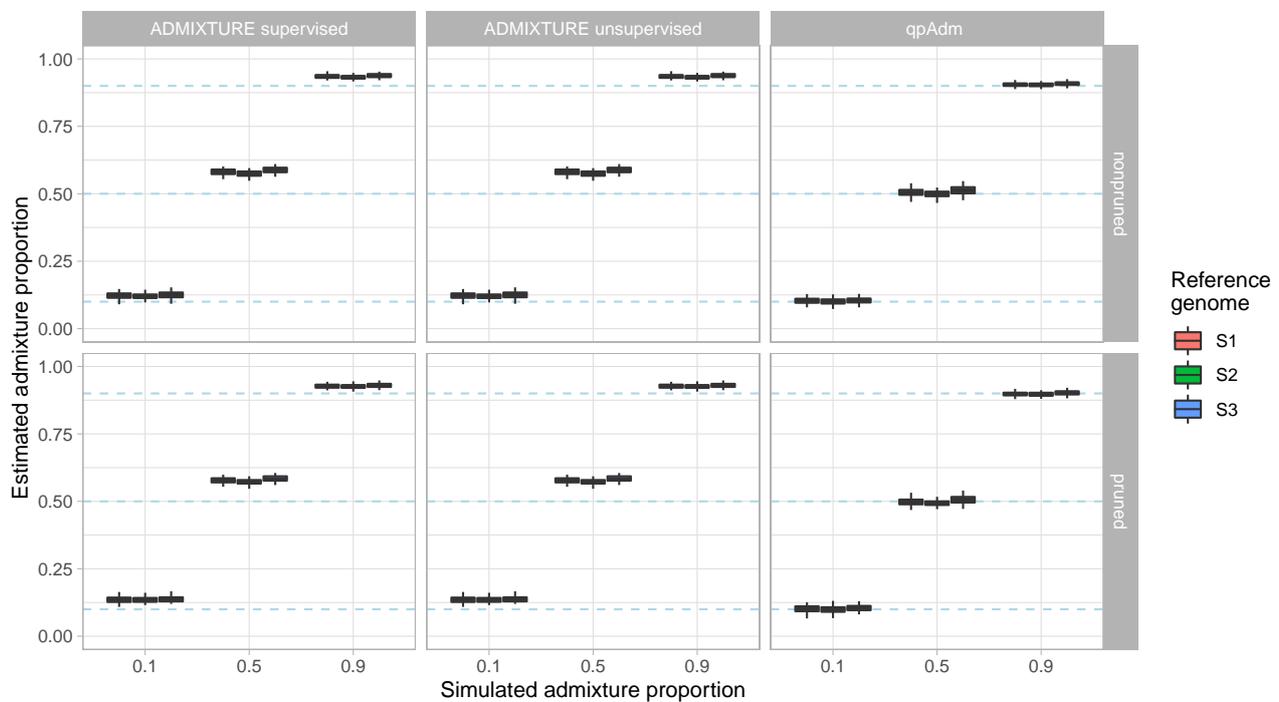


Figure S5: Simulation results for genotype call based methods using $t_{123} = 20000$ generations and a sequencing depth of 2.0X. Dashed blue lines represent the simulated admixture proportions.

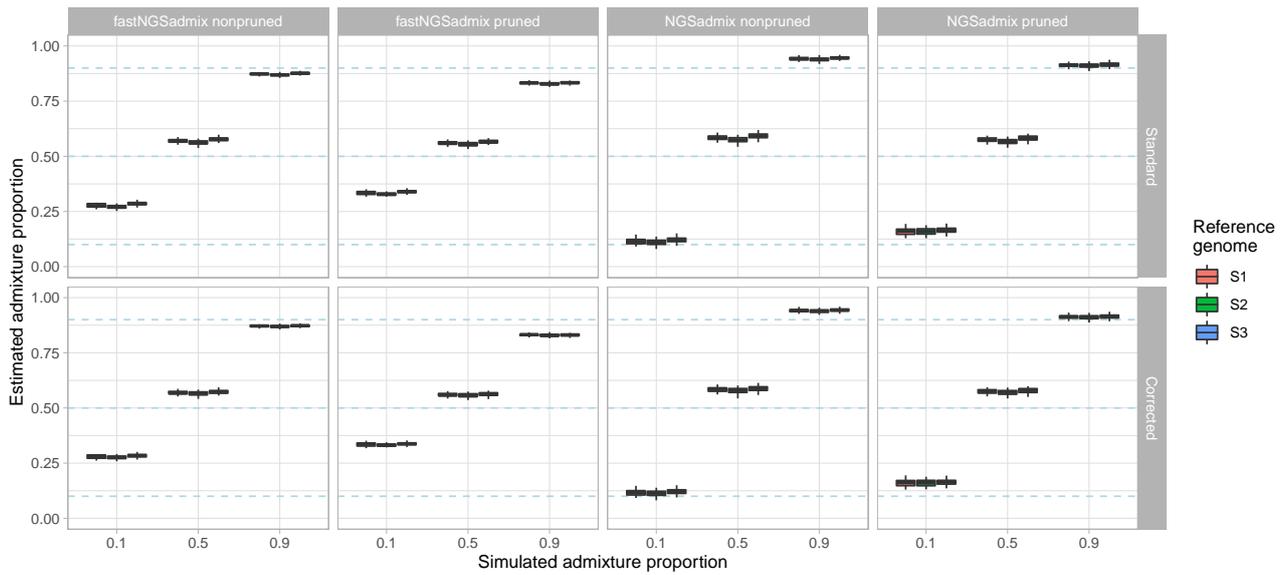


Figure S6: Simulation results for genotype likelihood based methods using $t_{123} = 20000$ generations and a sequencing depth of 0.5X. Dashed blue lines represent the simulated admixture proportions.

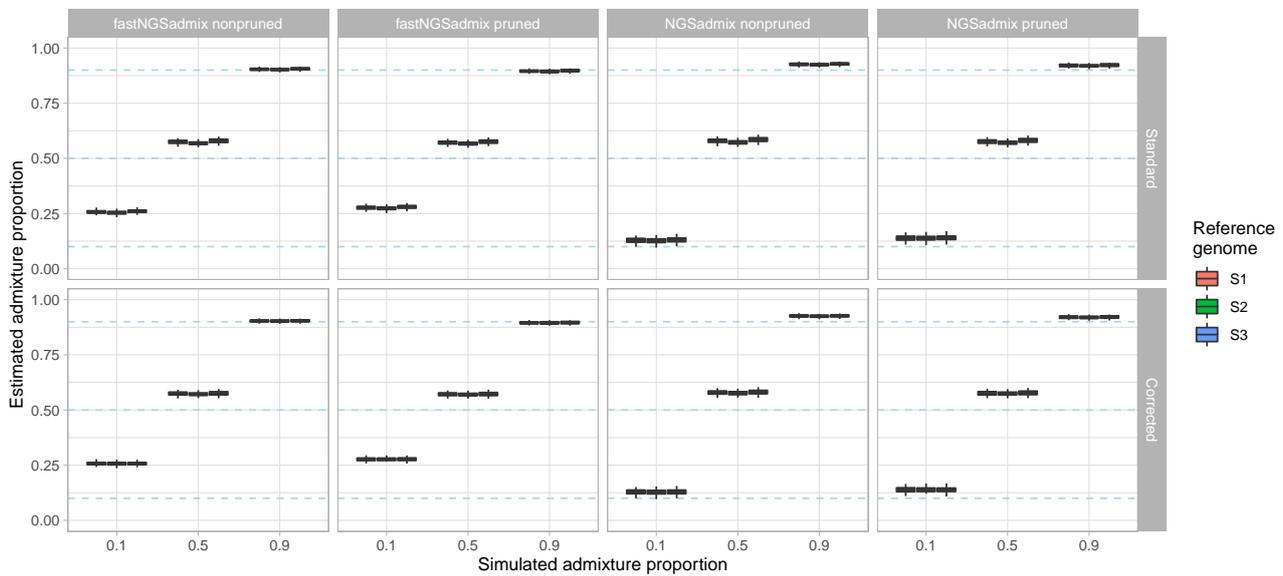


Figure S7: Simulation results for genotype likelihood based methods using $t_{123} = 20000$ generations and a sequencing depth of 2.0X. Dashed blue lines represent the simulated admixture proportions.

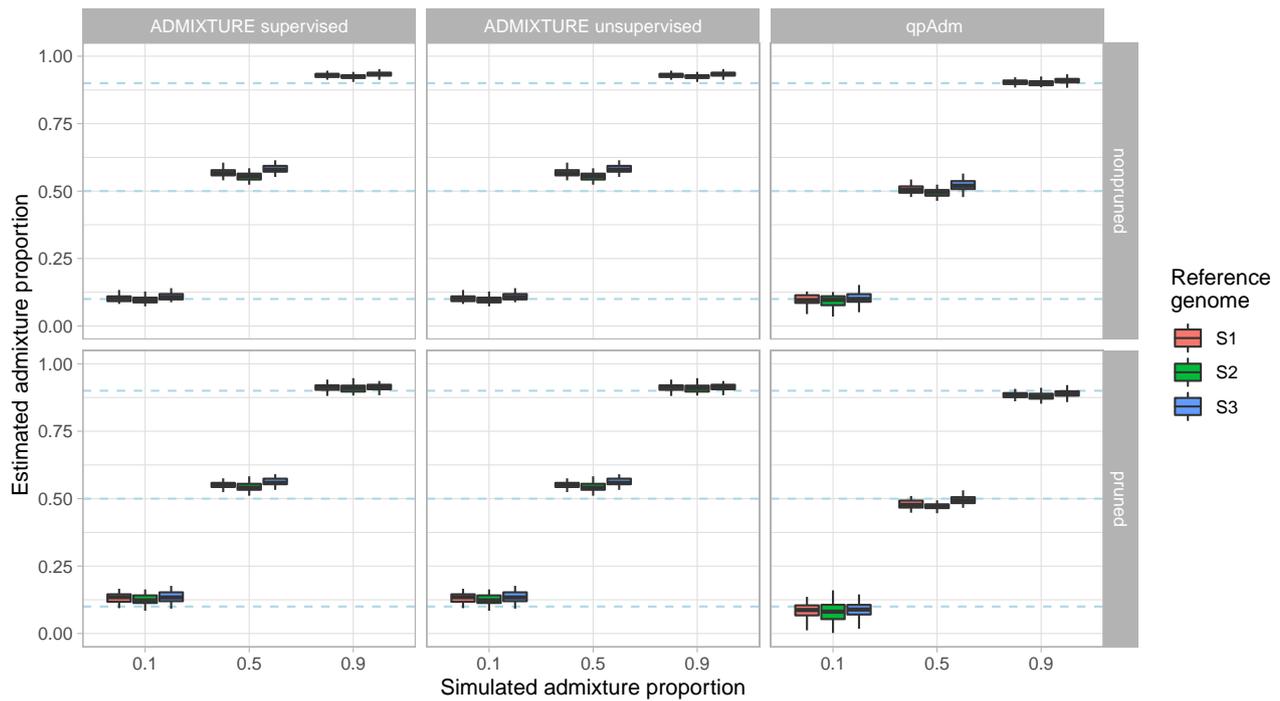


Figure S8: Simulation results for genotype call based methods using $t_{123} = 50000$ generations and a sequencing depth of 2.0X. Dashed blue lines represent the simulated admixture proportions.

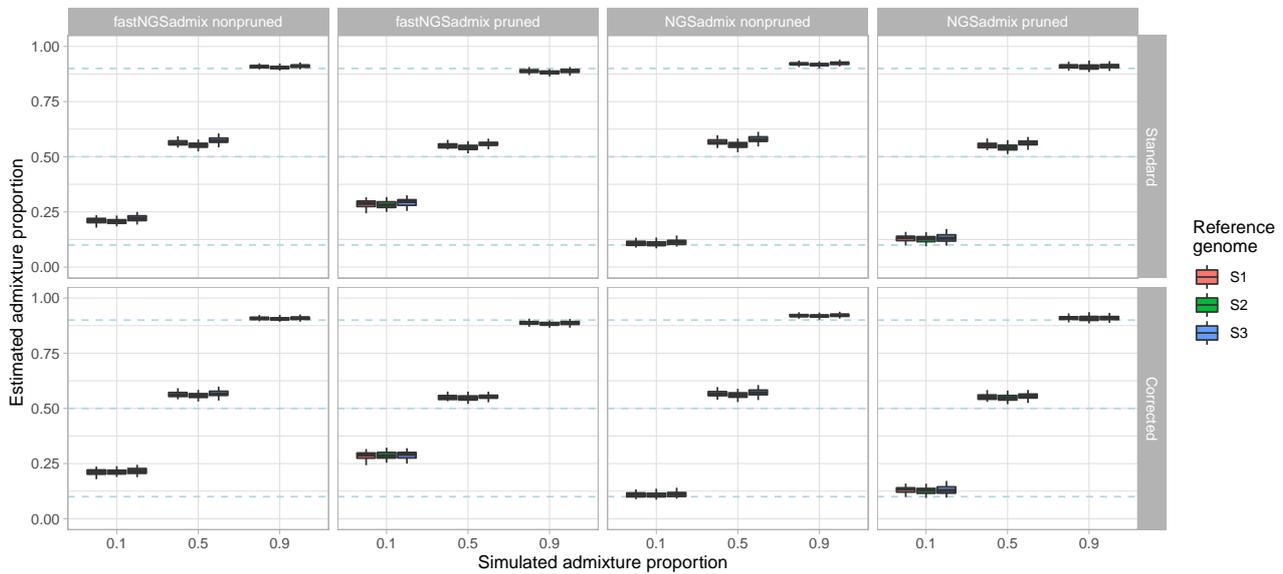


Figure S9: Simulation results for genotype likelihood based methods using $t_{123} = 50000$ generations and a sequencing depth of 2.0X. Dashed blue lines represent the simulated admixture proportions.

Supplementary Tables

Table S1: 1000 genomes individuals used for the analysis of empirical data.

individual	Population	Autosomal sequencing depth
HG00171	FIN	3.12803
HG00190	FIN	3.089
HG00272	FIN	3.61242
HG00277	FIN	3.86275
HG00284	FIN	4.08807
HG00323	FIN	2.80008
HG00330	FIN	13.9648
HG00380	FIN	3.45273
HG00177	FIN	3.43327
HG00189	FIN	3.48314
NA18853	YRI	2.56291
NA18923	YRI	4.42742
NA19197	YRI	4.19443
NA19200	YRI	4.22902
NA19236	YRI	4.21535
NA19248	YRI	4.24979
NA19116	YRI	3.03829
NA19130	YRI	4.97799
NA18520	YRI	3.99207
NA18522	YRI	2.55368

Table S2: Total number and percentage of SNPs with extreme differences ($\geq |0.5|$) between "True" and estimated allele frequencies.

Population	True vs default GL	True vs. corrected GL	True vs. Pseudohaploid
FIN	738 (0.118%)	608 (0.096%)	979 (0.157%)
YRI	829 (0.133%)	674 (0.108%)	947 (0.152%)