

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

Torsten Günther^{1,2,*}, [Amy Goldberg](#)³ & Joshua G. Schraiber^{3,4}

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

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

³[Department of Evolutionary Anthropology, Duke University, USA](#)

⁴[Department 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~~ [composition](#) 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 Using data from the 1000 Genomes Project, we find that our new method improves allele frequency estimation. To test a downstream application, we simulate](#) ancient DNA data with realistic post-mortem damage ~~;-we-to~~ 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 [more](#) 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 re-
- 2 sequencing studies employing short sequencing reads ([Orlando et al., 2013; Gopalakrishnan et al., 2017; Günther et al., 2019](#);
- 3 [Orlando et al., 2013; Gopalakrishnan et al., 2017; Günther and Nettelblad, 2019; Martiniano et al., 2020; Chen et al., 2020](#)).

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

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

47 To render genotype likelihoods and their downstream applications more robust to the presence of
48 mapping bias, we introduce a modified genotype likelihood, building off of the approach in Günther
49 and Nettelblad (2019). We ~~use modified reads carrying the other allele modify reads to carry both~~
50 ~~alleles~~ at biallelic SNP positions to assess the distribution of mapping bias and to obtain an empirical
51 quantification of the locus- and individual-specific mapping bias. We then calculate a modified geno-
52 type likelihood to account for mapping bias. The approach is similar to snpAD (Prüfer, 2018), with the
53 contrast that we are using a set of pre-ascertained biallelic SNPs because our aim is not to call geno-
54 types ~~all sites and we are using a set of~~ at all sites across the genome including potentially novel SNPs.

55 Restricting to known biallelic SNPs is a common practice in the population genomic analysis of ancient
56 DNA data as low-coverage and post-mortem damage usually limit the possibility of calling novel SNPs
57 for most individuals (see e.g. Günther and Jakobsson, 2019), and methods like `snpAD` are restricted to
58 very few high quality, high coverage individuals (Prüfer, 2018). Instead, most studies resort to using
59 pseudohaploid calls or genotype likelihoods at known variant sites (Günther and Jakobsson, 2019);
60 using ascertained biallelic SNPs ~~allowing~~ is particularly relevant when ancient DNA is enriched using
61 a SNP capture array (Rohland et al., 2022). This choice also allows us to estimate mapping bias
62 locus-specific rather than using one estimate across the full genome of the particular individual.

63 We examine two downstream applications of genetic data to determine the impact of mapping bias,
64 and assess the ability of our corrected genotype likelihood to ameliorate issues with mapping bias.
65 First, we look at a very high-level summary of genetic variation: allele frequencies. Because allele
66 frequencies can be estimated from high-quality SNP array data, we can use them as a control and
67 assess the impact of mapping bias and our corrected genotype likelihood in real short-read data.

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

87 2 Materials and Methods

88 2.1 Correcting genotype-likelihoods for mapping bias

89 Two versions of genotype likelihoods (Nielsen et al., 2011) were calculated for this study. First, we
90 use the direct method as included in the original version of GATK (McKenna et al., 2010) and also
91 implemented in ANGSD (Korneliussen et al., 2014). For a position ℓ covered by n reads, the genotype
92 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
93 genotype is A_1A_2 :

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

94 with

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

95 where $e_{\ell i}$ is the probability of a sequencing error of read i at position ℓ , calculated from the phred scaled
96 base quality score $Q_{\ell i}$, i.e. $e_{\ell i} = 10^{-Q_{\ell i}/10}$. The calculation of genotype likelihoods was implemented

97 in Python 3 using the pysam library (<https://github.com/pysam-developers/pysam>), a wrapper
98 around htslib and the samtools package (Li et al., 2009), or by calling samtools mpileup and parsing
99 the output in the Python script. Both corrected and default genotype likelihoods are calculated by
100 the same Python script.

101 To quantify the impact of mapping bias, we restrict the following analysis to a list of pre-defined
102 ascertained biallelic SNPs (list provided by the user) and modify each original read to carry the
103 other allele at the SNP position, as in Günther and Nettelblad (2019). The modified reads are then
104 remapped to the reference genome using the same mapping parameters. If there were no mapping
105 bias, all modified reads would map to the same position as the unmodified original read. Consequently,
106 when counting both original and modified reads together, we should observe half of our reads carrying
107 the reference allele and the other half carrying the alternative allele at the SNP position. We can
108 summarize the read balance at position ℓ as r_ℓ , which measures the proportion of reference alleles
109 among all original and modified reads mapping to the position. Without mapping bias, we would
110 observe $r_\ell = 0.5$. Under reference bias, we would observe $r_\ell > 0.5$ and under alternative bias $r_\ell < 0.5$.
111 We can see r_ℓ as an empirical quantification of the locus- and individual-specific mapping bias. Similar
112 to Prüfer (2018), we can then modify equation Equation 1 for heterozygous sites to

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

113 where R_ℓ is the reference allele at position ℓ and A_ℓ is the alternative allele. Note that when $r_\ell \equiv \frac{1}{2}$,
114 this recovers Equation 1. Genotype likelihood-based methods are tested with both genotype likelihood
115 versions. All code used in this study can be found under https://github.com/tgue/refbias_GL

116 2.2 Empirical Data

117 To estimate the effect of mapping bias in empirical data we obtained low coverage BAM files for ten
118 ~~FIN individuals and 10 YRI~~ (Finnish in Finland) individuals, ten JPT individuals (Japanese in Tokyo,
119 Japan) and ten YRI (Yoruba in Ibadan, Nigeria) individuals from the 1000 Genomes project (mostly
120 2-4x coverage; Table S1) (Auton et al., 2015). We also downloaded Illumina Omni2.5M chip genotype
121 calls for the same individuals ([http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/supporting/hd_genotype_chip/ALL.chip.omni_broad_sanger_combined.20140818.snps.genotypes.vcf.gz)
122 [supporting/hd_genotype_chip/ALL.chip.omni_broad_sanger_combined.20140818.snps.genotypes.](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/supporting/hd_genotype_chip/ALL.chip.omni_broad_sanger_combined.20140818.snps.genotypes.vcf.gz)
123 [vcf.gz](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/supporting/hd_genotype_chip/ALL.chip.omni_broad_sanger_combined.20140818.snps.genotypes.vcf.gz)). The SNP data was filtered to restrict to sites without missing data in the ~~20-30~~
124 individuals, a minor allele frequency of at least 0.2 in the reduced dataset (considering individuals from
125 ~~both all~~ populations together), ~~and excluding which makes it more likely that the SNPs are common~~
126 in all populations and both over- and underestimation of allele frequencies could be observed. We also
127 excluded A/T and C/G SNPs to avoid strand misidentification. Reads mapping to these positions
128 were extracted from the BAM files using samtools (Li et al., 2009). To make the sequence data
129 more similar to fragmented ancient DNA, each read was split into two halves at its mid-point and
130 each sub-read was re-mapped separately. For mapping, we used bwa aln (Li and Durbin, 2009) and
131 the non-default parameters -l 16500 (to avoid seeding), -n 0.01 and -o 2. Only reads with mapping
132 qualities of 30 or higher were kept for further analysis.

133 Pseudohaploid genotypes were called with ANGSD v0.933 (Korneliussen et al., 2014) by randomly
134 drawing one read per SNP ~~as described for the simulations below and only with a minimum base~~
135 quality of 30. This step was performed using ANGSD with the parameters -checkBamHeaders 0 (to
136 deactivate checking the headers of the BAM files) -doHaploCall 1 (to sample a single base only)
137 -doCounts 1 (needed to determine the most common base) -doGeno -4 (to format genotypes as bases
138 not integers in the output) -doPost 2 (estimate the posterior genotype probability assuming a uniform
139 prior, output files not used) -doPlink 2 (produce output in tfam/tped format) -minMapQ 30 (to set
140 the minimum mapping quality) -minQ 30 (to set the minimum base quality) -doMajorMinor 1 (to
141 infer major and minor from genotype likelihoods) -GL 2 (to calculate GATK genotype likelihood,
142 output files not used) -domaf 1 (calculate allele frequencies with fixed major and minor alleles). This

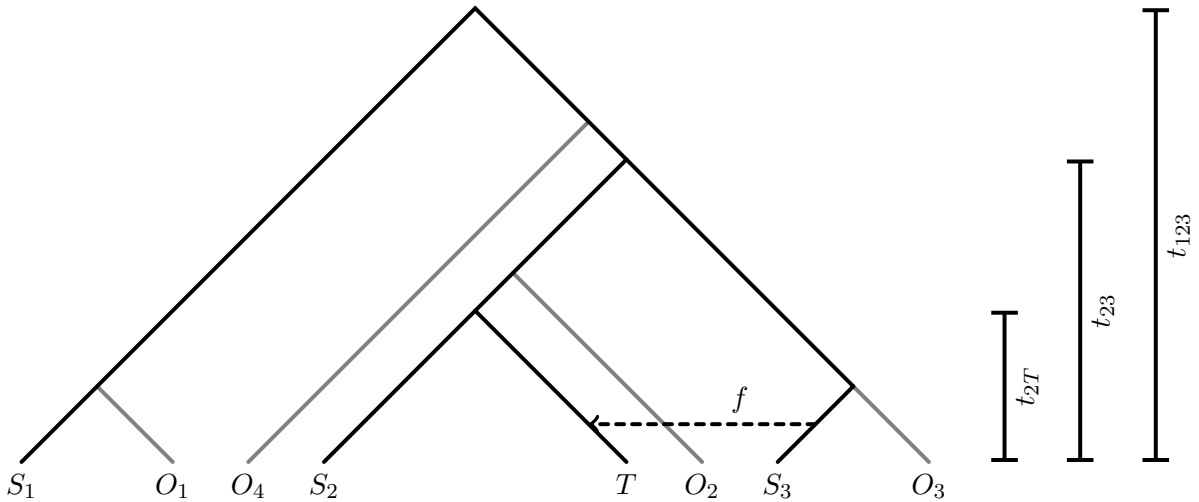


Figure 1: Illustration of the population relationships used in the simulations. Branch lengths are not to scale

143 ~~call also calculates genotype likelihoods in ANGSD but we used both default and corrected likelihoods~~
 144 ~~calculated from our own Python script to ensure consistency. Haplocall files were then converted to~~
 145 ~~Plink format using haploToPlink distributed with ANGSD (Korneliussen et al., 2014). Only SNPs~~
 146 ~~with the same two alleles in pseudohaploid and SNP chip data were included in all comparisons.~~
 147 ~~Remapping of modified reads and genotype likelihood calculation were performed as described above.~~
 148 ~~Allele frequencies were calculated from genotype likelihoods with ANGSD v0.933 (Korneliussen et al.,~~
 149 ~~2014) using -doMaf 4 and the human reference as “ancestral” allele (-anc) in order to calculate the~~
 150 ~~allele frequency of the reference alleles. SNP calls from the genotyping array and pseudohaploid calls~~
 151 ~~were converted to genotype likelihood files assuming no genotyping errors, ~~so the allele frequency~~~~
 152 ~~estimation (i.e. the genotype likelihood of the observed genotype was set to 1.0, others to 0.0 whereas~~
 153 ~~all three likelihoods were set to $\frac{1}{3}$ if data was missing for the site and individual). This allowed us to~~
 154 ~~also estimate allele frequency estimates for this data ~~could be based on with~~ ANGSD as well.~~

2.3 Simulation of genomic data

155
 156 ~~Population histories are~~ To test the methods while having control over the “true” admixture proportions,
 157 ~~population histories were~~ simulated using msprime v0.6.2 (Kelleher et al., 2016). We ~~simulate simulated~~
 158 a demographic history where a target population T receives a single pulse of admixture with propor-
 159 tion f from source $S3$ 50 generations ago. Furthermore, we ~~simulate simulated~~ population $S1$ which
 160 forms an outgroup and population $S2$ which is closer to T than $S3$ to serve as second source for
 161 estimating ancestry proportions (Figure 1). Finally, we ~~simulate simulated~~ populations $O1$, $O2$, $O3$,
 162 and $O4$ as populations not involved in the admixture events which split off internal branches of the
 163 tree to serve as “right” populations for qpAdm (Haak et al., 2015; Harney et al., 2021). Split times ~~are~~
 164 ~~were~~ scaled relative to the deepest split t_{123} : the split between $(S2, T)$ and $S3$, t_{23} , is set to $0.5 \times t_{123}$
 165 while the split between T and $S2$ ~~is was~~ set to $0.2 \times t_{123}$. ~~Different values~~ To set t_{123} , we considered
 166 a value of 20,000 generations, approximately falling in the range of the split of all human populations
 167 (Schlebusch et al., 2017) or the Neanderthal-Denisovan split (Rogers et al., 2017) i.e. approximating
 168 the divergence between distant populations or sub-species, and 50,000 generations ~~are tested for~~ t_{123}
 169 ~~approximately corresponding to divergence times within and between (sub-), corresponding to a~~
 170 ~~comparison between closely related~~ species. Mutation rate was set to 2.5×10^{-8} and recombination
 171 rate was set to 2×10^{-8} , ~~which are both in the upper part of the ranges for mammals and vertebrates~~
 172 ~~(Dumont and Payseur, 2008; Bergeron et al., 2023). The effective population size along all branches~~
 173 ~~is was~~ 10,000,000, a value often considered for humans (Charlesworth, 2009). For each popula-
 174 tion, 21 diploid individuals (i.e. 42 haploid chromosomes) with 5 chromosome pairs of 20,000,000 bp

175 (corresponding to a short mammalian chromosome arm, Lander et al. (2001)) each were simulated.
176 As msprime does not produce sequences but positions of derived alleles at each haploid chromosome,
177 we had to convert this information into a sequence. For each chromosome, a random ancestral sequence
178 was generated with a GC content of 41% corresponding to the GC content of the human genome
179 (Lander et al., 2001). Transversion polymorphisms were then placed along the sequence ~~according to~~
180 at the positions produced by the msprime simulations. The ~~first resulting sequences for each haploid~~
181 chromosome were then stored as FASTA files. One of the 42 simulated sequences from populations
182 *S1*, *S2* and *S3* were used as reference genomes. ~~Pairs of sequences~~ Out of the remaining sequences,
183 pairs of FASTA files were then considered as diploid individuals and used as input for gargamel
184 (Renaud et al., 2017) ~~was used to simulate to serve as endogenous sequences for the simulation of~~
185 next-generation sequencing data with ancient DNA damage. Data were simulated to mimic data generated
186 with an Illumina HiSeq 2500 sequencing machine assuming the post-mortem damage pattern observed
187 when sequencing Neandertals in Briggs et al. (2007). We simulated coverages of 0.5X and 2.0X. For
188 each individual, fragment sizes followed a log-normal distribution with a location between 3.3 and 3.8
189 (randomly drawn per individual from a uniform distribution) and a scale of 0.2, corresponding to an
190 average fragment length per individual between 27 and ~~46bp~~ 46 bp. Fragments shorter than ~~20bp~~ 30
191 bp were excluded. No contaminating sequences were simulated. Sequencing reads were then trimmed
192 and merged with AdapterRemoval (Schubert et al., 2016). ~~Reads~~ All reads (merged and the small
193 proportion of unmerged) were then mapped to the different reference genomes using `bwa aln v0.7.17`
194 (Li and Durbin, 2009) together with the commonly used non-default parameters `-l 16500` (to avoid
195 seeding), `-n 0.01` and `-o 2` (to allow for more mismatches and gaps due to post-mortem damages and
196 increased evolutionary distance to the reference) (Schubert et al., 2012; Oliva et al., 2021). BAM files
197 were handled using `samtools v1.5` (Li et al., 2009).

198 ~~Genotype calling and downstream analysis were performed separately for the three reference genomes~~
199 ~~originating from populations S1, S2 and S3. To avoid ascertainment bias, polymorphic SNPs were~~
200 ~~ascertained~~ To ascertain SNPs, we avoided the effect of damage, sequencing errors and genotype
201 callers, by identifying biallelic SNPs directly from the simulated ~~true genotypes~~ genotypes, prior to
202 the gargamel simulation of reads and mapping, and restricted to SNPs with a minimum allele fre-
203 quency of 10% in the outgroup population *S1*. This mimics an ascertainment procedure in which
204 SNPs are ascertained in an outgroup population, which may be common in many taxa. 100,000
205 SNPs were selected at random using `Plink v1.90` (Chang et al., 2015) `-thin-count`. Genotype calling
206 and downstream analysis were performed separately for the three reference genomes originating from
207 populations S1, S2 and S3. Pseudohaploid calls were then generated for all individuals at these sites
208 using `ANGSD v0.917` (Korneliussen et al., 2014) by randomly sampling a single read per position with
209 minimum base and mapping quality of at least 30. This step was performed using `ANGSD` with the pa-
210 rameters ~~`-checkBamHeaders 0 -doHaploCall 1 -doCounts 1 -doGeno 4 -doPost 2 -doPlink 2 -minMapQ`~~
211 ~~`30 -minQ 30 -doMajorMinor 1 -GL 1 -domaf 1.`~~ Files as described for the empirical data above and
212 files were then converted to `Plink` format using `haploToPlink` distributed with `ANGSD` (Korneliussen
213 et al., 2014). For downstream analyses, the set of SNPs was further restricted to sites with less than
214 50 % missing data and a minor allele frequency of at least 10% in *S1*, *S2*, *S3* and *T* together. Binary
215 and transposed `Plink` files were handled using `Plink v1.90` (Chang et al., 2015). `convertf` (Patterson
216 et al., 2006; Price et al., 2006) was used to convert between `Plink` and `EIGENSTRAT` file formats. `Plink`
217 was also used for linkage disequilibrium (LD) pruning with parameters `-indep-pairwise 200 25 0.7`.

2.4 Estimating admixture proportions

218
219 We used ~~five~~ four different approaches to estimate ancestry proportions in our target population *T*.
220 In addition to differences in the underlying model and ~~implementations, for users~~ implementation, the
221 tools differ in the type of their input data (genotype calls or genotype likelihoods) and whether their
222 approaches are unsupervised and/or supervised (Table 1).

223 All software was set to estimate ancestry assuming two source populations. Unless stated otherwise,
224 *S2* and *S3* were set as sources and *T* as the target population while no other individuals were included

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

225 in when running the software. ADMIXTURE (Alexander et al., 2009; Alexander and Lange, 2011) is the
 226 only included method that has both a supervised (i.e. with pre-defined source populations) and an
 227 unsupervised mode. Both options were tested using the `-haploid` option without multithreading as the
 228 genotype calls were pseudo-haploid. For qpAdm (Haak et al., 2015; Harney et al., 2021), populations
 229 *O1*, *O2*, *O3* and *O4* served as “right” populations. qpAdm was run with the options `allsnps: YES` and
 230 `details: YES`. For fastNGSadmix (Jørsboe et al., 2017), allele frequencies in the source populations
 231 were estimated using NGSadmix (Skotte et al., 2013) with the option `-printInfo 1`. fastNGSadmix
 232 was then run to estimate ancestry per individual without bootstrapping. NGSadmix (Skotte et al.,
 233 2013) was run in default setting. The mean ancestry proportions across all individuals in the target
 234 population was used as an ancestry estimate for the entire population. In the case of unsupervised
 235 approaches, the clusters belonging to the source populations were identified as those where individuals
 236 from *S2* or *S3* showed more than 90 % estimated ancestry.

237 3 Results

238 3.1 Mapping Impact of mapping bias on allele frequency estimates in empirical data

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

247 We first tested the effect of mapping bias on allele frequency estimates in empirical data. We selected
 248 low to medium coverage (mostly between 2 and 4x depth 2-4x coverage, except for one individual at
 249 14X14x, Table S1) for ten individuals from each of two-three 1000 Genomes populations (FIN, JPT
 250 and YRI) from different continents. All individuals show an empirical bias towards the reference allele
 251 as indicated by average $r_L > 0.5$ (Tables S1 and S2). We used ANGSD to estimate allele frequencies
 252 from genotype likelihoods based on short-read NGS data (read lengths reduced to 36-54 bp to better
 253 resemble fragmented aDNA data) and compare them to allele frequencies estimated from the same
 254 individuals genotyped using a SNP array and pseudohaploid genotype data. As the genotyping array
 255 does not involve a mapping step to a reference genome it should be less affected by mapping bias, we
 256 consider these estimates as “true” allele frequencies.

257 Overall, genotype likelihood-based point estimates of the allele frequencies tend towards more inter-
 258 mediate allele frequencies while pseudohaploid genotypes and “true” genotypes result in more alleles es-
 259 timated to have low and high alternative allele frequency (Figure 2A and B). In FIN, the pseudohaploid
 260 genotypes lead to a slight underestimation of the reference allele frequencies (Figure 2A), while this
 261 signal is reversed in YRI (Figure 2B), a pattern which could be related to the fact that most of the

Table 2: Pearson’s correlation coefficients comparing different allele frequency estimates in the three empirical populations. 95% confidence intervals are shown in parentheses.

<u>Population</u>	<u>True vs</u>
<u>FIN</u>	<u>0.8460 [0.9294, 0.9301]</u>
<u>YRI</u>	<u>0.8246 [0.9457, 0.9462]</u>
<u>JPT</u>	<u>0.8466 [0.8238, 0.8254]</u>

human reference genome has European ancestry (Green et al., 2010; Church et al., 2015; Günther and Nettelblad, 2019). In both S1). In all tested populations, the default version of genotype likelihood calculation produced an allele frequency distribution slightly shifted towards lower non-reference allele frequency estimates compared to the corrected genotype likelihood (Paired Wilcoxon test $p < 2.2 \times 10^{-22}$ in both all populations). The Consistently, the per-site allele frequencies estimated from the corrected genotype likelihoods exhibit a slightly better correlation with the “true” frequencies in both FIN (Pearson’s correlation coefficient 0.9297 (Table 2). Allele frequency estimates from pseudohaploid data display the best correlation with the “true” frequencies in all populations (Table 2).

Overall, the per-site differences between “true” frequencies in both FIN ($r = 0.8571$) and YRI ($r = 0.8344$) indicating that while the distribution of allele frequencies seems close to the true spectrum (Figure 2A and B), the estimates at individual loci are rather noisy.

Differences at individual sites, however, display some extreme outliers with $\sim 0.1\%$ of the SNPs showing more than 50% difference between estimates from SNP chips and sequencing data, which could hint at systematic technological differences between the two data types at those sites. This pattern of outliers is slightly less pronounced when using the corrected genotype likelihoods (Table ??) allele frequencies and all frequencies estimated from NGS data (genotype-likelihoods and pseudohaploid) show a trend towards lower estimated non-reference alleles in the NGS data (Figure 2A-C), suggesting an impact of mapping bias. Outliers even reach a difference of up to -1.0. Interestingly, despite the overall closer concordance between the pseudohaploid allele frequency spectrum and the SNP array allele frequency spectrum, there is significantly higher variation between pseudohaploid and true frequencies at any particular per-site (Figure 2A-C), suggesting that this is a general difference between NGS and SNP chip data. In Günther and Nettelblad (2019), we found that different parts of the human reference genome exhibit different types of mapping bias. We find a similar result here: the parts of the reference genome that can be attributed to African ancestry (Green et al., 2010) display a mean and median difference of nearly 0 in FIN but allele frequencies remain higher than array estimates in YRI (Figure S2). In contrast, the European and East Asian parts of the reference genome show a distribution of differences around 0 in YRI but positive means and median in FIN (Figures S3 and S4). This confirms the utility of reducing the effect of mapping bias by mapping against a reference genome from an outgroup allele frequency estimates from pseudohaploid calls are relatively noisy but also relatively unbiased. A consequence of the systematic over-estimation of the allele frequencies when using genotype likelihoods is that the population differentiation (here measured as f_2 statistic) is reduced compared to estimates from the SNP array or pseudohaploid genotype calls (Figure 2E-D-F). In Günther and Nettelblad (2019), we found that different parts of the human reference genome exhibit different types of mapping bias in the estimation of archaic ancestry which could be attributed to the fact that the human reference genome is a mosaic of different ancestries (Green et al., 2010; Church et al., 2015). Here, we do not find substantial differences in the allele frequency patterns between the different continental ancestries (Figures S2-S4).

3.2 Estimation of admixture proportions based on genotype calls in simulated data

We compare the accuracy of the different methods for estimating admixture proportion under a set of different population divergence times, sequencing depths, and with or without LD pruning of the SNP panel. Mapping to three different reference genomes, one from an outgroup (S1) and the two

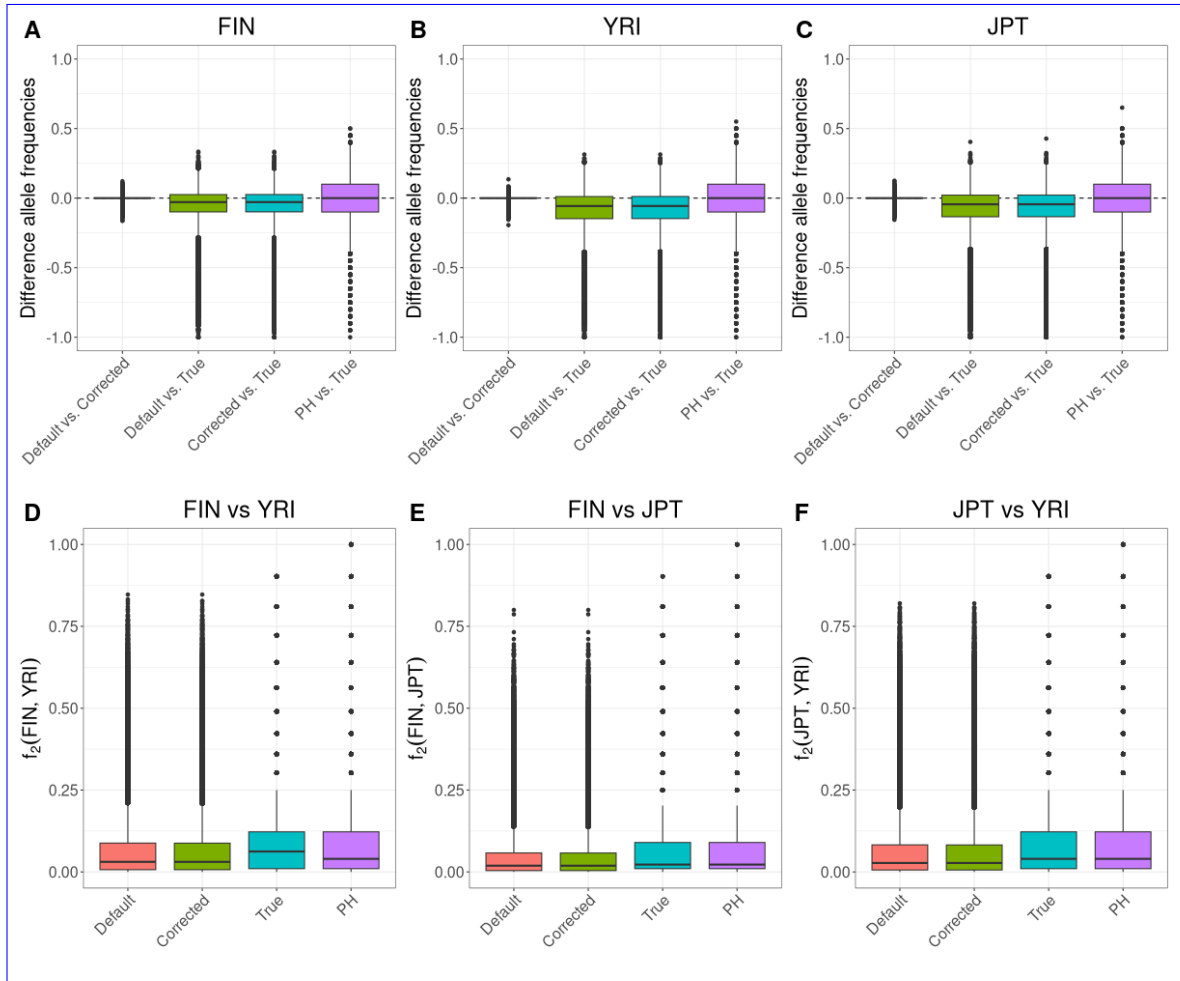


Figure 2: Differences in allele frequency estimates. Boxplots for the differences between default genotype likelihood-based estimates and corrected genotype likelihood-based estimates, default genotype likelihood-based estimates and SNP array-based estimates, corrected genotype likelihood-based estimates, pseudohaploid (PH) genotype-based and SNP array-based estimates (A) in the FIN population, (B) in the YRI population and (C) in the JPT population. (D-F) are showing boxplots of the pairwise per-site population differentiation (measured as f_2 statistic) for the four allele frequency estimates.

ingroups also representing the sources of the admixture event ($S2$ and $S3$), allows us to use $S1$ as a control which should not be affected by mapping bias and only other aspects of the data. We expect that mapping reads to one of the sources will cause a preference for reads carrying alleles from that population at heterozygous sites and, consequently, an overestimation of the ancestry proportion attributed to that population. The distance between the estimates when mapped to $S2$ or $S3$ (and their distances to the results when using $S1$) can then be seen as an estimate of the extent of mapping bias.

For most parts of this results section, we will focus on the scenario with an average sequencing depth of 0.5X where the deepest population split (t_{123}) was 50,000 generations ago and the split (t_{23}) between the relevant sources dating to 25,000 generations ago. Consequently, mapping the reads against a reference genome sequence from one or the other source would be equivalent to a study comparing (sub-)species where the reference genome originated from one of those populations. Results for other population divergences and sequencing depths are shown in Figures S5-S10.

We begin by assessing methods that require hard genotype calls, ADMIXTURE and qpAdm. For these approaches, we used single randomly drawn reads per individual and site to generate pseudo-haploid data in the target population. The popular implementation of the PSD (Pritchard et al., 2000) model working with SNP genotype calls, ADMIXTURE (Alexander et al., 2009; Alexander and Lange, 2011), has both supervised and unsupervised modes. Both modes show similar general patterns: low (10%) admixture proportions are estimated well while medium to high ($\geq 50\%$) admixture proportions are over-estimated (Figure 3). On the full SNP panel, the median estimated admixture proportion differs up to $\sim 4\%$ when mapping to reference genomes representing either of the two sources ($S2$ or $S3$) while mapping to the outgroup reference genome ($S1$) results in estimates intermediate between the two (Data S1). LD pruning slightly reduces mapping bias and reduces the overestimation, at least for high (90%) admixture proportions. qpAdm (Haak et al., 2015; Harney et al., 2021), on the other hand, estimated all admixture proportions accurately when the outgroup ($S1$) was used for the reference genome sequence and when the full SNP panel was used. The median estimates of admixture differed up to 3% between mapping to reference genomes from one of the source populations ($S2$ or $S3$). Notably, LD pruning increased the noise of the qpAdm estimates (probably due to the reduced number of SNPs) and led to all admixture proportions being slightly underestimated (Figure 3). The extent of mapping bias decreases with lower population divergence between the sources across all methods (Figure S5), as mapping bias should correlate with distance to the reference genome sequence. Conversely, increasing sequencing depth mostly reduced noise but not mapping bias (Figures S6 and S9) as the genotype-based methods benefit from the increased number of SNPs but the genotype calls do not increase certainty when multiple reads are mapping to the same position.

3.3 Estimation of admixture proportions based on genotype likelihoods in simulated data

We next examined the performance of genotype-likelihood-based approaches to estimate admixture proportions. In principle, genotype likelihoods should be able to make better use of all of the data in ancient DNA, because more than a single random read can be used per site. Moreover, we are able to explicitly incorporate our mapping bias correction into the genotype likelihood. We compared the supervised fastNGSadmix (Jørsboe et al., 2017) to the unsupervised NGSadmixmap (Skotte et al., 2013). fastNGSadmix shows the highest level of overestimation of low to medium admixture proportions ($\leq 50\%$) among all tested approaches while high admixture proportions (90%) are estimated well (Figure 4). Mapping bias caused differences of up to $\sim 3\%$ in the admixture estimates when mapping to the different reference genomes. LD pruning enhances the overestimation of low admixture proportions while leading to an underestimation of high admixture proportions (Data S1). Notably, when employing the corrected genotype-likelihood the estimated admixture proportions when mapping to $S2$ or $S3$ are slightly more similar than with the default formula without correction, showing that the correction makes the genome-wide estimates less dependent on the reference sequence used for mapping while not fully removing the effect. The estimates when using the outgroup $S1$ as reference are

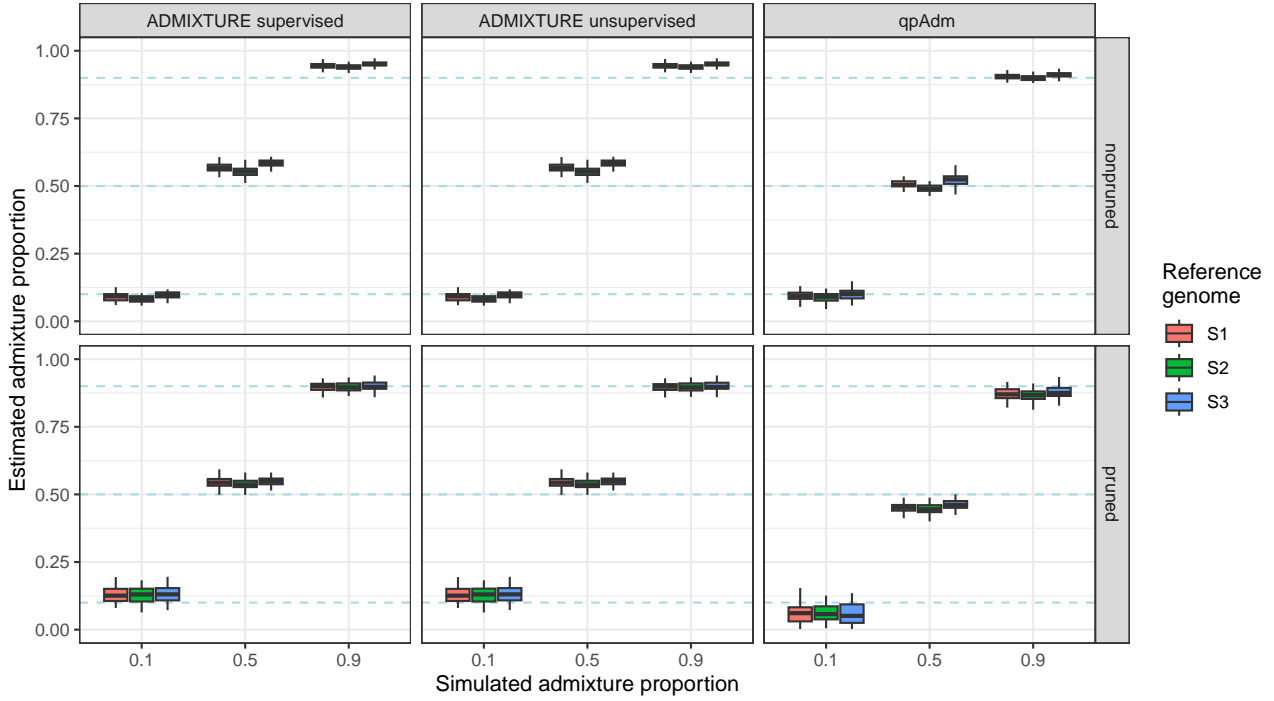


Figure 3: Simulation results for genotype call based methods using $t_{123} = 50000$ generations and a sequencing depth of 0.5X. Dashed blue lines represent the simulated admixture proportions.

353 slightly higher for high admixture proportions (90%). The results for *NGSadmix* show similar patterns
 354 to *ADMIXTURE* with a moderate overestimation of admixture proportions $\geq 50\%$ (Figure 4). Mapping
 355 bias caused differences of up to $\sim 4\%$ in the admixture estimates when mapping to the different
 356 reference genomes. After LD pruning, estimated admixture proportions for higher simulated values
 357 were closer to the simulated values. Furthermore, employing the mapping bias corrected genotype-
 358 likelihoods made the estimated admixture proportions less dependent on the reference genome used
 359 during mapping, particularly when using *NGSadmix* in pruned data, where all three reference genomes
 360 produce nearly identical results. Notably, the extent of over-estimation for both methods seems to
 361 be somewhat negatively correlated with population divergence (Figures S7 and 4), i.e. increased distan-
 362 ces between the source populations reduces the method bias. Further patterns are as expected:
 363 the extent of mapping bias is correlated with population divergence and increased sequencing depth
 364 reduces noise (Figures S7, 4, S8 and S10).

365 4 Discussion

366 We illustrate the impacts of mapping bias on downstream applications, such as allele frequency esti-
 367 mation and ancestry proportion estimation, and we introduced a new approach to recalibrate genotype
 368 likelihoods in the presence of mapping bias to alleviate its effects. The impact of mapping bias in
 369 our comparisons is small but pervasive suggesting that it can have an effect on the results of different
 370 types of analysis in empirical studies. In contrast to other approaches to alleviate mapping bias, such
 371 as employing pangenome variation graphs (Martiniano et al., 2020; Koptekin et al., 2023), it does not
 372 require establishing a separate pipeline. Instead, only reads mapping to a set of ascertained SNP
 373 positions need to be modified and remapped which only represents only a fraction of all reads and
 374 consequently will require a small proportion of the original mapping time. Our Python scripts used to
 375 calculate the genotype likelihoods could be optimized further, but this step is of minor computational
 376 costs compared to other parts of the general bioinformatic pipelines (~ 1 minute per individual in the
 377 empirical data analysis for this study) in ancient DNA research. The corrected genotype likelihoods

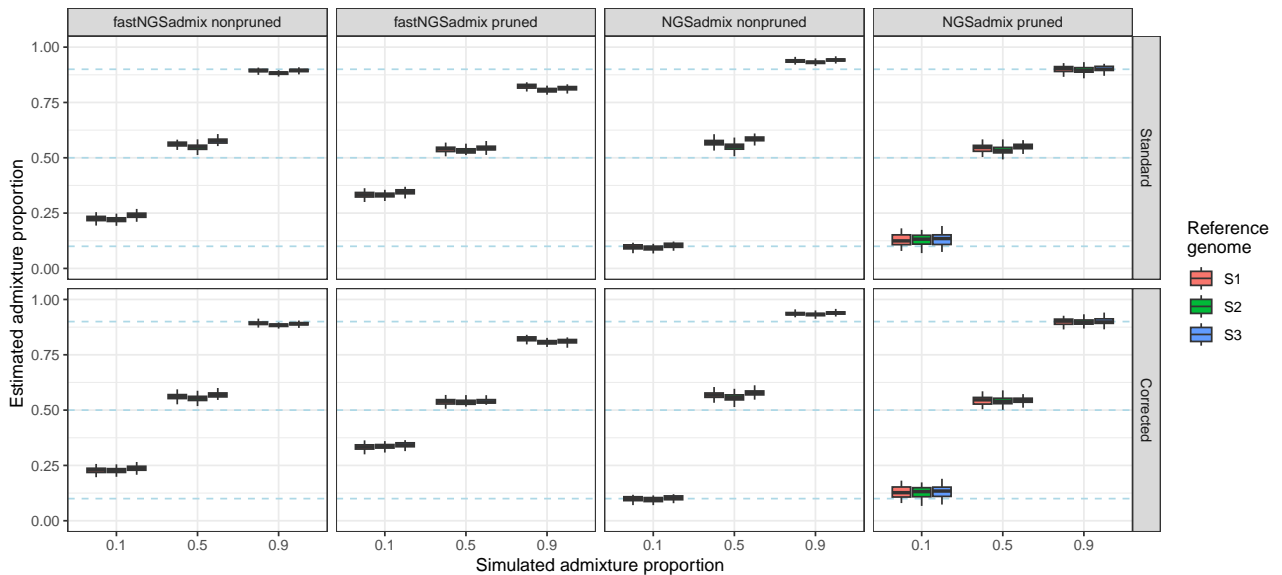


Figure 4: Simulation results for genotype likelihood based methods using $t_{123} = 50000$ generations and a sequencing depth of 0.5X. Dashed blue lines represent the simulated admixture proportions.

378 [can then be directly used in downstream analyses using the same file structures and formats as other](#)
 379 [genotype likelihood-based approaches.](#)

380 Increasing sample sizes in ancient DNA studies have motivated a number of studies aiming to detect
 381 selection in genome-wide scans or to investigate phenotypes in ancient populations (e.g. [Mathieson](#)
 382 [et al., 2015](#); [Cox et al., 2022](#); [Klunk et al., 2022](#); [Gopalakrishnan et al., 2022](#); [Mathieson and Terhorst,](#)
 383 [2022](#); [Davy et al., 2023](#); [Barton et al., 2023](#); [Hui et al., 2024](#)). Such investigations are potentially very
 384 sensitive to biases and uncertainties in genotype calls or allele frequencies at individual sites while
 385 certain effects will average out for genome-wide estimates such as ancestry proportions. Concerns
 386 about certain biases and how to estimate allele frequencies have even reduced confidence in the results
 387 of some studies [searching for loci under selection](#) ([Gopalakrishnan et al., 2022](#); [Barton et al., 2023](#)).
 388 Our results indicate that such concerns are valid as individual sites can show very strong deviations
 389 ~~in their allele frequencies when allele frequencies are~~ estimated from low-coverage sequencing data
 390 [\(Figure 2\)](#). This is due to a combination of effects, including mapping bias ~~and sampling artifacts.~~
 391 ~~Allele frequency point estimates from genotype likelihoods tend to be higher than true frequencies~~
 392 ~~because most alleles segregate at low frequencies, and thus appear most commonly in heterozygotes.~~
 393 [However, Without high coverage data,](#) genotype likelihood approaches without an allele frequency
 394 prior will naturally put some weight on ~~individuals being homozygous for the allele~~ [all three potential](#)
 395 [genotypes at a site](#), ultimately collectively driving ~~up allele frequency estimates~~ [allele frequency to more](#)
 396 [intermediate values](#). The risk is then that most downstream analyses will treat the allele frequency
 397 point estimates ~~as face values at face value,~~
 398 While our new approach to recalibrate genotype likelihoods reduces the number of outlier loci, there
 399 is still uncertainty in allele frequency estimates from low coverage data. Therefore, results heavily
 400 relying on allele frequency estimates or genotype calls at single loci from low-coverage sequencing data
 401 or even ancient DNA data need to be taken with a grain of salt.

402 The simulations in this study revealed a modest but ~~significant noticeable~~ effect of mapping bias on
 403 ancestry estimates as the difference between reference genomes never exceeded 5 percent. [In particular,](#)
 404 [we found that mapping bias and method bias even counteract each other in certain cases, leading to](#)
 405 [better estimates of the admixture proportion when mapping to one of the sources.](#) The differences
 406 seen in our simulations are likely underestimates of what might occur in empirical studies ~~as,~~ [because](#)
 407 real genomes are larger and more complex than what ~~was we~~ used in the simulations. For instance, we
 408 simulated five 20 megabase long chromosomes for a 100 megabase genome, while mammalian genomes

409 are one order of magnitude larger; the human genome is roughly 3 gigabases and the shortest human
410 chromosome alone is ~45 megabases long. Furthermore, the only added complexity when generating
411 the random sequences was a GC content of 41%. Real genomes also experience more complex mutation
412 events involving translocations and duplications, which, together with the increased length and the
413 presence of repetitive elements, should increase mapping bias in empirical studies. Finally, the range
414 of possible demographic histories including the relationships of targets and sources, ~~drift as well as~~
415 the amount of drift, and the timing and number of gene flow events is impossible to explore in a
416 simulation study. The restricted scenarios tested in this study should affect the quantitative results
417 but the qualitative interpretation of mapping bias impacting ancestry estimates should extend beyond
418 the specific model used in the simulations.

419 While the ancestry estimates depended slightly on the reference genome the reads were mapped to,
420 they seemed more influenced by the choice of method or software. Methods ~~easily~~ differed by more
421 than 10% in their ancestry estimates from the same source data. This highlights that other factors
422 and biases play major roles in the performance of these methods. Depending on the method, the
423 type of input data, and the implementation, they showed different sensitivities to e.g. linkage or the
424 amount of missing data ~~or linkage~~ (which was on average ~37% per SNP for the 0.5x and ~3% for
425 the 2.0x simulations). For non-pruned data, qpAdm performed best across all scenarios and did not
426 show any method-specific bias in certain ranges of simulated admixture proportions. ~~This supports~~
427 Multiple differences between the PSD and qpAdm methods may have contributed to the relative biases
428 we observed. PSD models may propagate allele-frequency misestimation more than qpAdm because
429 of their assumptions of linkage equilibrium and Hardy-Weinberg equilibrium. Indeed, we observed
430 that LD pruning improved the performance of PSD models, but they are known to be sensitive
431 to sample size and drift (e.g. Lawson et al., 2018; Toyama et al., 2020). More generally, because it
432 is based on Patterson’s f statistics (Patterson et al., 2012), qpAdm estimates ancestry from relative
433 differences. If mapping bias affects all populations similarly, then their relative relationships remain
434 more stable. In contrast, PSD models reconstruct exact allele frequencies for the putative source
435 populations therefore emphasizing the impact of mapping bias. Finally, the ancestry proportions of
436 PSD models are constrained to [0, 1] which is not the case for qpAdm. Indeed, we see negative estimates
437 in a small number of simulations (3 runs with 0.5X depth and 50,000 generations divergence). This
438 (biologically unrealistic) flexibility of qpAdm compared to PSD models drives the mean estimated
439 admixture admixture proportion down, which may account for some of the reduction in upward
440 method bias compared to the other methods.

441 Broadly speaking, our results support the common practice of using qpAdm in most human ancient
442 DNA studies. However, the requirement of data from additional, “right” populations, ~~might not~~
443 ~~make it applicable~~ may make it difficult to apply to many non-human species. Furthermore, qpAdm
444 only works with genotype calls, so it is influenced by mapping bias in similar ways as ADMIXTURE
445 and these methods cannot benefit from the newly introduced genotype likelihood estimation. We
446 also need to note that we tested qpAdm under almost ideal settings in our simulations with left and
447 right populations clearly separated and without gene flow between them. More thorough assessments
448 of the performance of qpAdm can be found elsewhere (Harney et al., 2021; Yüncü et al., 2023). In
449 our simulations, unsupervised PSD-model approaches (ADMIXTURE, NGSadmixmap) work as well as or even
450 better than supervised PSD-model approaches (ADMIXTURE, fastNGSadmixmap) in estimating the ancestry
451 proportions in the target population. ADMIXTURE and NGSadmixmap benefit from LD pruning while LD
452 pruning increases the method bias for fastNGSadmixmap and introduces method bias for qpAdm.

453 Genotype likelihood-based methods for estimating ancestry proportions are not commonly used
454 in human ancient DNA studies (but ~~they genotype likelihoods~~ are popular as input for imputation
455 pipelines). This may be surprising, because genotype-likelihood-based approaches are targeted at low
456 coverage data, exactly as seen in ancient DNA studies. However, the definition of “low coverage” differs
457 between fields. While most working with modern DNA would understand ~~2-4X-2-4x~~ as “low depth”,
458 the standards for ancient DNA researchers are ~~usually a lot~~ typically much lower due to limited DNA
459 preservation. Genotype likelihood methods perform much better with ~~>1X-1x~~ coverage, an amount

460 of data that is not within reach for most ancient DNA samples investigated so far (Mallick et al.,
461 2023). The large body of known, common polymorphic sites in human populations allows the use of
462 pseudohaploid calls at those positions instead. Nonetheless, this study highlights that unsupervised
463 methods employing genotype-likelihoods (NGSadmix) can reach similar accuracies as methods such
464 as ADMIXTURE that require (pseudo-haploid) genotype calls. Moreover, methods that incorporate
465 genotype likelihoods have the added benefit that the modified genotype likelihood estimation approach
466 can be used to reduce the effect of mapping bias. Furthermore, if some samples in the dataset have ~~>1x~~
467 1x depth, genotype likelihood-based approaches will benefit from the additional data and provide more
468 precise estimates of ancestry proportions while pseudo-haploid data will not gain any information from
469 more than one read at a position. Finally, genotype likelihoods are very flexible and can be adjusted for
470 many other aspects of the data. For example, variations of genotype likelihood estimators exist that
471 incorporate the effect of post-mortem damage (Hofmanová et al., 2016; Link et al., 2017; Kousathanas
472 et al., 2017) allowing ~~to~~-use of all sequence data without filtering for potentially damaged sites or
473 enzymatic repair of the damages in the wet lab.

474 As the main aim of this study was to show the general impact of mapping bias and introduce
475 a modified genotype likelihood, we opted for a comparison of some of the most popular meth-
476 ods with a limited set of settings. This was done in part to limit the computational load of this
477 study. We also decided to not set this up as a systematic assessment of different factors influenc-
478 ing mapping bias. The effects of fragmentation (~~Günther and Nettelblad, 2019~~) and ~~deamination~~
479 ~~damage~~ (~~Martiniano et al., 2020~~) (shorter fragments increasing bias, Günther and Nettelblad, 2019),
480 deamination damage (deamination increasing the number of mismatches and bias, Martiniano et al., 2020)
481 and mapping algorithm/parameters (Dolenz et al., 2024) on mapping bias have been explored else-
482 where. Our simulations were restricted to one mapping software (*bwa aln*) and the commonly used
483 mapping quality threshold of 30. Mapping quality calculations differ substantially between tools
484 and algorithms making their impact on mapping bias not directly comparable (Dolenz et al., 2024)
485 . For *bwa aln* (Li and Durbin, 2009), it has been suggested that a mapping quality threshold of
486 25 (the value assigned when the maximum number of mismatches is reached) reduces mapping
487 bias (e.g. Martiniano et al., 2020; Dolenz et al., 2024), and we also see a reduction in mapping bias
488 when using these thresholds (Figures S11-S14). Therefore, a general suggestion for users of *bwa*
489 *aln* should be to use 25 as the mapping quality cutoff. However, many users are using other mappers
490 (e.g. *bowtie*, Langmead and Salzberg, 2012) in their research, and adjusted genotype likelihoods allow
491 correcting for mapping bias independent of the mapping software and its specifics in calculating
492 mapping quality values. Our results reiterate that mapping bias can skew results in studies using
493 low-coverage data as is the case in most ancient DNA studies. Different strategies exist for mitigating
494 these effects and we added a modified genotype likelihood approach to the population genomic toolkit.
495 Nevertheless, none of these methods will be the ideal solution in all cases and they will not always
496 fully remove the potential effect of mapping bias, making proper verification and critical presentation
497 of all results crucial.

498 Acknowledgements

499 ~~We are extremely grateful to Amy Goldberg for numerous discussions during the initial phase of~~
500 ~~this project. We thank~~ thank Kay Prüfer for feedback on the preprint and Gabriel Renaud for
501 making code for connecting *msprime* and *gargammel* available on Github. The computations were
502 enabled by resources in projects SNIC 2017/7-259, SNIC 2018/8-6, SNIC 2021/2-17, SNIC 2022/22-
503 874, NAISS 2023/22-883, sllstore2017087, UPPMAX 2023/2-30 and NAISS 2023/2-19 provided by the
504 National Academic Infrastructure for Supercomputing in Sweden (NAISS) and the Swedish National
505 Infrastructure for Computing (SNIC) at Uppmax, partially funded by Uppsala University and the
506 Swedish Research Council through grant agreements no. 2022-06725 and no. 2018-05973.

Funding

507

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

Conflict of interest disclosure

510

511 The authors declare they have no conflict of interest relating to the content of this article. Torsten
512 Günther is a recommender for PCI Genomics and PCI Evolutionary Biology.

Data, script and code availability

513

514 Raw data for the boxplots can be found in Data S1. Code used in this study can be found under
515 https://github.com/tgue/refbias_GL with a snapshot of the version used for this revision available
516 on Zenodo (<https://doi.org/10.5281/zenodo.14505750>). Empirical data from the 1000 genomes
517 project is available from their resources: SNP array data ([http://ftp.1000genomes.ebi.ac.uk/](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/supporting/hd_genotype_chip/ALL.chip.omni_broad_sanger_combined.20140818.snps.genotypes.vcf.gz)
518 [vol1/ftp/release/20130502/supporting/hd_genotype_chip/ALL.chip.omni_broad_sanger_combined.](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/supporting/hd_genotype_chip/ALL.chip.omni_broad_sanger_combined.20140818.snps.genotypes.vcf.gz)
519 [20140818.snps.genotypes.vcf.gz](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/supporting/hd_genotype_chip/ALL.chip.omni_broad_sanger_combined.20140818.snps.genotypes.vcf.gz)) and low coverage sequencing data ([https://ftp.1000genomes.](https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/data/)
520 [ebi.ac.uk/vol1/ftp/phase3/data/](https://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/data/)).

References

521

522 D. H. Alexander and K. Lange. Enhancements to the ADMIXTURE algorithm for individ-
523 ual ancestry estimation. *BMC Bioinformatics*, 12(1):246, June 2011. ISSN 1471-2105. doi:
524 10.1186/1471-2105-12-246. URL <https://doi.org/10.1186/1471-2105-12-246>.

525 D. H. Alexander, J. Novembre, and K. Lange. Fast model-based estimation of ancestry in unrelated
526 individuals. *Genome research*, 19(9):1655–1664, 2009. ISSN 1088-9051. Number: 9 Publisher: Cold
527 Spring Harbor Lab.

528 A. Auton, G. R. Abecasis, D. M. Altshuler, R. M. Durbin, G. R. Abecasis, D. R. Bentley,
529 A. Chakravarti, A. G. Clark, P. Donnelly, E. E. Eichler, P. Flicek, S. B. Gabriel, R. A. Gibbs,
530 E. D. Green, M. E. Hurles, B. M. Knoppers, J. O. Korbel, E. S. Lander, C. Lee, H. Lehrach, E. R.
531 Mardis, G. T. Marth, G. A. McVean, D. A. Nickerson, J. P. Schmidt, S. T. Sherry, J. Wang, R. K.
532 Wilson, R. A. Gibbs, E. Boerwinkle, H. Doddapaneni, Y. Han, V. Korchina, C. Kovar, S. Lee,
533 D. Muzny, J. G. Reid, Y. Zhu, Y. Chang, Q. Feng, X. Fang, X. Guo, M. Jian, H. Jiang, X. Jin,
534 T. Lan, G. Li, J. Li, Y. Li, S. Liu, X. Liu, Y. Lu, X. Ma, M. Tang, B. Wang, G. Wang, H. Wu, R. Wu,
535 X. Xu, Y. Yin, D. Zhang, W. Zhang, J. Zhao, M. Zhao, X. Zheng, E. S. Lander, D. M. Altshuler,
536 S. B. Gabriel, N. Gupta, N. Gharani, L. H. Toji, N. P. Gerry, A. M. Resch, P. Flicek, J. Barker,
537 L. Clarke, L. Gil, S. E. Hunt, G. Kelman, E. Kulesha, R. Leinonen, W. M. McLaren, R. Rad-
538 hakrishnan, A. Roa, D. Smirnov, R. E. Smith, I. Streeter, A. Thormann, I. Toneva, B. Vaughan,
539 X. Zheng-Bradley, D. R. Bentley, R. Grocock, S. Humphray, T. James, Z. Kingsbury, H. Lehrach,
540 R. Sudbrak, M. W. Albrecht, V. S. Amstislavskiy, T. A. Borodina, M. Lienhard, F. Mertes, M. Sul-
541 tan, B. Timmermann, M.-L. Yaspo, E. R. Mardis, R. K. Wilson, L. Fulton, R. Fulton, S. T. Sherry,
542 V. Ananiev, Z. Belaia, D. Beloslyudtsev, N. Bouk, C. Chen, D. Church, R. Cohen, C. Cook, J. Gar-
543 ner, T. Hefferon, M. Kimelman, C. Liu, J. Lopez, P. Meric, C. O’Sullivan, Y. Ostapchuk, L. Phan,
544 S. Ponomarov, V. Schneider, E. Shekhtman, K. Sirotkin, D. Slotta, H. Zhang, G. A. McVean, R. M.
545 Durbin, S. Balasubramaniam, J. Burton, P. Danecek, T. M. Keane, A. Kolb-Kokocinski, S. Mc-
546 Carthy, J. Stalker, M. Quail, J. P. Schmidt, C. J. Davies, J. Gollub, T. Webster, B. Wong, Y. Zhan,
547 A. Auton, C. L. Campbell, Y. Kong, A. Marcketta, R. A. Gibbs, F. Yu, L. Antunes, M. Bainbridge,
548 D. Muzny, A. Sabo, Z. Huang, J. Wang, L. J. M. Coin, L. Fang, X. Guo, X. Jin, G. Li, Q. Li,
549 Y. Li, Z. Li, H. Lin, B. Liu, R. Luo, H. Shao, Y. Xie, C. Ye, C. Yu, F. Zhang, H. Zheng, H. Zhu,
550 C. Alkan, E. Dal, F. Kahveci, G. T. Marth, E. P. Garrison, D. Kural, W.-P. Lee, W. Fung Leong,

- 551 M. Stromberg, A. N. Ward, J. Wu, M. Zhang, M. J. Daly, M. A. DePristo, R. E. Handsaker, D. M.
552 Altshuler, E. Banks, G. Bhatia, G. del Angel, S. B. Gabriel, G. Genovese, N. Gupta, H. Li, S. Kashin,
553 E. S. Lander, S. A. McCarroll, J. C. Nemes, R. E. Poplin, S. C. Yoon, J. Lihm, V. Makarov, A. G.
554 Clark, S. Gottipati, A. Keinan, J. L. Rodriguez-Flores, J. O. Korbel, T. Rausch, M. H. Fritz, A. M.
555 Stütz, P. Flicek, K. Beal, L. Clarke, A. Datta, J. Herrero, W. M. McLaren, G. R. S. Ritchie,
556 R. E. Smith, D. Zerbino, X. Zheng-Bradley, P. C. Sabeti, I. Shlyakhter, S. F. Schaffner, J. Vitti,
557 D. N. Cooper, E. V. Ball, P. D. Stenson, D. R. Bentley, M. Bauer, R. Keira Cheetham, A. Cox,
558 M. Eberle, S. Humphray, S. Kahn, L. Murray, J. Peden, R. Shaw, E. E. Kenny, M. A. Batzer, M. K.
559 Konkel, J. A. Walker, D. G. MacArthur, M. Lek, R. Sudbrak, V. S. Amstislavskiy, R. Herwig,
560 E. R. Mardis, L. Ding, D. C. Koboldt, D. Larson, and K. Ye. A global reference for human genetic
561 variation. *Nature*, 526(7571):68–74, Oct. 2015. ISSN 1476-4687. doi: 10.1038/nature15393. URL
562 <https://www.nature.com/articles/nature15393>. Publisher: Nature Publishing Group.
- 563 A. R. Barton, C. G. Santander, P. Skoglund, I. Moltke, D. Reich, and I. Mathieson. Insuffi-
564 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:
565 Contradictory Results.
566
- 567 L. A. Bergeron, S. Besenbacher, J. Zheng, P. Li, M. F. Bertelsen, B. Quintard, J. I. Hoffman,
568 Z. Li, J. St. Leger, C. Shao, J. Stiller, M. T. P. Gilbert, M. H. Schierup, and G. Zhang. Evo-
569 lution of the germline mutation rate across vertebrates. *Nature*, 615(7951):285–291, Mar. 2023.
570 ISSN 1476-4687. doi: 10.1038/s41586-023-05752-y. URL <https://www.nature.com/articles/s41586-023-05752-y>. Publisher: Nature Publishing Group.
571
- 572 A. W. Briggs, U. Stenzel, P. L. Johnson, R. E. Green, J. Kelso, K. Prüfer, M. Meyer, J. Krause,
573 M. T. Ronan, M. Lachmann, and others. Patterns of damage in genomic DNA sequences from a
574 Neandertal. *Proceedings of the National Academy of Sciences*, 104(37):14616–14621, 2007.
- 575 C. C. Chang, C. C. Chow, L. C. Tellier, S. Vattikuti, S. M. Purcell, and J. J. Lee. Second-generation
576 PLINK: rising to the challenge of larger and richer datasets. *Gigascience*, 4(1):s13742–015, 2015.
577 ISSN 2047-217X. Number: 1 Publisher: Oxford University Press.
- 578 B. Charlesworth. Effective population size and patterns of molecular evolution and variation. *Nature*
579 *Reviews Genetics*, 10(3):195–205, Mar. 2009. ISSN 1471-0064. doi: 10.1038/nrg2526. URL <https://www.nature.com/articles/nrg2526>. Publisher: Nature Publishing Group.
580
- 581 N.-C. Chen, B. Solomon, T. Mun, S. Iyer, and B. Langmead. Reference flow: reducing reference bias
582 using multiple population genomes. *Genome Biology*, 22(1):8, Jan. 2021. ISSN 1474-760X. doi:
583 10.1186/s13059-020-02229-3. URL <https://doi.org/10.1186/s13059-020-02229-3>.
- 584 D. M. Church, V. A. Schneider, K. M. Steinberg, M. C. Schatz, A. R. Quinlan, C.-S. Chin, P. A. Kitts,
585 B. Aken, G. T. Marth, M. M. Hoffman, J. Herrero, M. L. Z. Mendoza, R. Durbin, and P. Flicek.
586 Extending reference assembly models. *Genome Biology*, 16(1):13, Jan. 2015. ISSN 1465-6906. doi:
587 10.1186/s13059-015-0587-3. URL <https://doi.org/10.1186/s13059-015-0587-3>.
- 588 S. L. Cox, H. M. Moots, J. T. Stock, A. Shbat, B. D. Bitarello, N. Nicklisch, K. W. Alt,
589 W. Haak, E. Rosenstock, C. B. Ruff, and I. Mathieson. Predicting skeletal stature using ancient
590 DNA. *American Journal of Biological Anthropology*, 177(1):162–174, 2022. ISSN 2692-7691. doi:
591 10.1002/ajpa.24426. URL <https://onlinelibrary.wiley.com/doi/abs/10.1002/ajpa.24426>.
592 eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/ajpa.24426>.
- 593 T. Davy, D. Ju, I. Mathieson, and P. Skoglund. Hunter-gatherer admixture facilitated natural selection
594 in Neolithic European farmers. *Current Biology*, 33(7):1365–1371.e3, Apr. 2023. ISSN 0960-9822.
595 doi: 10.1016/j.cub.2023.02.049. URL <https://www.sciencedirect.com/science/article/pii/S0960982223001896>.
596

- 597 S. Dolenz, T. van der Valk, C. Jin, J. Oppenheimer, M. B. Sharif, L. Orlando, B. Shapiro, L. Dalén,
598 and P. D. Heintzman. Unravelling reference bias in ancient DNA datasets. *Bioinformatics*, 40(7):
599 btae436, July 2024. ISSN 1367-4811. doi: 10.1093/bioinformatics/btae436. URL <https://doi.org/10.1093/bioinformatics/btae436>.
600
- 601 B. L. Dumont and B. A. Payseur. EVOLUTION OF THE GENOMIC RATE OF RECOMBINATION
602 IN MAMMALS. *Evolution*, 62(2):276–294, Feb. 2008. ISSN 0014-3820. doi: 10.1111/j.1558-5646.
603 2007.00278.x. URL <https://doi.org/10.1111/j.1558-5646.2007.00278.x>.
- 604 D. Falush, M. Stephens, and J. K. Pritchard. Inference of population structure using multilocus
605 genotype data: linked loci and correlated allele frequencies. *Genetics*, 164(4):1567–1587, 2003.
606 ISSN 0016-6731. Number: 4.
- 607 D. Falush, M. Stephens, and J. K. Pritchard. Inference of population structure using multilocus
608 genotype data: dominant markers and null alleles. *Molecular ecology notes*, 7(4):574–578, 2007.
609 ISSN 1471-8278. Number: 4.
- 610 M. Fumagalli, F. G. Vieira, T. S. Korneliussen, T. Linderoth, E. Huerta-Sánchez, A. Albrechtsen,
611 and R. Nielsen. Quantifying Population Genetic Differentiation from Next-Generation Sequencing
612 Data. *Genetics*, 195(3):979–992, Nov. 2013. ISSN 1943-2631. doi: 10.1534/genetics.113.154740.
613 URL <https://doi.org/10.1534/genetics.113.154740>.
- 614 S. Gopalakrishnan, J. A. Samaniego Castruita, M.-H. S. Sinding, L. F. K. Kuderna, J. Räikkönen,
615 B. Petersen, T. Sicheritz-Ponten, G. Larson, L. Orlando, T. Marques-Bonet, A. J. Hansen, L. Dalén,
616 and M. T. P. Gilbert. The wolf reference genome sequence (*Canis lupus lupus*) and its implications
617 for *Canis* spp. population genomics. *BMC Genomics*, 18:495, June 2017. ISSN 1471-2164. doi:
618 10.1186/s12864-017-3883-3. URL <https://doi.org/10.1186/s12864-017-3883-3>.
- 619 S. Gopalakrishnan, S. S. Ebenesersdóttir, I. K. C. Lundstrøm, G. Turner-Walker, K. H. S. Moore,
620 P. Luisi, A. Margaryan, M. D. Martin, M. R. Ellegaard, Magnússon, Sigursson, S. Snorrardóttir,
621 D. N. Magnúsdóttir, J. E. Laffoon, L. van Dorp, X. Liu, I. Moltke, M. C. Ávila Arcos, J. G.
622 Schraiber, S. Rasmussen, D. Juan, P. Gelabert, T. de Dios, A. K. Fotakis, M. Iraeta-Orbegozo,
623 J. Vågane, S. D. Denham, A. Christophersen, H. K. Stenøien, F. G. Vieira, S. Liu, T. Günther,
624 T. Kivisild, O. G. Moseng, B. Skar, C. Cheung, M. Sandoval-Velasco, N. Wales, H. Schroeder, P. F.
625 Campos, V. B. Gumundsdóttir, T. Sicheritz-Ponten, B. Petersen, J. Halgunset, E. Gilbert, G. L.
626 Cavalleri, E. Hovig, I. Kockum, T. Olsson, L. Alfredsson, T. F. Hansen, T. Werge, E. Willerslev,
627 F. Balloux, T. Marques-Bonet, C. Lalueza-Fox, R. Nielsen, K. Stefánsson, A. Helgason, and M. T. P.
628 Gilbert. The population genomic legacy of the second plague pandemic. *Current Biology*, 32
629 (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>.
630
- 631 R. E. Green, J. Krause, A. W. Briggs, T. Maricic, U. Stenzel, M. Kircher, N. Patterson, H. Li, W. Zhai,
632 and M. H.-Y. Fritz. A draft sequence of the Neandertal genome. *science*, 328(5979):710–722, 2010.
633 ISSN 0036-8075. Number: 5979 Publisher: American Association for the Advancement of Science.
- 634 T. Günther and M. Jakobsson. Population genomic analyses of DNA from ancient remains. In
635 *Handbook of statistical genomics*, pages 295–324. John Wiley & Sons, 4th edition, 2019. ISBN
636 1-119-42914-5.
- 637 T. Günther and C. Nettelblad. The presence and impact of reference bias on population genomic stud-
638 ies of prehistoric human populations. *PLOS Genetics*, 15(7):e1008302, July 2019. ISSN 1553-7404.
639 doi: 10.1371/journal.pgen.1008302. URL [https://journals.plos.org/plosgenetics/article?
640 id=10.1371/journal.pgen.1008302](https://journals.plos.org/plosgenetics/article?id=10.1371/journal.pgen.1008302).

- 641 W. Haak, I. Lazaridis, N. Patterson, N. Rohland, S. Mallick, B. Llamas, G. Brandt, S. Nordenfelt,
642 E. Harney, and K. Stewardson. Massive migration from the steppe was a source for Indo-European
643 languages in Europe. *Nature*, 522(7555):207–211, 2015. ISSN 1476-4687. Number: 7555 Publisher:
644 Nature Publishing Group.
- 645 K. Hanghøj, I. Moltke, P. A. Andersen, A. Manica, and T. S. Korneliussen. Fast and accu-
646 rate relatedness estimation from high-throughput sequencing data in the presence of inbreed-
647 ing. *GigaScience*, 8(5), May 2019. ISSN 2047-217X. doi: 10.1093/gigascience/giz034. URL
648 <https://doi.org/10.1093/gigascience/giz034>.
- 649 E. Harney, N. Patterson, D. Reich, and J. Wakeley. Assessing the performance of qpAdm: a statistical
650 tool for studying population admixture. *Genetics*, 217(4), Apr. 2021. ISSN 1943-2631. doi: 10.
651 1093/genetics/iyaa045. URL <https://doi.org/10.1093/genetics/iyaa045>.
- 652 P. D. Heintzman, G. D. Zazula, R. D. MacPhee, E. Scott, J. A. Cahill, B. K. McHorse, J. D. Kapp,
653 M. Stiller, M. J. Wooller, L. Orlando, J. Southon, D. G. Froese, and B. Shapiro. A new genus of
654 horse from Pleistocene North America. *eLife*, 6, 2017. ISSN 2050-084X. doi: 10.7554/eLife.29944.
- 655 Z. Hofmanová, S. Kreutzer, G. Hellenthal, C. Sell, Y. Diekmann, D. Díez-del Molino, L. van Dorp,
656 S. López, A. Kousathanas, V. Link, and others. Early farmers from across Europe directly descended
657 from Neolithic Aegeans. *Proceedings of the National Academy of Sciences*, page 201523951, 2016.
- 658 L. Huang, V. Popic, and S. Batzoglou. Short read alignment with populations of genomes. *Bioin-*
659 *formatics*, 29(13):i361–i370, July 2013. ISSN 1367-4803. doi: 10.1093/bioinformatics/btt215. URL
660 <https://doi.org/10.1093/bioinformatics/btt215>.
- 661 M. J. Hubisz, D. Falush, M. Stephens, and J. K. Pritchard. Inferring weak population structure with
662 the assistance of sample group information. *Molecular ecology resources*, 9(5):1322–1332, 2009. ISSN
663 1755-098X. Number: 5.
- 664 R. Hui, C. L. Scheib, E. D’Atanasio, S. A. Inskip, C. Cessford, S. A. Biagini, A. W. Wohms, M. Q.
665 Ali, S. J. Griffith, A. Solnik, H. Niinemäe, X. J. Ge, A. K. Rose, O. Beneker, T. C. O’Connell, J. E.
666 Robb, and T. Kivisild. Genetic history of Cambridgeshire before and after the Black Death. *Science*
667 *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
668 Science.
- 669
- 670 E. Jørsboe, K. Hanghøj, and A. Albrechtsen. fastNGSadmix: admixture proportions and principal
671 component analysis of a single NGS sample. *Bioinformatics*, 33(19):3148–3150, 2017.
- 672 J. Kelleher, A. M. Etheridge, and G. McVean. Efficient coalescent simulation and genealogical analysis
673 for large sample sizes. *PLoS computational biology*, 12(5):e1004842, 2016.
- 674 J. Klunk, T. P. Vilgalys, C. E. Demeure, X. Cheng, M. Shiratori, J. Madej, R. Beau, D. Elli, M. I.
675 Patino, R. Redfern, S. N. DeWitte, J. A. Gamble, J. L. Boldsen, A. Carmichael, N. Varlik, K. Eaton,
676 J.-C. Grenier, G. B. Golding, A. Devault, J.-M. Rouillard, V. Yotova, R. Sindeaux, C. J. Ye,
677 M. Bikaran, A. Dumaine, J. F. Brinkworth, D. Missiakas, G. A. Rouleau, M. Steinrücken, J. Pizarro-
678 Cerdá, H. N. Poinar, and L. B. Barreiro. Evolution of immune genes is associated with the Black
679 Death. *Nature*, 611(7935):312–319, Nov. 2022. ISSN 1476-4687. doi: 10.1038/s41586-022-05349-x.
680 URL <https://www.nature.com/articles/s41586-022-05349-x>. Number: 7935 Publisher: Na-
681 ture Publishing Group.
- 682 D. Koptekin, E. Yapar, K. B. Vural, E. Sağlıcan, N. E. Altınışık, A.-S. Malaspinas, C. Alkan, and
683 M. Somel. Pre-processing of paleogenomes: Mitigating reference bias and postmortem damage in
684 ancient genome data, Nov. 2023. URL [https://www.biorxiv.org/content/10.1101/2023.11.](https://www.biorxiv.org/content/10.1101/2023.11.11.566695v1)
685 [11.566695v1](https://www.biorxiv.org/content/10.1101/2023.11.11.566695v1). Pages: 2023.11.11.566695 Section: New Results.

- 686 T. S. Korneliussen and I. Moltke. NgsRelate: a software tool for estimating pairwise relatedness
687 from next-generation sequencing data. *Bioinformatics*, 31(24):4009–4011, 2015. ISSN 1460-2059.
688 Number: 24 Publisher: Oxford University Press.
- 689 T. S. Korneliussen, I. Moltke, A. Albrechtsen, and R. Nielsen. Calculation of Tajima’s D and other
690 neutrality test statistics from low depth next-generation sequencing data. *BMC bioinformatics*, 14:
691 289, Oct. 2013. ISSN 1471-2105. doi: 10.1186/1471-2105-14-289.
- 692 T. S. Korneliussen, A. Albrechtsen, and R. Nielsen. ANGSD: Analysis of Next Generation Sequencing
693 Data. *BMC bioinformatics*, 15(1):356, 2014. ISSN 1471-2105. doi: 10.1186/s12859-014-0356-4.
- 694 A. Kousathanas, C. Leuenberger, V. Link, C. Sell, J. Burger, and D. Wegmann. Inferring Heterozygosity
695 from Ancient and Low Coverage Genomes. *Genetics*, 205(1):317–332, Jan. 2017. ISSN 0016-6731,
696 1943-2631. doi: 10.1534/genetics.116.189985. URL [http://www.genetics.org/content/205/1/
697 317](http://www.genetics.org/content/205/1/317).
- 698 E. S. Lander, L. M. Linton, B. Birren, C. Nusbaum, M. C. Zody, J. Baldwin, K. Devon, K. De-
699 war, M. Doyle, W. FitzHugh, R. Funke, D. Gage, K. Harris, A. Heaford, J. Howland, L. Kann,
700 J. Lehoczy, R. LeVine, P. McEwan, K. McKernan, J. Meldrim, J. P. Mesirov, C. Miranda, W. Mor-
701 ris, J. Naylor, C. Raymond, M. Rosetti, R. Santos, A. Sheridan, C. Sougnez, N. Stange-Thomann,
702 N. Stojanovic, A. Subramanian, D. Wyman, J. Rogers, J. Sulston, R. Ainscough, S. Beck, D. Bent-
703 ley, J. Burton, C. Clee, N. Carter, A. Coulson, R. Deadman, P. Deloukas, A. Dunham, I. Dunham,
704 R. Durbin, L. French, D. Grafham, S. Gregory, T. Hubbard, S. Humphray, A. Hunt, M. Jones,
705 C. Lloyd, A. McMurray, L. Matthews, S. Mercer, S. Milne, J. C. Mullikin, A. Mungall, R. Plumb,
706 M. Ross, R. Shownkeen, S. Sims, R. H. Waterston, R. K. Wilson, L. W. Hillier, J. D. McPherson,
707 M. A. Marra, E. R. Mardis, L. A. Fulton, A. T. Chinwalla, K. H. Pepin, W. R. Gish, S. L. Chissoe,
708 M. C. Wendl, K. D. Delehaunty, T. L. Miner, A. Delehaunty, J. B. Kramer, L. L. Cook, R. S. Fulton,
709 D. L. Johnson, P. J. Minx, S. W. Clifton, T. Hawkins, E. Branscomb, P. Predki, P. Richardson,
710 S. Wenning, T. Slezak, N. Doggett, J.-F. Cheng, A. Olsen, S. Lucas, C. Elkin, E. Uberbacher, M. Fra-
711 zier, R. A. Gibbs, D. M. Muzny, S. E. Scherer, J. B. Bouck, E. J. Sodergren, K. C. Worley, C. M.
712 Rives, J. H. Gorrell, M. L. Metzker, S. L. Naylor, R. S. Kucherlapati, D. L. Nelson, G. M. Wein-
713 stock, Y. Sakaki, A. Fujiyama, M. Hattori, T. Yada, A. Toyoda, T. Itoh, C. Kawagoe, H. Watanabe,
714 Y. Totoki, T. Taylor, J. Weissenbach, R. Heilig, W. Saurin, F. Artiguenave, P. Brottier, T. Bruls,
715 E. Pelletier, C. Robert, P. Wincker, A. Rosenthal, M. Platzer, G. Nyakatura, S. Taudien, A. Rump,
716 D. R. Smith, L. Doucette-Stamm, M. Rubenfield, K. Weinstock, H. M. Lee, J. Dubois, H. Yang,
717 J. Yu, J. Wang, G. Huang, J. Gu, L. Hood, L. Rowen, A. Madan, S. Qin, R. W. Davis, N. A. Feder-
718 spiel, A. P. Abola, M. J. Proctor, B. A. Roe, F. Chen, H. Pan, J. Ramser, H. Lehrach, R. Reinhardt,
719 W. R. McCombie, M. de la Bastide, N. Dedhia, H. Blöcker, K. Hornischer, G. Nordsiek, R. Agarwala,
720 L. Aravind, J. A. Bailey, A. Bateman, S. Batzoglou, E. Birney, P. Bork, D. G. Brown, C. B. Burge,
721 L. Cerutti, H.-C. Chen, D. Church, M. Clamp, R. R. Copley, T. Doerks, S. R. Eddy, E. E. Eichler,
722 T. S. Furey, J. Galagan, J. G. R. Gilbert, C. Harmon, Y. Hayashizaki, D. Haussler, H. Hermjakob,
723 K. Hokamp, W. Jang, L. S. Johnson, T. A. Jones, S. Kasif, A. Kasprzyk, S. Kennedy, W. J. Kent,
724 P. Kitts, E. V. Koonin, I. Korf, D. Kulp, D. Lancet, T. M. Lowe, A. McLysaght, T. Mikkelsen, J. V.
725 Moran, N. Mulder, V. J. Pollara, C. P. Ponting, G. Schuler, J. Schultz, G. Slater, A. F. A. Smit,
726 E. Stupka, J. Szustakowki, D. Thierry-Mieg, J. Thierry-Mieg, L. Wagner, J. Wallis, R. Wheeler,
727 A. Williams, Y. I. Wolf, K. H. Wolfe, S.-P. Yang, R.-F. Yeh, F. Collins, M. S. Guyer, J. Peterson,
728 A. Felsenfeld, K. A. Wetterstrand, R. M. Myers, J. Schmutz, M. Dickson, J. Grimwood, D. R.
729 Cox, M. V. Olson, R. Kaul, C. Raymond, N. Shimizu, K. Kawasaki, S. Minoshima, G. A. Evans,
730 M. Athanasiou, R. Schultz, A. Patrinos, M. J. Morgan, International Human Genome Sequencing
731 Consortium, C. f. G. R. Whitehead Institute for Biomedical Research, The Sanger Centre:, Wash-
732 ington University Genome Sequencing Center, US DOE Joint Genome Institute:, Baylor College
733 of Medicine Human Genome Sequencing Center:, RIKEN Genomic Sciences Center:, Genoscope
734 and CNRS UMR-8030:, I. o. M. B. Department of Genome Analysis, GTC Sequencing Center:,

735 Beijing Genomics Institute/Human Genome Center:, T. I. f. S. B. Multimegabase Sequencing Cen-
736 ter, Stanford Genome Technology Center:, University of Oklahoma’s Advanced Center for Genome
737 Technology:, Max Planck Institute for Molecular Genetics:, L. A. H. G. C. Cold Spring Harbor Lab-
738 oratory, GBF—German Research Centre for Biotechnology:, a. i. i. l. u. o. h. *Genome Analysis
739 Group (listed in alphabetical order, U. N. I. o. H. Scientific management: National Human Genome
740 Research Institute, Stanford Human Genome Center:, University of Washington Genome Center:,
741 K. U. S. o. M. Department of Molecular Biology, University of Texas Southwestern Medical Center
742 at Dallas:, U. D. o. E. Office of Science, and The Wellcome Trust:. Initial sequencing and analysis of
743 the human genome. *Nature*, 409(6822):860–921, Feb. 2001. ISSN 1476-4687. doi: 10.1038/35057062.
744 URL <https://www.nature.com/articles/35057062>. Number: 6822 Publisher: Nature Publishing
745 Group.

746 B. Langmead and S. L. Salzberg. Fast gapped-read alignment with Bowtie 2. *Nature Methods*, 9(4):
747 357–359, Apr. 2012. ISSN 1548-7105. doi: 10.1038/nmeth.1923. URL [https://www.nature.com/](https://www.nature.com/articles/nmeth.1923)
748 [articles/nmeth.1923](https://www.nature.com/articles/nmeth.1923). Number: 4 Publisher: Nature Publishing Group.

749 D. J. Lawson, L. van Dorp, and D. Falush. A tutorial on how not to over-interpret STRUCTURE and
750 ADMIXTURE bar plots. *Nature Communications*, 9(1):3258, Aug. 2018. ISSN 2041-1723. doi:
751 10.1038/s41467-018-05257-7. URL <https://www.nature.com/articles/s41467-018-05257-7>.
752 Number: 1 Publisher: Nature Publishing Group.

753 H. Li and R. Durbin. Fast and accurate short read alignment with Burrows–Wheeler transform.
754 *bioinformatics*, 25(14):1754–1760, 2009. ISSN 1367-4803. Number: 14 Publisher: Oxford University
755 Press.

756 H. Li, B. Handsaker, A. Wysoker, T. Fennell, J. Ruan, N. Homer, G. Marth, G. Abecasis, R. Durbin,
757 and 1000 Genome Project Data Processing Subgroup. The Sequence Alignment/Map format and
758 SAMtools. *Bioinformatics (Oxford, England)*, 25(16):2078–2079, Aug. 2009. ISSN 1367-4811. doi:
759 10.1093/bioinformatics/btp352.

760 V. Link, A. Kousathanas, K. Veeramah, C. Sell, A. Scheu, and D. Wegmann. ATLAS: analysis tools
761 for low-depth and ancient samples. *bioRxiv*, page 105346, 2017.

762 R. N. Lou, A. Jacobs, A. P. Wilder, and N. O. Therikildsen. A beginner’s guide to low-coverage
763 whole genome sequencing for population genomics. *Molecular Ecology*, 30(23):5966–5993, 2021.
764 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)
765 [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>.

766 S. Mallick, A. Micco, M. Mah, H. Ringbauer, I. Lazaridis, I. Olalde, N. Patterson, and D. Reich. The
767 Allen Ancient DNA Resource (AADR): A curated compendium of ancient human genomes, Apr.
768 2023. URL <https://www.biorxiv.org/content/10.1101/2023.04.06.535797v1>.

769 R. Martiniano, E. Garrison, E. R. Jones, A. Manica, and R. Durbin. Removing reference bias and
770 improving indel calling in ancient DNA data analysis by mapping to a sequence variation graph.
771 *Genome Biology*, 21(1):250, Sept. 2020. ISSN 1474-760X. doi: 10.1186/s13059-020-02160-7. URL
772 <https://doi.org/10.1186/s13059-020-02160-7>.

773 I. Mathieson and J. Terhorst. Direct detection of natural selection in Bronze Age Britain. *Genome*
774 *Research*, 32(11-12):2057–2067, Nov. 2022. ISSN 1088-9051, 1549-5469. doi: 10.1101/gr.276862.122.
775 URL <https://genome.cshlp.org/content/32/11-12/2057>. Company: Cold Spring Harbor Lab-
776 oratory Press Distributor: Cold Spring Harbor Laboratory Press Institution: Cold Spring Harbor
777 Laboratory Press Label: Cold Spring Harbor Laboratory Press Publisher: Cold Spring Harbor Lab.

778 I. Mathieson, I. Lazaridis, N. Rohland, S. Mallick, N. Patterson, S. A. Roodenberg, E. Harney, K. Stew-
779 ardson, D. Fernandes, M. Novak, and others. Genome-wide patterns of selection in 230 ancient
780 Eurasians. *Nature*, 528(7583):499–503, 2015.

- 781 I. Mathieson, F. Abascal, L. Vinner, P. Skoglund, C. Pomilla, P. Mitchell, C. Arthur, D. Gurdasani,
782 E. Willerslev, M. S. Sandhu, and G. Dewar. An Ancient Baboon Genome Demonstrates Long-Term
783 Population Continuity in Southern Africa. *Genome Biology and Evolution*, 12(4):407–412, Apr.
784 2020. ISSN 1759-6653. doi: 10.1093/gbe/evaa019. URL <https://doi.org/10.1093/gbe/evaa019>.
- 785 A. McKenna, M. Hanna, E. Banks, A. Sivachenko, K. Cibulskis, A. Kernytsky, K. Garimella, D. Alt-
786 shuler, S. Gabriel, M. Daly, and M. A. DePristo. The Genome Analysis Toolkit: A MapReduce
787 framework for analyzing next-generation DNA sequencing data. *Genome Research*, 20(9):1297–1303,
788 Sept. 2010. ISSN 1088-9051. doi: 10.1101/gr.107524.110. URL [https://www.ncbi.nlm.nih.gov/
789 pmc/articles/PMC2928508/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2928508/).
- 790 J. Meisner and A. Albrechtsen. Inferring population structure and admixture proportions in low-
791 depth NGS data. *Genetics*, 210(2):719–731, 2018. ISSN 1943-2631. Number: 2 Publisher: Oxford
792 University Press.
- 793 R. Nielsen, J. S. Paul, A. Albrechtsen, and Y. S. Song. Genotype and SNP calling from next-generation
794 sequencing data. *Nature Reviews Genetics*, 12(6):443, 2011.
- 795 A. K. Nøhr, K. Hanghøj, G. Garcia-Erill, Z. Li, I. Moltke, and A. Albrechtsen. NGSremix: a soft-
796 ware tool for estimating pairwise relatedness between admixed individuals from next-generation
797 sequencing data. *G3*, (jkab174), May 2021. ISSN 2160-1836. doi: 10.1093/g3journal/jkab174. URL
798 <https://doi.org/10.1093/g3journal/jkab174>.
- 799 A. Oliva, R. Tobler, A. Cooper, B. Llamas, and Y. Souilmi. Systematic benchmark of ancient DNA
800 read mapping. *Briefings in Bioinformatics*, (bbab076), Apr. 2021. ISSN 1477-4054. doi: 10.1093/
801 bib/bbab076. URL <https://doi.org/10.1093/bib/bbab076>.
- 802 L. Orlando, A. Ginolhac, G. Zhang, D. Froese, A. Albrechtsen, M. Stiller, M. Schubert, E. Cap-
803 pellini, B. Petersen, I. Moltke, P. L. F. Johnson, M. Fumagalli, J. T. Vilstrup, M. Raghavan,
804 T. Korneliussen, A.-S. Malaspinas, J. Vogt, D. Szklarczyk, C. D. Kelstrup, J. Vinther, A. Dolocan,
805 J. Stenderup, A. M. V. Velazquez, J. Cahill, M. Rasmussen, X. Wang, J. Min, G. D. Zazula,
806 A. Seguin-Orlando, C. Mortensen, K. Magnussen, J. F. Thompson, J. Weinstock, K. Gregersen,
807 K. H. Røed, V. Eisenmann, C. J. Rubin, D. C. Miller, D. F. Antczak, M. F. Bertelsen, S. Brunak,
808 K. A. S. Al-Rasheid, O. Ryder, L. Andersson, J. Mundy, A. Krogh, M. T. P. Gilbert, K. Kjær,
809 T. Sicheritz-Ponten, L. J. Jensen, J. V. Olsen, M. Hofreiter, R. Nielsen, B. Shapiro, J. Wang, and
810 E. Willerslev. Recalibrating Equus evolution using the genome sequence of an early Middle Pleis-
811 tocene horse. *Nature*, 499(7456):74–78, July 2013. ISSN 1476-4687. doi: 10.1038/nature12323.
812 URL <https://www.nature.com/articles/nature12323>. Bandiera_abtest: a Cg_type: Nature Re-
813 search Journals Number: 7456 Primary_atype: Research Publisher: Nature Publishing Group Sub-
814 ject_term: Evolutionary genetics Subject_term_id: evolutionary-genetics.
- 815 L. Orlando, R. Allaby, P. Skoglund, C. Der Sarkissian, P. W. Stockhammer, M. C. Ávila Arcos,
816 Q. Fu, J. Krause, E. Willerslev, A. C. Stone, and C. Warinner. Ancient DNA analysis. *Nature*
817 *Reviews Methods Primers*, 1(1):1–26, Feb. 2021. ISSN 2662-8449. doi: 10.1038/s43586-020-00011-0.
818 URL <https://www.nature.com/articles/s43586-020-00011-0>. Number: 1 Publisher: Nature
819 Publishing Group.
- 820 N. Patterson, A. L. Price, and D. Reich. Population structure and eigenanalysis. *PLoS genetics*, 2
821 (12):e190, 2006. ISSN 1553-7390. Number: 12 Publisher: Public Library of Science San Francisco,
822 USA.
- 823 N. Patterson, P. Moorjani, Y. Luo, S. Mallick, N. Rohland, Y. Zhan, T. Genschoreck, T. Webster, and
824 D. Reich. Ancient admixture in human history. *Genetics*, 192(3):1065–1093, 2012. ISSN 1943-2631.
825 Number: 3 Publisher: Oxford University Press.

- 826 A. Prasad, E. D. Lorenzen, and M. V. Westbury. Evaluating the role of reference-genome phy-
827 logenetic distance on evolutionary inference. *Molecular Ecology Resources*, 22(1):45–55, 2022.
828 ISSN 1755-0998. doi: 10.1111/1755-0998.13457. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/1755-0998.13457>.
829 eprint: [https://onlinelibrary.wiley.com/doi/pdf/10.1111/1755-](https://onlinelibrary.wiley.com/doi/pdf/10.1111/1755-0998.13457)
830 0998.13457.
- 831 A. L. Price, N. J. Patterson, R. M. Plenge, M. E. Weinblatt, N. A. Shadick, and D. Reich. Principal
832 components analysis corrects for stratification in genome-wide association studies. *Nature genetics*,
833 38(8):904–909, 2006. ISSN 1546-1718. Number: 8 Publisher: Nature Publishing Group.
- 834 J. K. Pritchard, M. Stephens, and P. Donnelly. Inference of population structure using multilocus
835 genotype data. *Genetics*, 155(2):945–959, 2000. ISSN 0016-6731. Number: 2.
- 836 K. Prüfer. snpAD: An ancient DNA genotype caller. *Bioinformatics*, 2018. doi: 10.1093/
837 bioinformatics/bty507. URL [https://academic.oup.com/bioinformatics/advance-article/](https://academic.oup.com/bioinformatics/advance-article/doi/10.1093/bioinformatics/bty507/5042170)
838 [doi/10.1093/bioinformatics/bty507/5042170](https://academic.oup.com/bioinformatics/advance-article/doi/10.1093/bioinformatics/bty507/5042170).
- 839 G. Renaud, K. Hanghøj, E. Willerslev, and L. Orlando. gargammel: a sequence simulator for ancient
840 DNA. *Bioinformatics*, 33(4):577–579, Feb. 2017. ISSN 1367-4803. doi: 10.1093/bioinformatics/
841 btw670. URL <https://academic.oup.com/bioinformatics/article/33/4/577/2608651>.
- 842 A. R. Rogers, R. J. Bohlender, and C. D. Huff. Early history of Neanderthals and Deniso-
843 vans. *Proceedings of the National Academy of Sciences*, 114(37):9859–9863, Sept. 2017. doi:
844 10.1073/pnas.1706426114. URL <https://www.pnas.org/doi/10.1073/pnas.1706426114>. Pub-
845 lisher: Proceedings of the National Academy of Sciences.
- 846 N. Rohland, S. Mallick, M. Mah, R. Maier, N. Patterson, and D. Reich. Three assays for in-solution
847 enrichment of ancient human DNA at more than a million SNPs. *Genome Research*, 32(11-12):
848 2068–2078, Nov. 2022. ISSN 1088-9051, 1549-5469. doi: 10.1101/gr.276728.122. URL <https://genome.cshlp.org/content/32/11-12/2068>.
849 Company: Cold Spring Harbor Laboratory Press Distributor: Cold Spring Harbor Laboratory Press Institution: Cold Spring Harbor Laboratory
850 Press Label: Cold Spring Harbor Laboratory Press Publisher: Cold Spring Harbor Lab.
- 852 S. Rubinacci, D. M. Ribeiro, R. J. Hofmeister, and O. Delaneau. Efficient phasing and imputa-
853 tion of low-coverage sequencing data using large reference panels. *Nature Genetics*, 53(1):120–126,
854 Jan. 2021. ISSN 1546-1718. doi: 10.1038/s41588-020-00756-0. URL [https://www.nature.com/](https://www.nature.com/articles/s41588-020-00756-0)
855 [articles/s41588-020-00756-0](https://www.nature.com/articles/s41588-020-00756-0). Number: 1 Publisher: Nature Publishing Group.
- 856 C. M. Schlebusch, H. Malmström, T. Günther, P. Sjödin, A. Coutinho, H. Edlund, A. R. Munters,
857 M. Vicente, M. Steyn, H. Soodyall, M. Lombard, and M. Jakobsson. Southern African ancient
858 genomes estimate modern human divergence to 350,000 to 260,000 years ago. *Science*, 358(6363):
859 652–655, Nov. 2017. ISSN 0036-8075, 1095-9203. doi: 10.1126/science.aao6266. URL [http://](http://science.sciencemag.org/content/358/6363/652)
860 science.sciencemag.org/content/358/6363/652.
- 861 M. Schubert, A. Ginolhac, S. Lindgreen, J. F. Thompson, K. A. AL-Rasheid, E. Willerslev, A. Krogh,
862 and L. Orlando. Improving ancient DNA read mapping against modern reference genomes. *BMC*
863 *Genomics*, 13:178, May 2012. ISSN 1471-2164. doi: 10.1186/1471-2164-13-178. URL <https://doi.org/10.1186/1471-2164-13-178>.
864
- 865 M. Schubert, S. Lindgreen, and L. Orlando. AdapterRemoval v2: rapid adapter trimming, identifica-
866 tion, and read merging. *BMC research notes*, 9(1):1–7, 2016. ISSN 1756-0500. Number: 1 Publisher:
867 BioMed Central.
- 868 L. Skotte, T. S. Korneliussen, and A. Albrechtsen. Estimating individual admixture proportions from
869 next generation sequencing data. *Genetics*, 195(3):693–702, 2013. ISSN 1943-2631. Number: 3
870 Publisher: Oxford University Press.

- 871 D.-M. J. Thorburn, K. Sagonas, M. Binzer-Panchal, F. J. J. Chain, P. G. D. Feulner, E. Bornberg-
872 Bauer, T. B. H. Reusch, I. E. Samonte-Padilla, M. Milinski, T. L. Lenz, and C. Eizaguirre.
873 Origin matters: Using a local reference genome improves measures in population genomics.
874 *Molecular Ecology Resources*, 23(7):1706–1723, 2023. ISSN 1755-0998. doi: 10.1111/1755-0998.
875 13838. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/1755-0998.13838>. eprint:
876 <https://onlinelibrary.wiley.com/doi/pdf/10.1111/1755-0998.13838>.
- 877 K. S. Toyama, P.-A. Crochet, and R. Leblois. Sampling schemes and drift can bias admix-
878 ture proportions inferred by structure. *Molecular Ecology Resources*, 20(6):1769–1785, 2020.
879 ISSN 1755-0998. doi: 10.1111/1755-0998.13234. URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/1755-0998.13234>. eprint: [https://onlinelibrary.wiley.com/doi/pdf/10.1111/1755-](https://onlinelibrary.wiley.com/doi/pdf/10.1111/1755-0998.13234)
880 [0998.13234](https://onlinelibrary.wiley.com/doi/pdf/10.1111/1755-0998.13234).
- 882 T. van der Valk, C. M. Gonda, H. Silegowa, S. Almanza, I. Sifuentes-Romero, T. B. Hart, J. A. Hart,
883 K. M. Detwiler, and K. Guschanski. The Genome of the Endangered Dryas Monkey Provides New
884 Insights into the Evolutionary History of the Vervets. *Molecular Biology and Evolution*, 37(1):183–
885 194, Jan. 2020. ISSN 0737-4038. doi: 10.1093/molbev/msz213. URL [https://doi.org/10.1093/](https://doi.org/10.1093/molbev/msz213)
886 [molbev/msz213](https://doi.org/10.1093/molbev/msz213).
- 887 E. Yüncü, U. Işıldak, M. P. Williams, C. D. Huber, L. A. Vyazov, P. Changmai, and P. Flegontov. False
888 discovery rates of qpAdm-based screens for genetic admixture. *bioRxiv*, page 2023.04.25.538339, Apr.
889 2023. doi: 10.1101/2023.04.25.538339. URL [https://www.biorxiv.org/content/10.1101/2023.](https://www.biorxiv.org/content/10.1101/2023.04.25.538339v1)
890 [04.25.538339v1](https://www.biorxiv.org/content/10.1101/2023.04.25.538339v1). Section: New Results.

Supplementary Figures

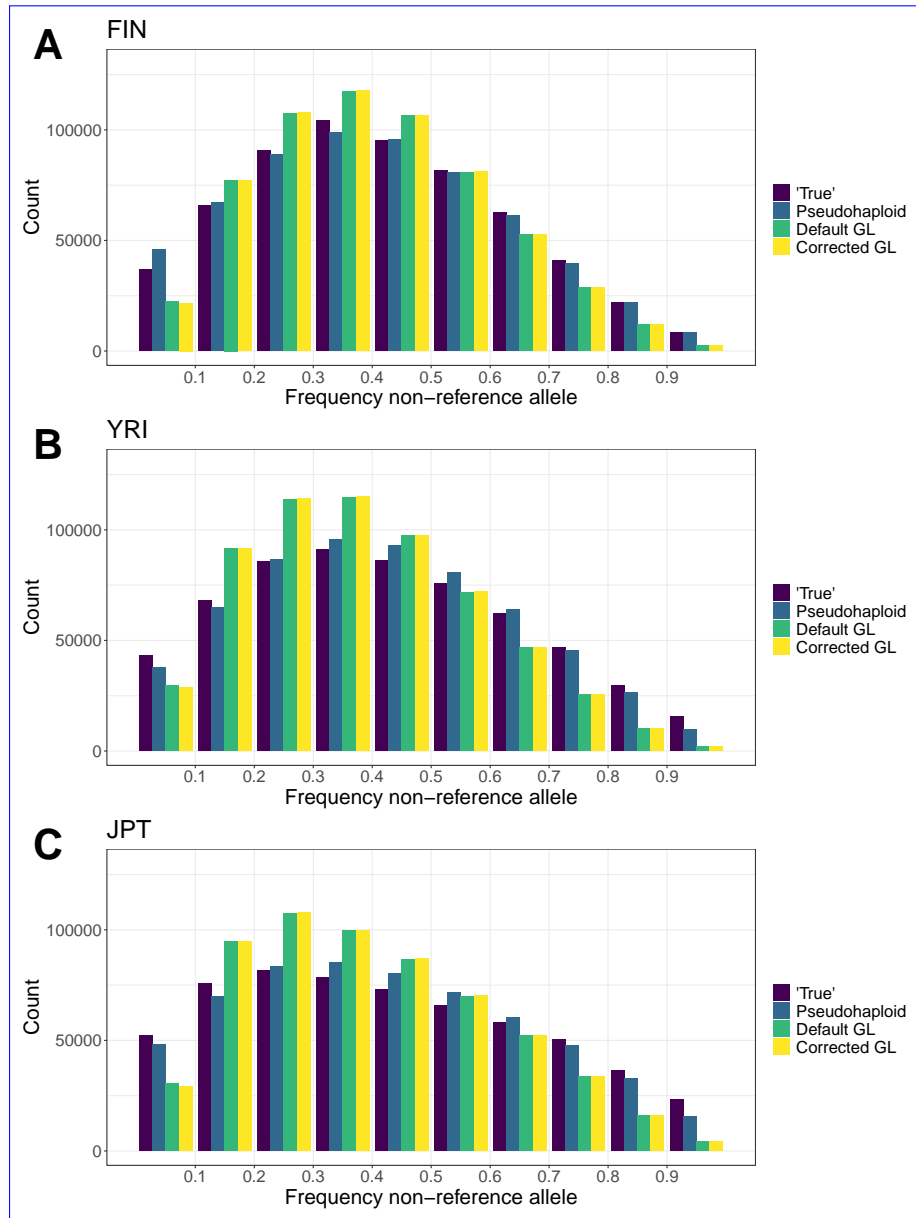


Figure S1: Binned spectrum of non-reference alleles in FIN (A), YRI (B) and JPT (C) for the four different estimation methods. Note that the specific ascertainment of common SNPs in the joint genotyping data contributes to the enrichment of variants with (true) intermediate frequencies.

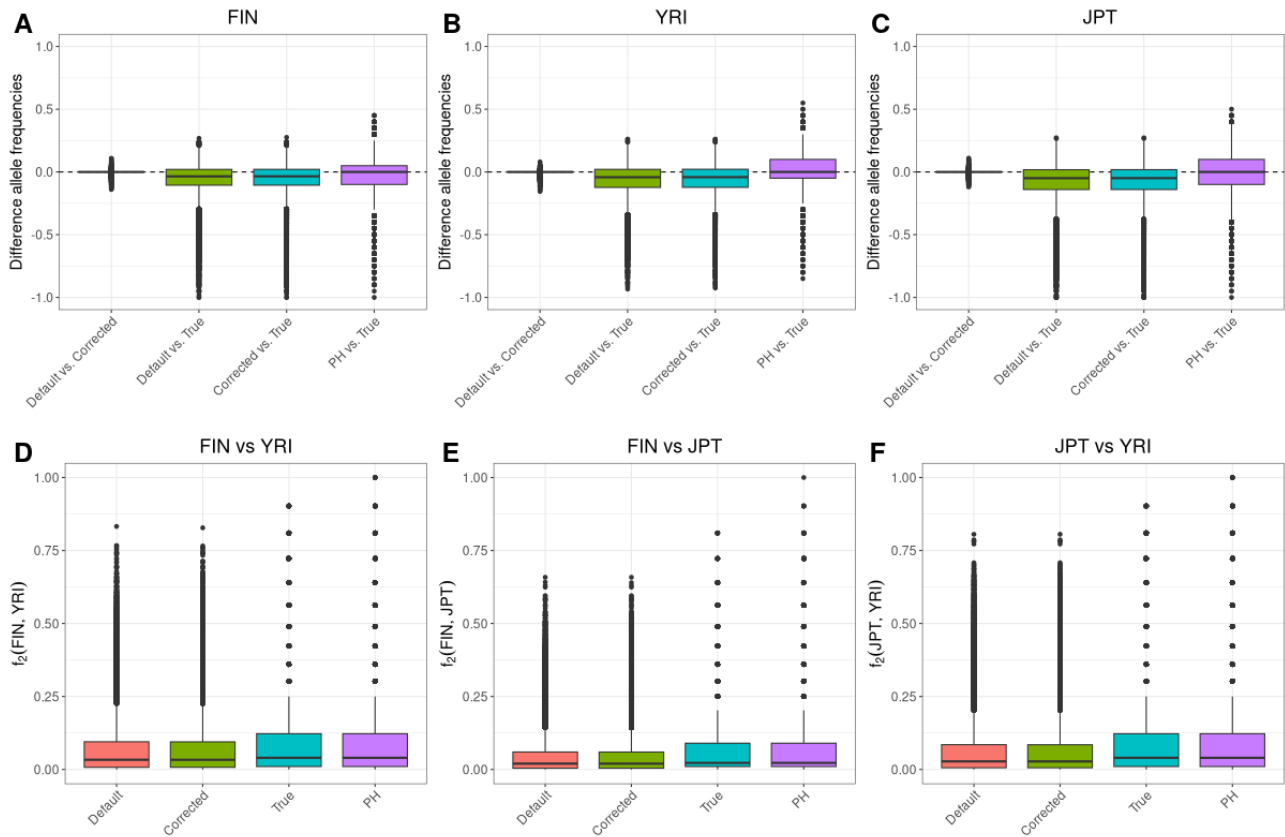


Figure S2: Differences in allele frequency estimates in the parts of the reference genome attributed to African ancestry. Boxplots for the differences between default genotype likelihood-based estimates and corrected genotype likelihood-based estimates, default genotype likelihood-based estimates and 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.

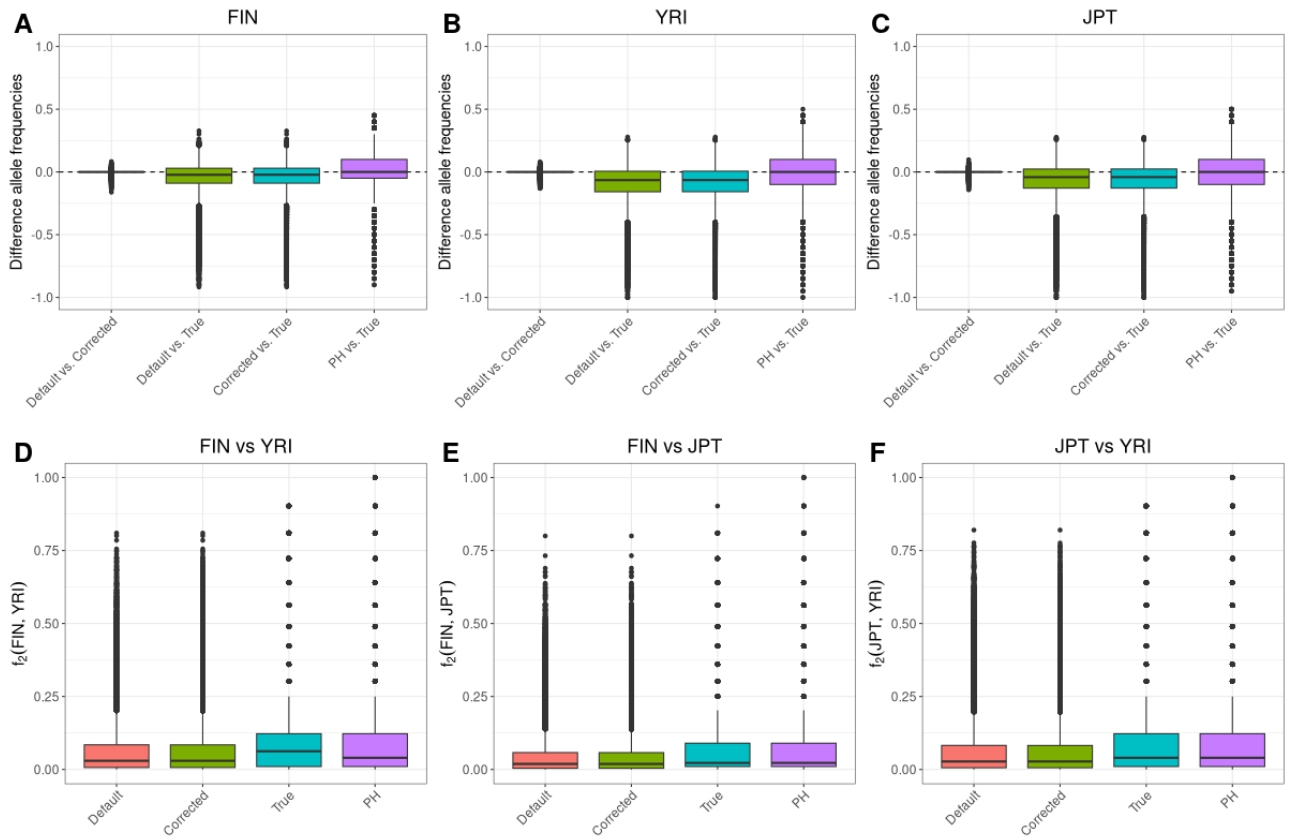


Figure S3: Differences in allele frequency estimates in the parts of the reference genome attributed to European ancestry. Boxplots for the differences between default genotype likelihood-based estimates and corrected genotype likelihood-based estimates, default genotype likelihood-based estimates and 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.

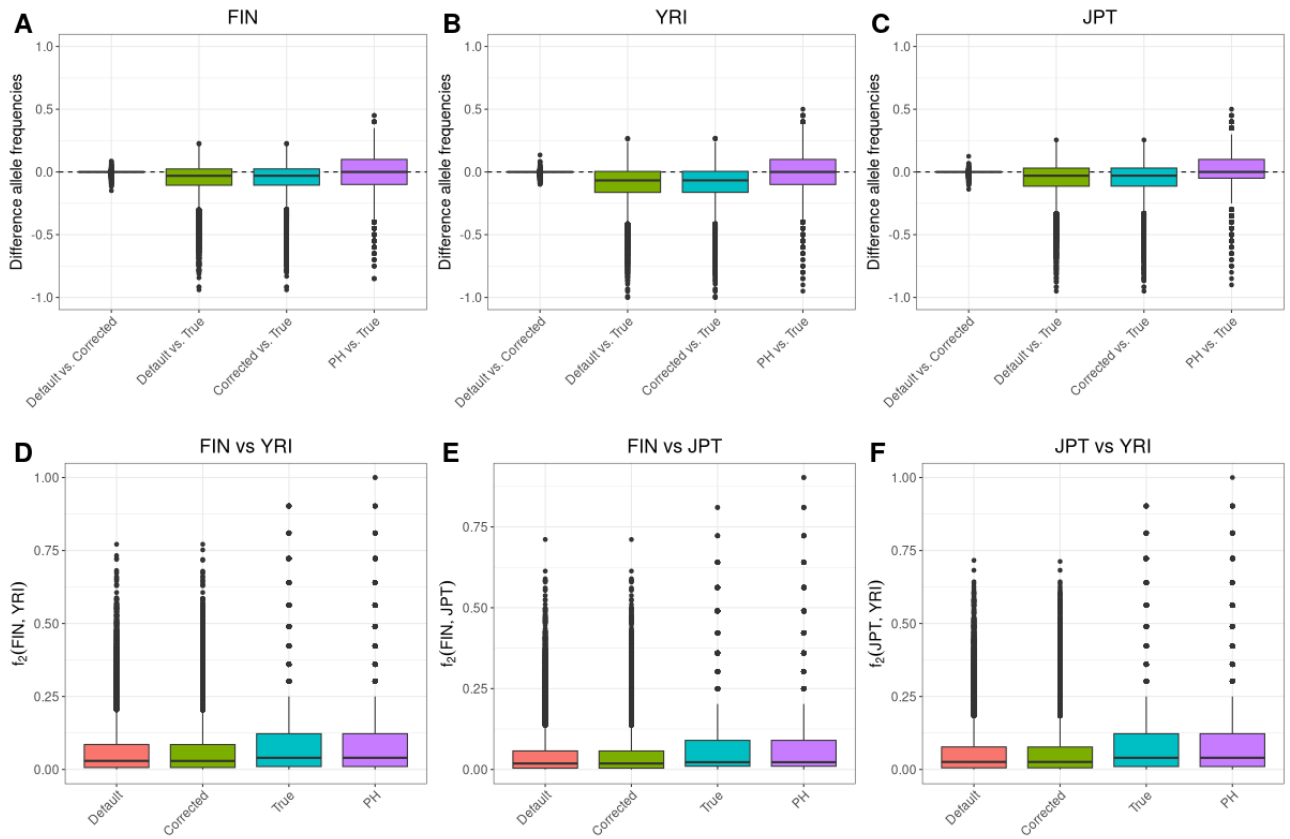


Figure S4: Differences in allele frequency estimates in the parts of the reference genome attributed to East Asian ancestry. Boxplots for the differences between default genotype likelihood-based estimates and corrected genotype likelihood-based estimates, default genotype likelihood-based estimates 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.

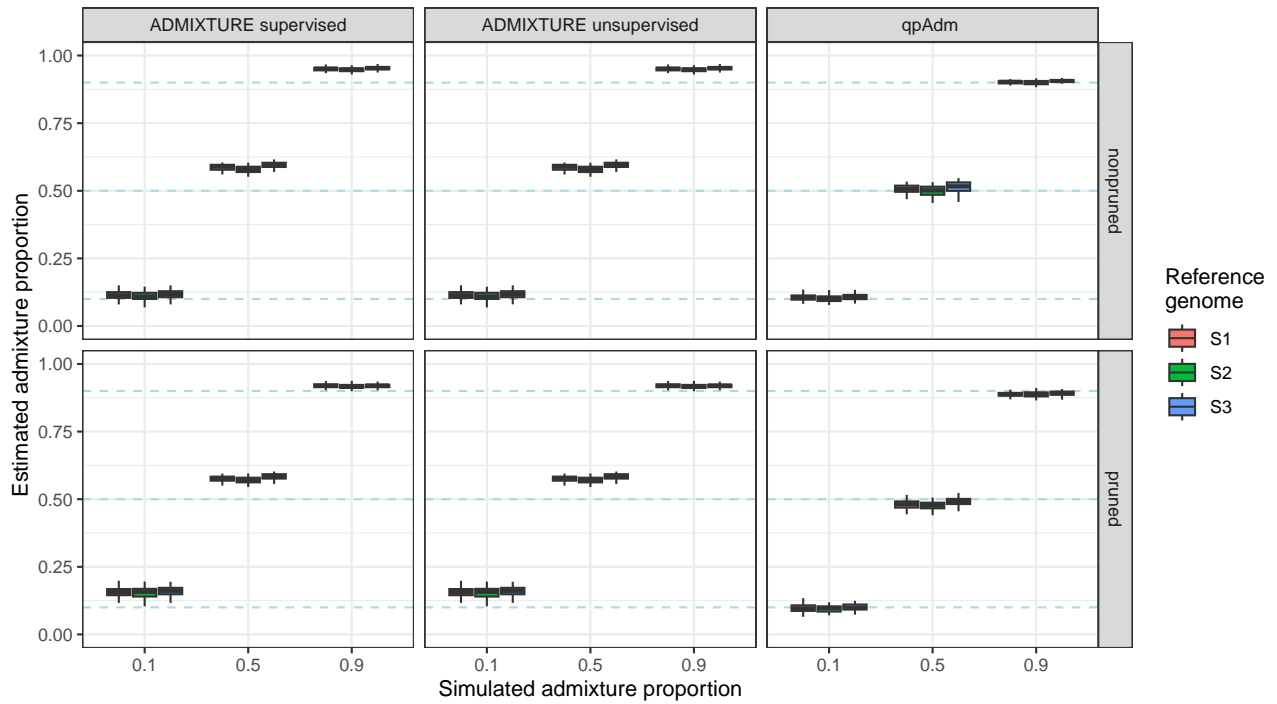


Figure S5: Simulation results for genotype call based methods using $t_{123} = 20000$ generations and a sequencing depth of 0.5X. Dashed blue lines represent the simulated admixture proportions.

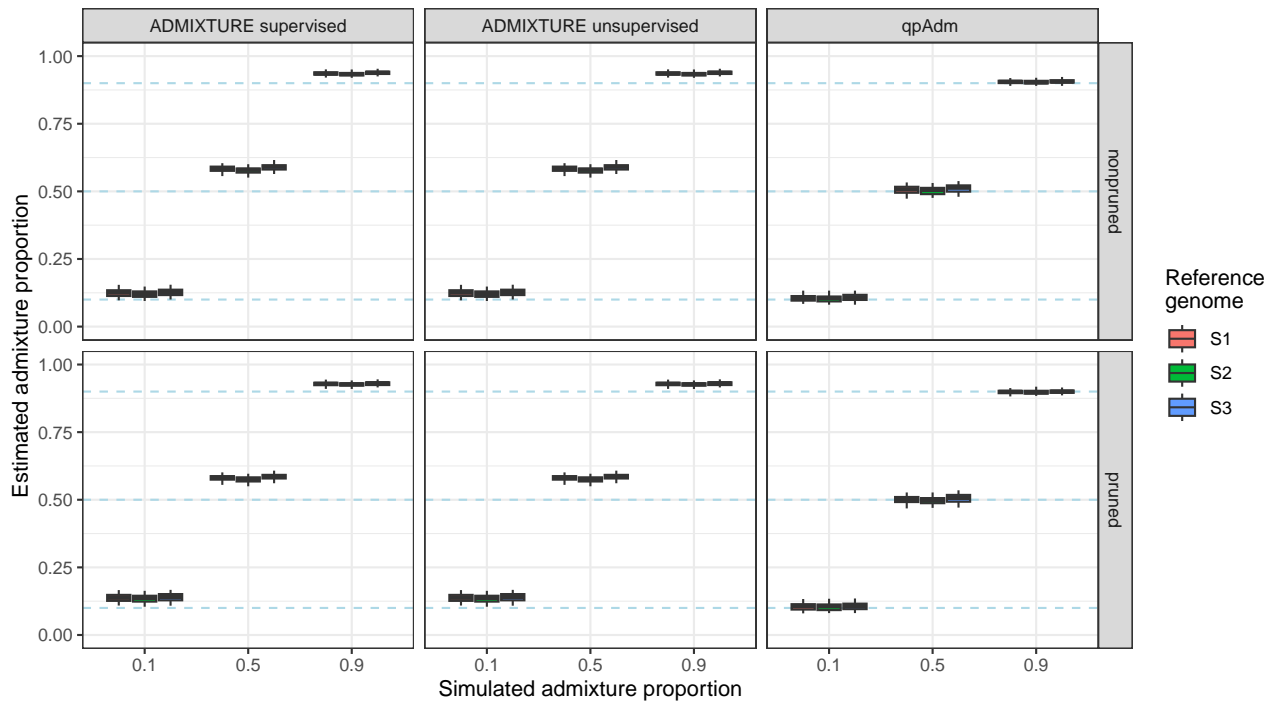


Figure S6: Simulation results for genotype call based methods using $t_{123} = 20000$ generations and a sequencing depth of 2.0X. Dashed blue lines represent the simulated admixture proportions.

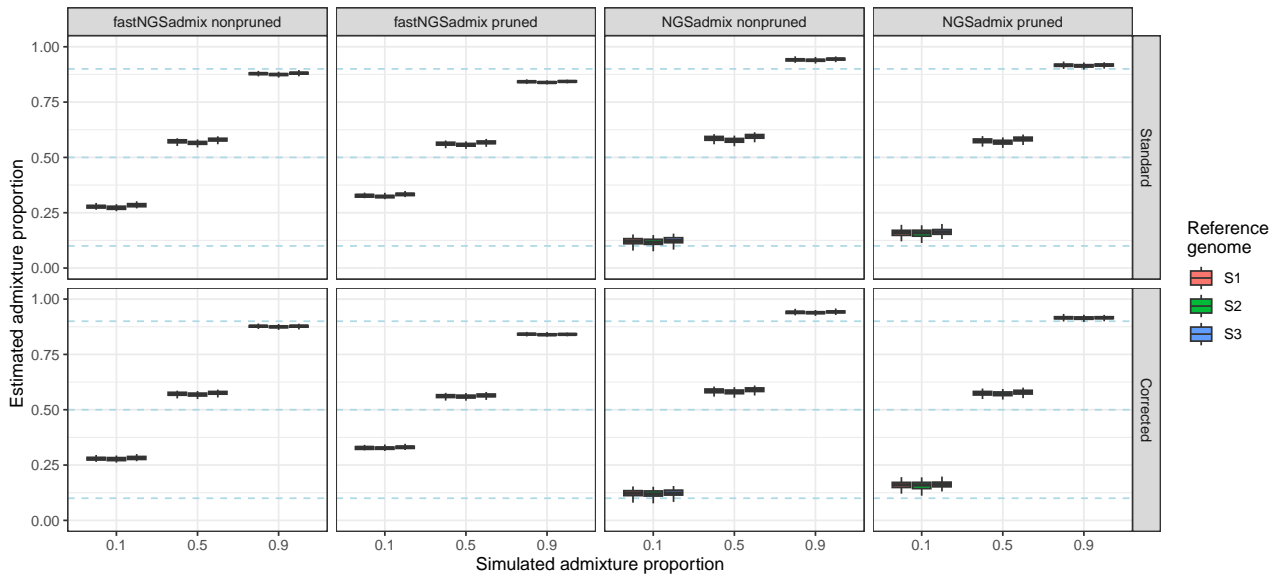


Figure S7: Simulation results for genotype likelihood based methods using $t_{123} = 20000$ generations and a sequencing depth of 0.5X. Dashed blue lines represent the simulated admixture proportions.

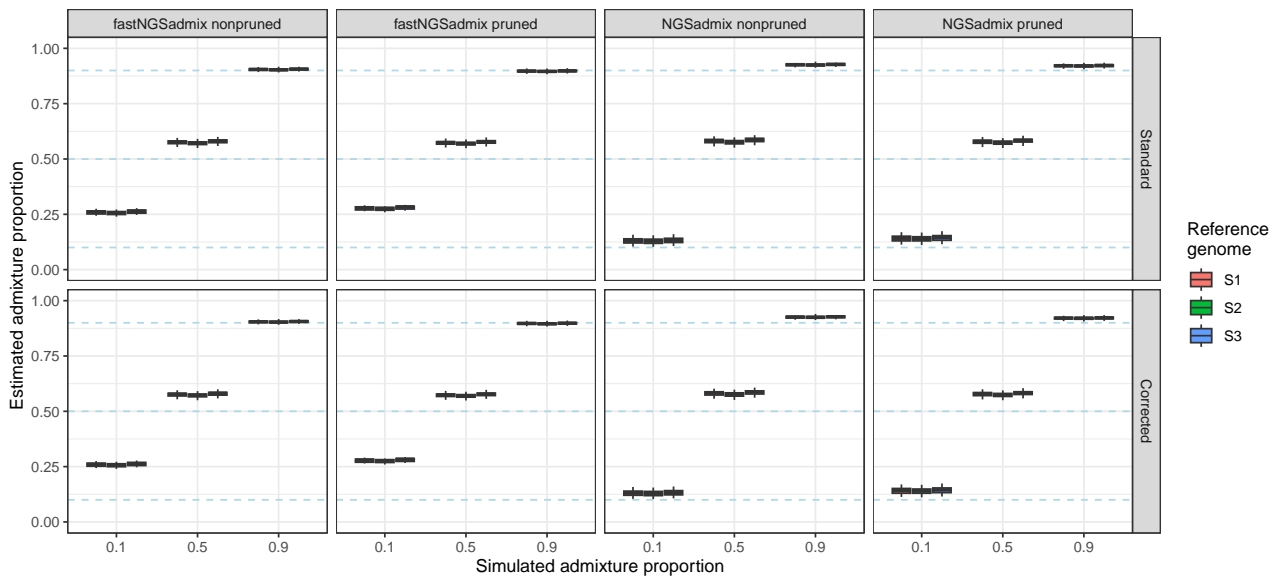


Figure S8: Simulation results for genotype likelihood based methods using $t_{123} = 20000$ generations and a sequencing depth of 2.0X. Dashed blue lines represent the simulated admixture proportions.

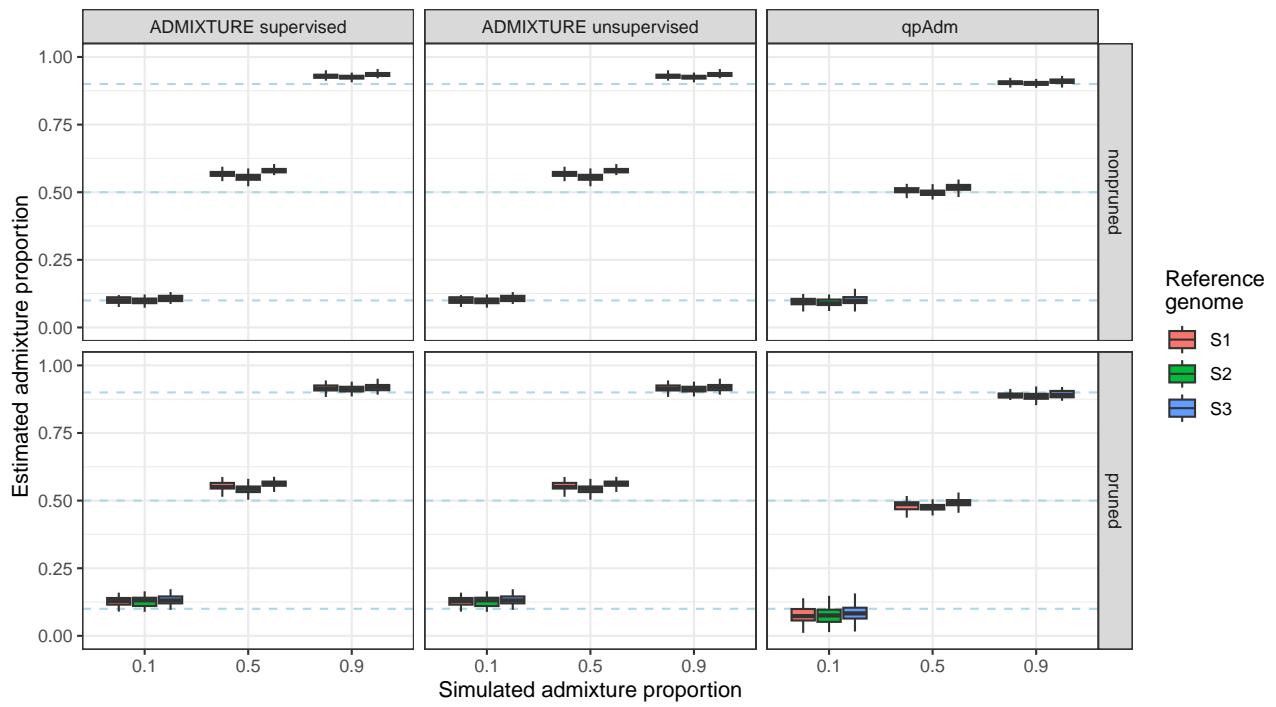


Figure S9: Simulation results for genotype call based methods using $t_{123} = 50000$ generations and a sequencing depth of 2.0X. Dashed blue lines represent the simulated admixture proportions.

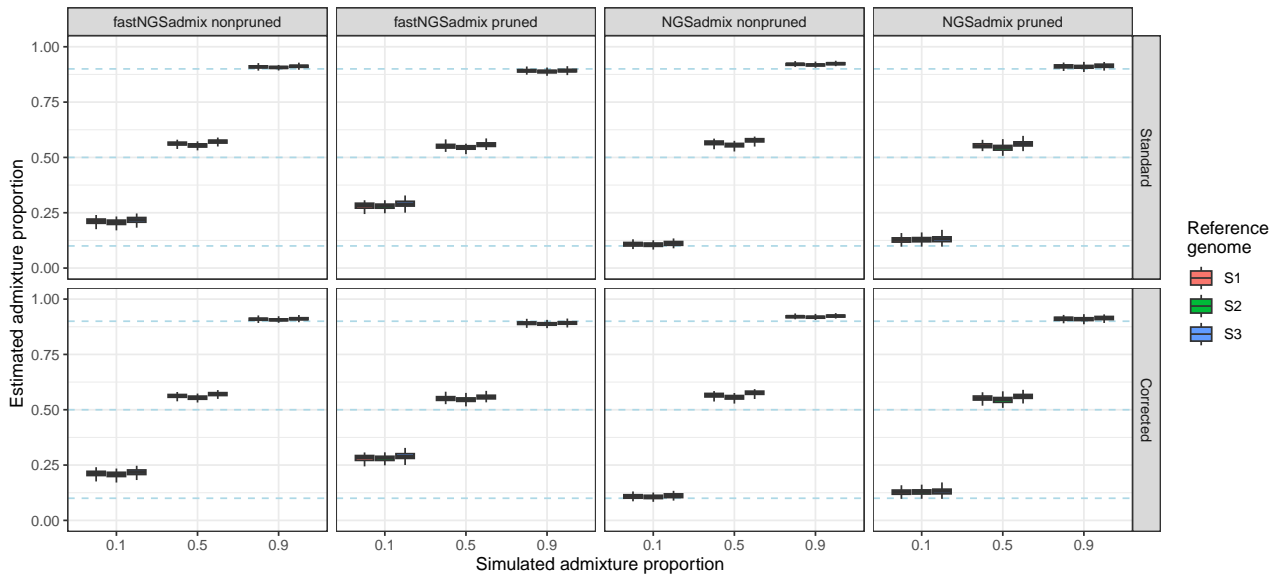


Figure S10: Simulation results for genotype likelihood based methods using $t_{123} = 50000$ generations and a sequencing depth of 2.0X. Dashed blue lines represent the simulated admixture proportions.

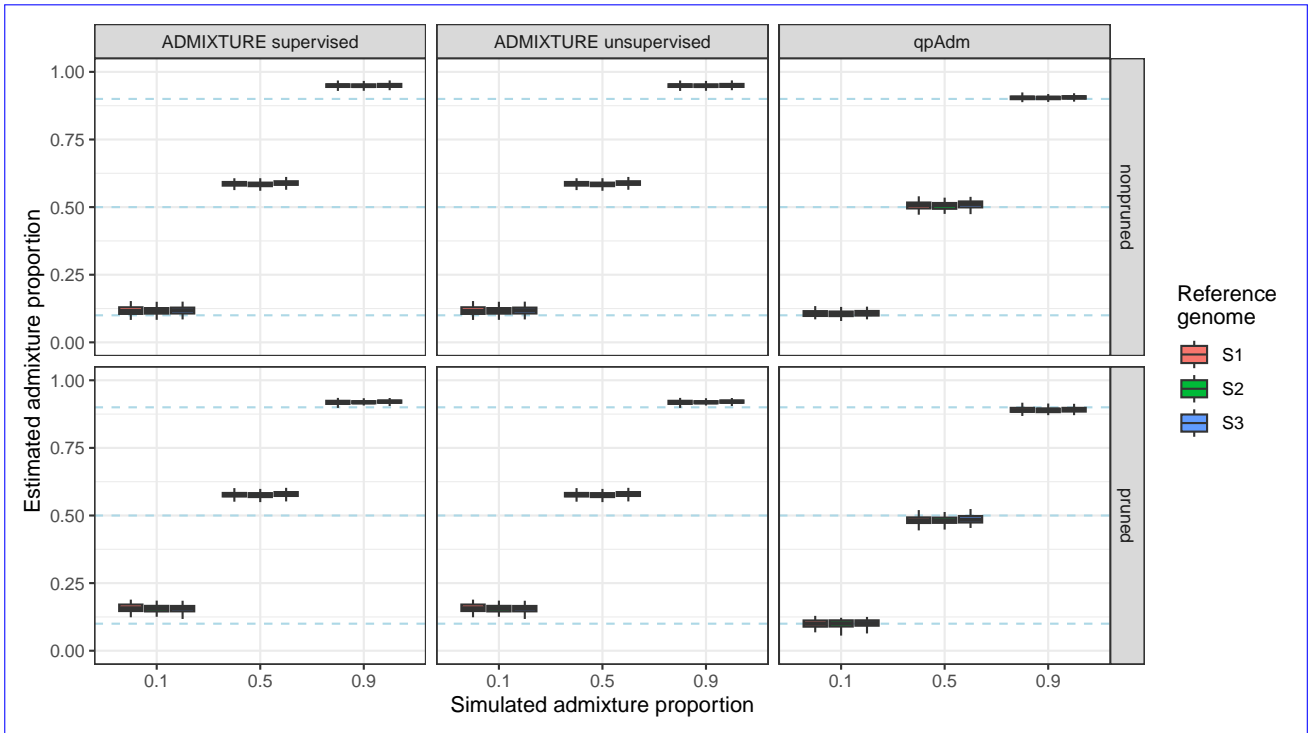


Figure S11: [Simulation results for genotype call based methods using \$t_{123} = 20000\$ generations and a sequencing depth of 0.5X. Dashed blue lines represent the simulated admixture proportions. For this run, the mapping quality threshold was set to 25 instead of 30 as in all other runs.](#)

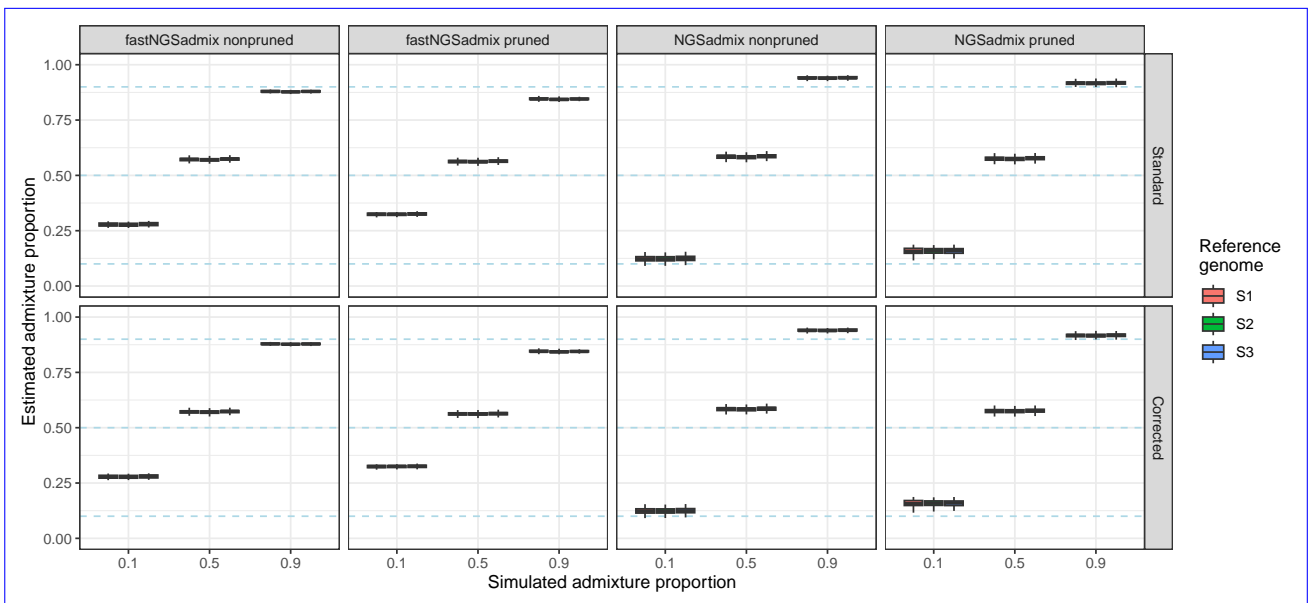


Figure S12: [Simulation results for genotype likelihood based methods using \$t_{123} = 20000\$ generations and a sequencing depth of 0.5X. Dashed blue lines represent the simulated admixture proportions. For this run, the mapping quality threshold was set to 25 instead of 30 as in all other runs.](#)

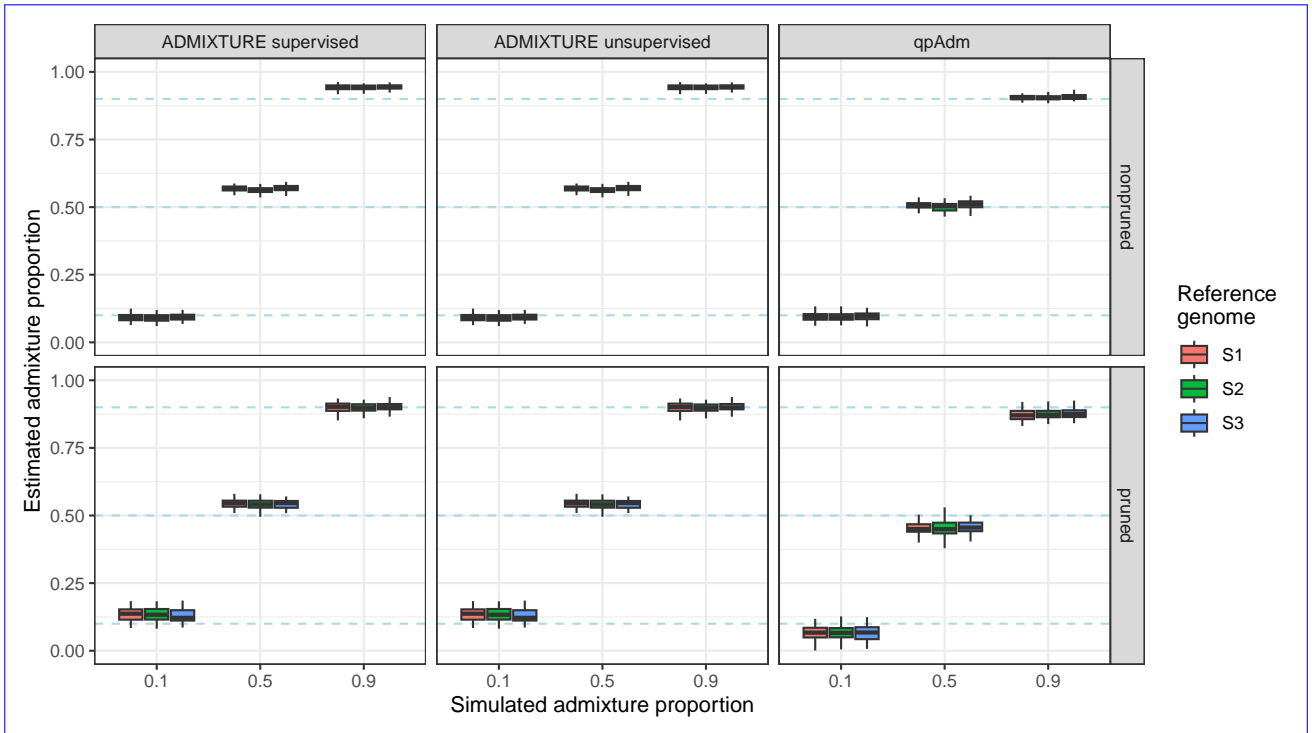


Figure S13: [Simulation results for genotype call based methods using \$t_{123} = 50000\$ generations and a sequencing depth of 0.5X. Dashed blue lines represent the simulated admixture proportions. For this run, the mapping quality threshold was set to 25 instead of 30 as in all other runs.](#)

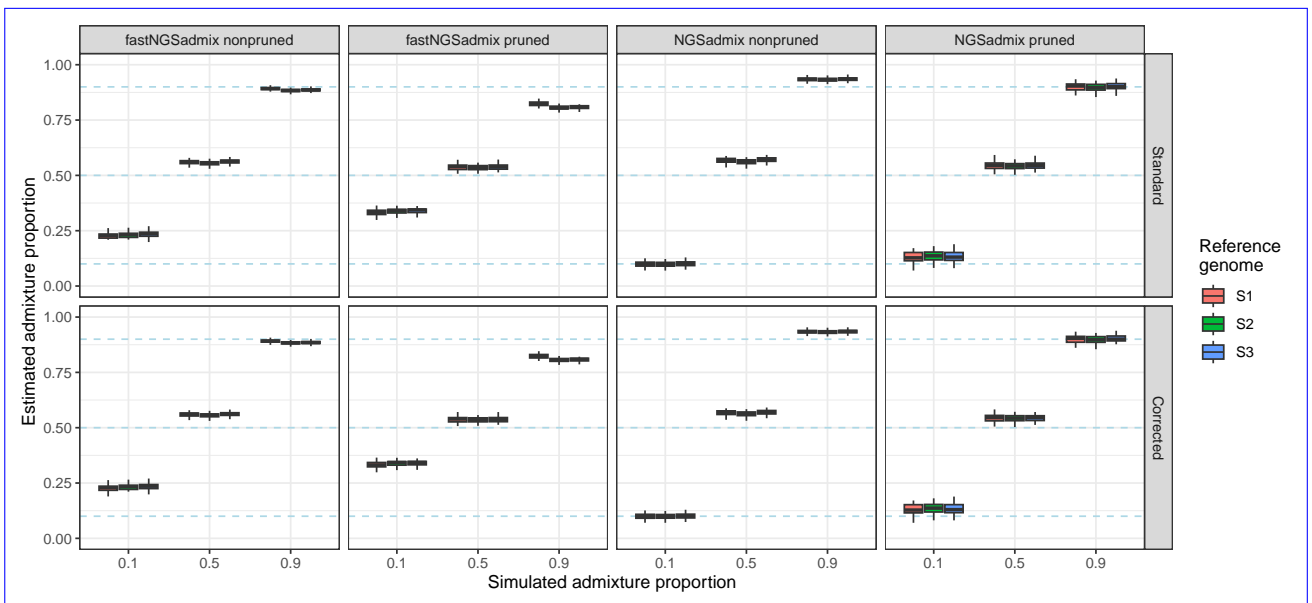


Figure S14: [Simulation results for genotype likelihood based methods using \$t_{123} = 50000\$ generations and a sequencing depth of 0.5X. Dashed blue lines represent the simulated admixture proportions. For this run, the mapping quality threshold was set to 25 instead of 30 as in all other runs.](#)

Supplementary Tables

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

height <u>individual</u>	Individual <u>Individual</u>	Population	Autosomal sequencing depth	<u>Average original read length</u>	<u>Average r_L</u>
HG00171		FIN	3.12803	108	0.5031
<u>HG00177</u>		<u>FIN</u>	<u>3.43327</u>	108	<u>0.5023</u>
<u>HG00189</u>		<u>FIN</u>	<u>3.48314</u>	108	<u>0.5026</u>
HG00190		FIN	3.089	108	0.5023
HG00272		FIN	3.61242	108	0.5027
HG00277		FIN	3.86275	76	0.5052
HG00284		FIN	4.08807	76	0.5052
HG00323		FIN	2.80008	89.19	0.5035
HG00330		FIN	13.9648	90.22	0.5045
HG00380		FIN	3.45273	100	0.502
HG00177 <u>NA18961</u>		FIN <u>JPT</u>	3.43327 <u>3.48611</u>	76	0.5067
HG00189 <u>NA18964</u>		FIN <u>JPT</u>	3.48314 <u>3.333</u>	76	0.5052
NA18853 <u>NA18969</u>		YRI <u>JPT</u>	2.56291 <u>2.6653</u>	100	0.5026
NA18923 <u>NA18970</u>		YRI <u>JPT</u>	4.42742 <u>4.47082</u>	100	0.502
NA19197 <u>NA19009</u>		<u>JPT</u>	<u>3.94626</u>	108	0.5033
<u>NA19076</u>		<u>JPT</u>	<u>3.50604</u>	108	0.5029
<u>NA19080</u>		<u>JPT</u>	<u>3.84401</u>	108	0.5055
<u>NA19081</u>		<u>JPT</u>	<u>2.60827</u>	108	0.5034
<u>NA19082</u>		<u>JPT</u>	<u>3.58866</u>	108	0.5018
<u>NA19084</u>		<u>JPT</u>	<u>4.37475</u>	108	0.5026
<u>NA18520</u>		YRI	<u>4.19443</u> <u>3.99207</u>	76	0.5057
NA19200 <u>NA18522</u>		YRI	4.22902 <u>2.55368</u>	76	0.5066
NA19236 <u>NA18853</u>		YRI	4.21535 <u>2.56291</u>	76	0.5099
NA19248 <u>NA18923</u>		YRI	4.24979 <u>4.42742</u>	100	0.5019
NA19116		YRI	3.03829	82.51	0.5056
NA19130		YRI	4.97799	76	0.5061
NA18520 <u>NA19197</u>		YRI	3.99207 <u>4.19443</u>	100	0.5021
NA18522 <u>NA19200</u>		YRI	2.55368 <u>4.22902</u>	100	0.502
<u>NA19236</u>		YRI	<u>4.21535</u>	76	0.5055
<u>NA19248</u>		YRI	<u>4.24979</u>	76	0.5058

Table S2: ~~Total number and percentage~~ Average read balances for the 1000 genomes populations used for the analysis of SNPs with extreme differences ($\geq |0.5|$) between "True" and estimated allele frequencies ~~empirical data.~~

height <u>Population</u>	True vs default GL <u>True vs. corrected GL</u>	True vs. Pseudohaploid <u>Average r_L</u>
FIN	738 (0.118%)608 (0.096%)	<u>0.50334</u>
<u>JPT</u>	979 (0.157%)	<u>0.5036</u>
YRI	829 (0.133%)674 (0.108%)947 (0.152%)	<u>0.50512</u>
height		