

# Modeling biases from low-pass genome sequencing to enable accurate population genetic inferences

Emanuel M. Fonseca<sup>\*</sup>, Linh N. Tran, Hannah Mendoza, Ryan N. Gutenkunst<sup>\*</sup>

Department of Molecular and Cellular Biology, University of Arizona, Tucson, AZ 85721, USA

**\*Corresponding author:** E-mail: [emanuelfonseca@arizona.edu](mailto:emanuelfonseca@arizona.edu); [rgutenk@arizona.edu](mailto:rgutenk@arizona.edu)

The correction for low-pass sequencing is performed using the publicly available dadi Python package, which can be accessed at <https://bitbucket.org/gutenkunstlab/dadi>. Additionally, the codebase for creating and analyzing both simulated and empirical datasets, ensuring reproducibility, is readily accessible on GitHub at <https://github.com/emanuelfonseca/low-coverage-sfs> and [https://github.com/lntan26/low-coverage-sfs/tree/main/empirical\\_analysis](https://github.com/lntan26/low-coverage-sfs/tree/main/empirical_analysis). Furthermore, we provide illustrative examples to assist users in implementing our methodology.

## Abstract

Low-pass genome sequencing is cost-effective and enables analysis of large cohorts. However, it introduces biases by reducing heterozygous genotypes and low-frequency alleles, impacting subsequent analyses such as demographic history inference. We developed a probabilistic model of low-pass biases from the Genome Analysis Toolkit (GATK) multi-sample calling pipeline, and we implemented it in the population genomic inference software dadi. We evaluated the model using simulated low-pass datasets and found that it alleviated low-pass biases in inferred demographic parameters. We further validated the model by downsampling 1000 Genomes Project data, demonstrating its effectiveness on real data. Our model is widely applicable and substantially improves model-based inferences from low-pass population genomic data.

**Key words:** demography inference, inbreeding, low-pass sequencing, allele frequency spectrum, GATK multi-sample calling

## Introduction

Enabled by reduced sequencing costs, population genetics has experienced a revolution, from focusing on a limited number of loci to now encompassing entire genomes (Maddison et al. 1992; Reid et al. 2016; Marchi et al. 2022). Yet researchers must still trade off a) the extent of the genome to be sequenced, b) the depth of coverage for each sample, and c) the number of sequenced samples (Lou et al. 2021; Martin et al. 2021; Duckett et al. 2023). One way to address this trade off is to sequence one reference sample at high coverage depth while sequencing others at lower depth (Lou et al. 2021). Low-pass sequencing, in which the genome is sequenced at a lower depth of coverage, avoids many of the financial, methodological, and computational challenges of high-pass sequencing (Li et al. 2011). Furthermore, limited availability of DNA can also make high depth impractical, especially for ancient samples and museum or herbarium specimens (Mota et al. 2023).

Despite its advantages, low-pass sequencing may lead to an incomplete and biased representation of genetic diversity within a population (e.g., Vieira et al. 2013; Fox et al. 2019). Low-frequency genomic variants may not be detected (Fumagalli 2013), and genotypes may be less accurate (Nielsen et al. 2011). Low-pass sequencing increases the likelihood of miscalling heterozygous loci as homozygous (Duitama et al. 2011; Gorjanc et al. 2015), due to a lack of sufficient reads on homologous chromosomes to distinguish between different alleles at a given locus. These issues can then bias downstream analyses. It is thus important for analysis methods to accommodate low-pass sequencing (see Carstens et al. 2022 for a discussion of related issues).

The allele frequency spectrum (AFS) is a powerful summary of population genomic data (Sawyer & Hartl 1992; Wakeley 2009). Briefly, the AFS is matrix which records the number alleles observed at given

43 frequencies in a sample of individuals from one or more populations. The AFS is often the basis for inferring  
44 demographic history (Gutenkunst et al. 2009) or distributions of fitness effects (Kim et al. 2017). In low-pass  
45 sequencing, the loss of alleles and the excess of homozygosity can bias the estimation of the AFS (Fumagalli  
46 2013) and thus those inferences.

47 To address the challenges of low-pass data, several tools have emerged (Bryc et al. 2013; Blischak et al.  
48 2018; Meisner & Albrechtsen 2018), with one of the most widely adopted being ANGSD (Korneliussen  
49 et al. 2014). ANGSD offers a diverse range of analyses tailored for low-pass sequencing data. To infer an  
50 AFS, ANGSD uses sample allele frequency likelihoods, which can be computed either directly from raw  
51 data or, more frequently, from genotype likelihoods (Nielsen et al. 2012). These likelihoods quantify the  
52 probability of observing the complete set of read data for multiple individuals at specific genomic sites,  
53 given particular sample allele frequencies (Nielsen et al. 2012; Korneliussen et al. 2014), enabling ANGSD  
54 to estimate allele frequencies. While ANGSD has proven its utility, limitations exist. For example, many  
55 analyses rely on distinguishing different types of variant sites (such a synonymous versus nonsynonymous)  
56 which the developers of ANGSD recommend against. Moreover, in some cases unbiased estimation of the  
57 AFS may be difficult or impossible.

58 Rather than attempting to estimate an unbiased AFS from low-pass data, we developed a proba-  
59 bilistic model of low-pass AFS biases and incorporated it into the population genomic inference soft-  
60 ware dadi (Gutenkunst et al. 2009). Our model is based on the multi-sample genotype calling pipeline  
61 of the Genome Analysis Toolkit (GATK), the most widely used tool for calling variants from read data  
62 (McKenna et al. 2010; Auwera & O’Connor 2020). We assessed the accuracy of our model using sim-  
63 ulated low-depth data as well as subsampled data from the 1000 Genomes Project (Fairley et al. 2020,  
64 <https://www.internationalgenome.org/>). We found that our model accurately captures low-depth biases in  
65 the AFS and enables accurate inference of demographic history from low-pass data.

## 66 Model for Low-pass Biases

67 When biases arises from low-pass sequencing, the AFS may be affected by both the loss of low-frequency  
68 variants and the misidentification of heterozygous individuals as homozygous. These two effects result in  
69 a deficit of variant sites and misleading shifts in allele frequencies, respectively. Moreover, the data must  
70 often be subsampled to generate an AFS for analysis, because not all individuals will be called at all sites.  
71 We account for these distortions by sequentially modeling the probabilities of a variable site being called, of  
72 that site having enough called individuals for subsampling, and of having its allele frequency misestimated.

73 The specific choices in our model are motivated by the default GATK multi-sample calling algorithm, in  
74 which information from all samples is used to identify whether a site is variant. In particular, we assume  
75 that a site will only be called as variant if at least two alternate allele reads are observed. Once a site is  
76 identified as variant, an individual will be called as missing if zero reads are observed, homozygous if all  
77 reads correspond to a single allele, and heterozygous if at least one reference and one alternate read are  
78 observed. For simplicity, we first describe the case of sequencing  $n_{seq}$  individuals from a single population.

79 Consider a site in which the true alternate allele count within our sample of  $n_{seq}$  individuals is  $f$ . Those  
80  $f$  alternate alleles can be distributed among the  $2n_{seq}$  sampled alleles in many ways. To quantify those ways,  
81 we define the partition function  $\mathbb{P}_{n_{seq}}(f)$ , which is an array of integer partitions with  $n$  entries that sum  
82 to the allele frequency  $f$  such that all entries in the partition are 0, 1, or 2 (corresponding to the possible  
83 genotype values). For example, the partitions defined by  $\mathbb{P}_4(3)$  are  $[2, 1, 0, 0]$  and  $[1, 1, 1, 0]$ . Each possible  
84 partition within  $\mathbb{P}_{n_{seq}}(f)$  can occur in  $\frac{n!}{n_0!n_1!n_2!}2^{n_1}$  ways, where  $n_0$ ,  $n_1$ , and  $n_2$  denote the number of partition  
85 entries equal to 0, 1, or 2. (The factor of  $2^{n_1}$  accounts for the two possible haplotypes the alternate allele  
86 could lie on in each heterozygote.) The corresponding probability of each partition within  $\mathbb{P}_{n_{seq}}(f)$  is then  
87 the number of ways it can occur divided by the total over all partitions within  $\mathbb{P}_{n_{seq}}(f)$ .

88 Let  $\mathbb{D}$  denote the distribution of read depth  $d$  within the population sample, which we assume to be shared  
89 among all individuals. For an individual homozygous for the alternate allele, the probability of observing  $a$   
90 alternate reads is simply  $P_a^{hom}(a) = \mathbb{D}(a)$ . For a heterozygous individual, the probability of zero alternate

91 reads is

$$P_a^{het}(0) = \sum_d \mathbb{D}(d) \left(\frac{1}{2}\right)^d. \quad (1)$$

92 Here we sum over the distribution of depths, and at each depth each read has a  $1/2$  chance of containing the  
93 reference allele, so the probability of all reads being reference is  $(1/2)^d$ . Similarly, the probability of exactly  
94 one alternate read is

$$P_a^{het}(1) = \sum_d d \mathbb{D}(d) \left(\frac{1}{2}\right) \left(\frac{1}{2}\right)^{d-1}. \quad (2)$$

95 Note that for depth  $d$ , there are  $d$  possible configurations with one alternate read and  $d - 1$  reference reads.

96 For a given partition within  $\mathbb{P}_{n_{seq}}(f)$  that has true genotype counts  $n_0$ ,  $n_1$ , and  $n_2$ , there are multiple  
97 ways of failing to identify the variant site. The probability of zero reads supporting the alternate allele is

$$P_a^{part}(0) = P_a^{het}(0)^{n_1} P_a^{hom}(0)^{n_2}. \quad (3)$$

98 The probability of exactly one read supporting the alternate allele is

$$P_a^{part}(1) = n_1 P_a^{het}(1) P_a^{het}(0)^{n_1-1} P_a^{hom}(0)^{n_2} + P_a^{het}(0)^{n_1} n_2 P_a^{hom}(1) P_a^{hom}(0)^{n_2-1}. \quad (4)$$

99 Here the two terms account for the probability that the alternate read occurs in one of the heterozygotes or  
100 homozygotes, respectively. The overall probability of not calling a variant site for a given partition is thus  
101  $P_a^{part}(0) + P_a^{part}(1)$ . And the overall probability of not calling a variant site with a given true allele frequency  
102  $f$  is the sum of these probabilities over partitions  $\mathbb{P}_{n_{seq}}(f)$ , weighted by the partition probabilities. For any  
103 given coverage distribution, the probability of calling a variant site increases rapidly with allele frequency  $f$   
104 (Fig. S1).

105 When analyzing low-pass data, generating an AFS for the full sample size  $n_{seq}$  may result in the loss  
106 of many sites where not all individuals were called. Consequently, it is common to subsample the data to  
107 some lower sample size  $n_{sub}$ ; only sites with calls for at least  $n_{sub}$  individuals can then be analyzed. The  
108 probability a site can be analyzed is independent of the allele frequency and is

$$\sum_{c=n_{sub}}^{n_{seq}} \frac{n_{seq}!}{c!(n_{seq}-c)!} \mathbb{D}(0)^{n_{seq}-c} (1 - \mathbb{D}(0))^c. \quad (5)$$

109 Here we sum the probability that exactly  $c$  individuals have at least one read at this site over all potential  
110 values of the number of covered individuals  $n_{sub} \leq c \leq n_{seq}$ . From this point onward, we consider partitions  
111  $\mathbb{P}_{n_{sub}}(f)$  over the subsampled individuals.

112 Once a site is called as variant, low-pass sequencing can bias the estimation of the allele frequency at  
113 that site, if one or more heterozygotes are miscalled because all their reads are reference or alternate. For  
114 each heterozygous individual, this occurs with total probability

$$P_{mis}^{het} = 2 \sum_d \mathbb{D}(d) \left(\frac{1}{2}\right)^d = 2 P_a^{het}(0). \quad (6)$$

115 For a partition with  $n_1$  true heterozygotes, the number of miscalled heterozygotes  $N_{mis}^{het}$  is binomially  
116 distributed with mean  $n_1 P_{mis}^{het}$ . Each miscalled heterozygote has equal chance of being called as homozygous  
117 reference or alternate, so the number of miscalls to homozygous reference  $N_{\rightarrow ref}^{het}$  is binomially distributed  
118 with mean  $N_{mis}^{het}/2$ , and the number of miscalls to homozygous alternate is  $N_{\rightarrow alt}^{het} = N_{mis}^{het} - N_{\rightarrow ref}^{het}$ . The net  
119 change in estimated alternative allele frequency is then  $N_{\rightarrow alt}^{het} - N_{\rightarrow ref}^{het}$ .

120 The biases caused by low-pass sequencing do not depend on the underlying AFS; for each true allele  
121 frequency a given fraction will always, on average, be miscalled as any given other allele frequency. The  
122 correction above can be thus be calculated once for a given data set then applied to all model AFS generated,  
123 for example, during demographic parameter optimization. For efficiency, we calculate and cache an  $n_{seq}$  by

124  $n_{sub}$  transition matrix that can be multiplied by any given model AFS for  $n_{seq}$  individuals to apply the low  
125 coverage correction. When analyzing multiple populations, we calculate and apply transitions matrices for  
126 each population, because variant calling is independent among populations once a variant has been identified.  
127 Variant identification is, however, not independent among populations, which we address using simulated  
128 calling described next.

129 When calculating the probability of miscalling a heterozygote (Eq. 6), the correct distribution of depth  
130 is not simply  $\mathbb{D}(d)$ ; rather it is the distribution conditional on the site being identified as variant. The lower  
131 the true allele frequency, the more these distributions will differ. The conditional distribution is complex to  
132 calculate, particularly when multiple populations are involved. Instead, for true allele frequencies for which  
133 the probability of not identifying is above a user-defined threshold (by default  $10^{-2}$ ), we simulate the calling  
134 process rather than using our analytic results. For multiple populations, we calculate this threshold assuming  
135 that a variant must be identified independently in all populations, which gives a lower bound on the true  
136 probability of not identifying. To simulate calling, for a given true allele frequency (or combination in the  
137 multi-population case) we simulate reads (default 1000) using the coverage distribution  $\mathbb{D}(d)$  and simulate  
138 variant identification and genotype calling for each potential partition of genotypes across the populations,  
139 proportional to its probability. For each combination of input true allele frequencies simulated, we estimate  
140 and store probability of each potential output allele frequency. These distortions are then applied in place  
141 of the transition matrices from the analytic model.

142 For inbred populations, there is an excess of homozygotes compared to the Hardy-Weinberg expectation,  
143 which reduces biases associated with low-pass sequencing. In this case, we follow [Blischak et al. \(2020\)](#)  
144 and within each genotype partition calculate the probability of reference homozygotes, heterozygotes, and  
145 alternate homozygotes using results from [Balding & Nichols \(1995, 1997\)](#), given the inbreeding coefficient  $F$ .  
146 The partition probability is then multinomial given these probabilities. In these calculations, we approximate  
147 the population allele frequency by the true sample allele frequency. Because calculation of the low-pass  
148 correction is expensive compared to typical normal model AFS calculation, we pre-calculate and cache  
149 transition matrices and calling simulations. But inbreeding is often an inferred model parameter, to be  
150 optimized during analysis. In this case, users can specify an assumed inbreeding parameter for the low-pass  
151 model, optimize the inbreeding parameter in their demographic model, update the inbreeding coefficient  
152 assumed in the low-pass model, and iterate until convergence.

## 153 Results

### 154 Low-pass sequencing biases the AFS

155 We used simulated data to assess the biases introduced by low-pass sequencing with GATK multi-sample  
156 calling, along with our model of those biases. For a simulated population undergoing growth (Fig. S2A),  
157 low-pass sequencing reduces the number of observed low-frequency alleles (Fig. 1). Our model accurately  
158 captures these biases (Fig. 1). In contrast with our model, ANGSD attempts to reconstruct the true AFS  
159 from low-pass data. In our simulations, ANGSD reconstructed the mean shape of the AFS well, but it  
160 introduced dramatic fluctuations into the reconstructed AFS at low depth (Fig. 2).

161 When a pair of populations undergoing a split and isolation (Fig. S2B) is analyzed through a joint AFS,  
162 similar low coverage biases occur (Fig. S3). Again, our model corrects those biases well (Fig. S3). Similar  
163 to the single-population case, ANGSD also introduces large fluctuations in the joint AFS S4).

164 Low-pass biases are expected to be smaller in inbred populations, due to the reduction of heterozygosity.  
165 In a simulated population recovering from a bottleneck with inbreeding (Fig. S2C), biases are still observed,  
166 which our model corrects (Fig. S5). Again, ANGSD introduced large fluctuations in low-pass AFS, beyond  
167 those expected from inbreeding (Fig. S6).

### 168 Demographic history inference from low-pass AFS

169 To assess effects on inference, we first fit demographic models to single-population data simulated under  
170 the same growth model as our prior simulations (Fig. S2A). When not modeling low-pass biases, the final

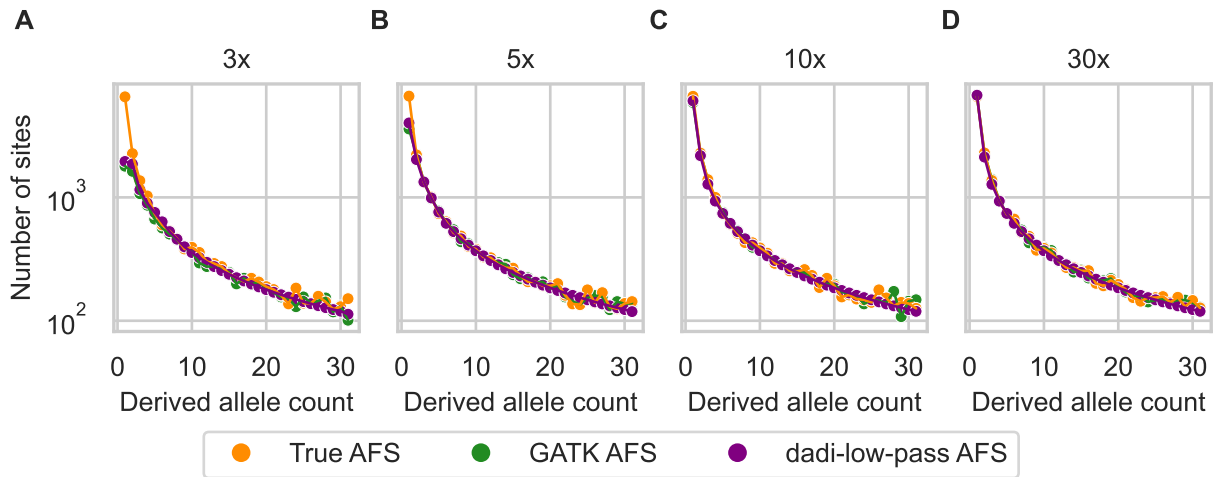


Figure 1: The low-pass AFS is biased, which our model captures. Simulated sequence data from an exponential growth demographic model for 20 individuals were called by GATK and subsampled to 16 individuals (to accommodate missing data at low depth). The GATK-called AFS (green) is biased compared to the true AFS (orange), and our dadi model for low-pass sequencing (purple) fits those biases well. Coverage was (A) 3 $\times$ , (B) 5 $\times$ , (C) 10 $\times$ , and (D) 30 $\times$ .

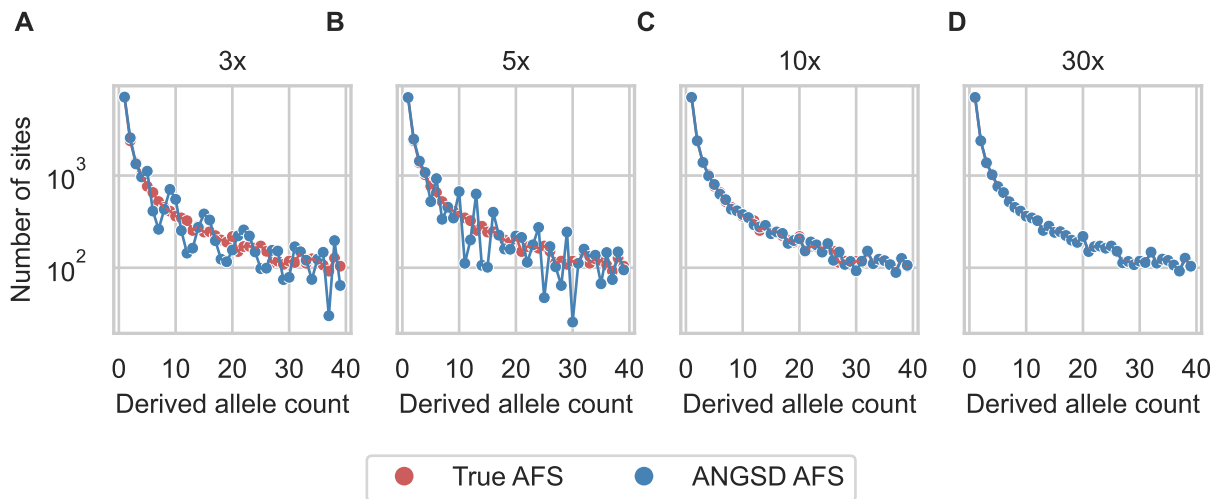


Figure 2: ANGSD corrects for low-pass bias of the AFS, but introduces fluctuations. For the same simulations as Fig. 1, ANGSD (blue) was used to reconstruct the true AFS (red). Coverage was (A) 3 $\times$ , (B) 5 $\times$ , (C) 10 $\times$ , and (D) 30 $\times$ .

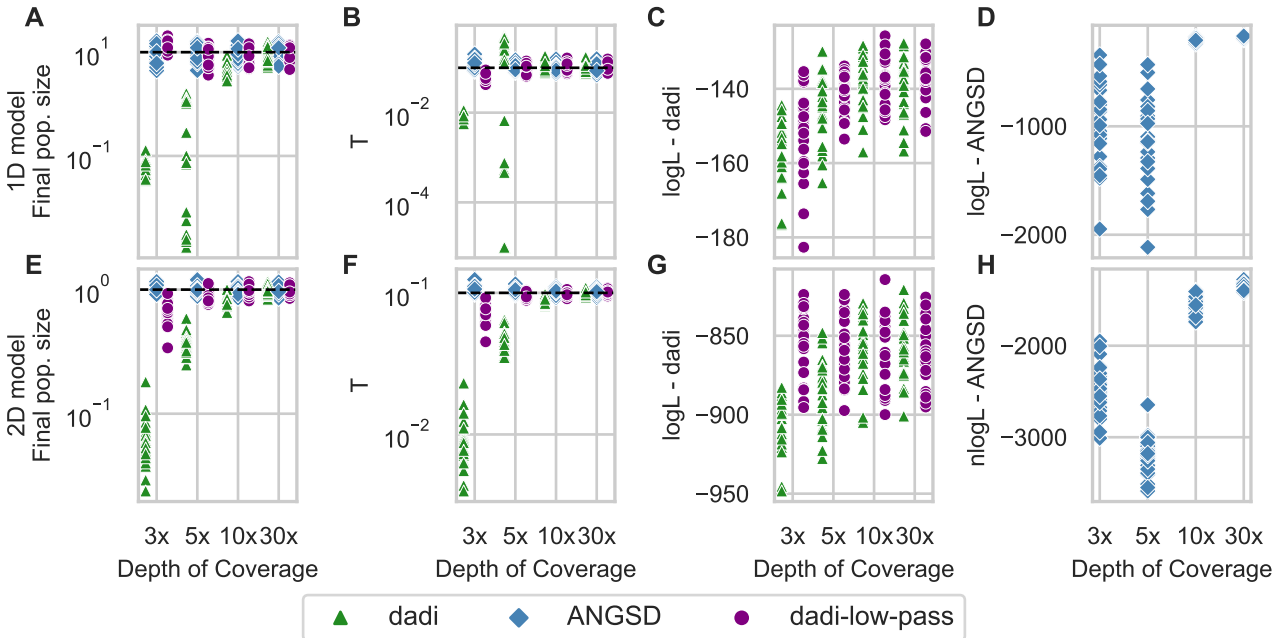


Figure 3: Our low-pass model and ANGSD enable accurate demographic parameter inference. A&B: From data simulated under a single-population growth model, the final population size and time of growth onset ( $T$ ) were accurately inferred using our low-pass model and a GATK-called AFS or using normal dadi and an ANGSD-called AFS. But they were biased if low depth was not accounted for when fitting a GATK-called AFS. (Dashed horizontal lines are simulated true values.) C: The likelihoods using the GATK AFS were similar whether or not low-pass biases were modeled. D: The fluctuations introduced into the AFS by ANGSD caused low likelihoods at low depth. E-H: For a two-population split with isolation model, similar results were found, although inferences from our low-pass model were slightly biased at  $3\times$  coverage.

171 population size was underestimated (Fig. 3A), consistent with a deficit of low-frequency alleles. The timing  
 172 of growth onset was also inaccurately inferred, underestimated at  $3\times$  depth and overestimated at  $5\times$  depth  
 173 (Fig. 3B). When the same data were fit with our low-pass model, both model parameters were accurately  
 174 recovered (Fig. 3A&B) even at the lowest depth. Fits to the AFS reconstructed by ANGSD also yielded  
 175 accurate model parameters (Fig. 3A&B).

176 The logarithm of the likelihood is commonly used to assess the quality of model fit. ANGSD reconstructs  
 177 the AFS for the full sequenced sample size, while we subsample in our approach to deal with missing geno-  
 178 types, so the likelihoods are not directly comparable. The likelihoods of models fit to the subsampled GATK  
 179 data were similar whether or not low-pass biases were modeled (Fig. 3C), suggesting that the likelihood itself  
 180 cannot be used to detect unmodeled low-pass bias. When fitting AFS estimated by ANGSD, likelihoods  
 181 were much lower at low coverage than high coverage (Fig. 3D), likely driven by the fluctuations ANGSD  
 182 introduced into the estimated AFS (Fig. 2).

183 For two-population data simulated under an isolation model (Fig. S2B), similar results were found. Fitting  
 184 the observed low-pass AFS with our model enabled accurate parameter inference (although there was some  
 185 bias at  $3\times$  coverage) as did fitting the AFS estimated by ANGSD (Fig. 3D). As in the single-population case,  
 186 likelihoods were substantially lower when fitting the ANGSD-estimated AFS, consistent with introduced  
 187 fluctuations in the AFS (Fig. S4).

188 For one-population data simulated under a growth model with inbreeding (Fig. S2C), failing to correct  
 189 for low-pass biases at low inbreeding ( $F = 0.1$  or  $F = 0.5$ ) led to similar biases as with no inbreeding, which  
 190 our low-pass model corrected (Fig. S7). For high inbreeding ( $F = 0.9$ ), the impact of low-pass sequencing



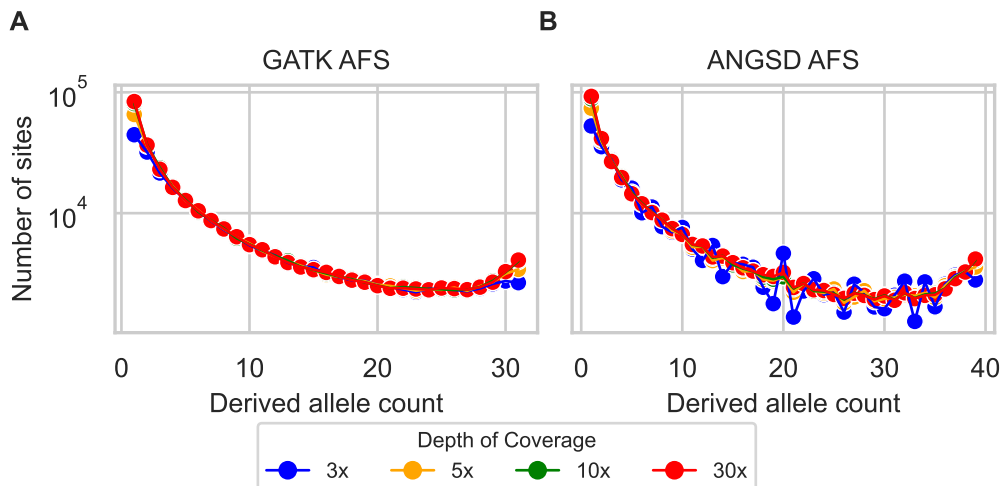


Figure 4: Allele frequency spectra from 20 YRI samples versus subsampled sequencing depth. A: Spectra generated using the GATK pipeline and subsampled to 32 haplotypes to accommodate missing genotypes. B: Spectra generated using ANGSD genotype likelihood optimization with BAM files input.

191 on accuracy was smaller, because inbreeding reduces heterozygosity (Fig. S7).

192 When applying our low-pass bias correction, the user must specify a value for inbreeding, while they may  
193 separately estimate it during demographic parameter optimization. We tested the impact of misspecifying  
194 inbreeding in the low-bias correction using data simulated with moderate inbreeding of  $F = 0.5$ . Large  
195 inbreeding values were inferred if inbreeding was initially underestimated in the low-coverage model, and  
196 small values were inferred if inbreeding was initially overestimated (Fig. S8C). A substantial difference  
197 between the inbreeding coefficient used for correction and the inferred value thus suggests that the assumed  
198 inbreeding coefficient was not optimal. Users can thus iterate and update the value assumed in the low-pass  
199 correction to converge to a best inference of inbreeding.

## 200 Analysis of human data

201 To empirically validate our approach and compare with ANGSD, we used chromosome 20 sequencing data  
202 from the 1000 Genomes Project, focusing on two sets of samples: Yoruba from Ibadan, Nigeria (YRI) and  
203 Utah residents of Northern and Western European ancestry (CEU). We inferred a single-population two-  
204 epoch demographic model (Fig. S9A) from the YRI samples, and a two-population isolation-with-migration  
205 model (Fig. S9B) from the combined YRI and CEU samples. To mimic low-pass sequencing, we subsampled  
206 the original high-depth data (which averaged  $30\times$  per site per individual) to create data with low to medium  
207 depth.

208 As with simulated data, the observed AFS from low-pass subsampled data was biased compared to  
209 high-pass data (Fig. 4A). Using the GATK pipeline, low-pass data yielded few low- and high-frequency  
210 derived alleles. In contrast to the simulated data, on these real data ANGSD failed to recover the correct  
211 number of low-frequency alleles at  $3\times$  and  $5\times$  depth, while still introducing large fluctuations at intermediate  
212 frequencies (Fig. 4B).

213 If low-pass biases were corrected for, we expected the inferred demographic parameters from subsampled  
214 low-pass data to match those from the original high-pass data. For the two-epoch model fit to YRI data, we  
215 found that with a GATK-called AFS and no low-pass model (Table 1), the inferred population sizes were  
216 biased downward and the times were inaccurate, similar to the growth model fit to simulated data. With  
217 the low-pass model, inferred values for low depth were similar to those for high depth, with some deviation  
218 at  $3\times$  (Table 1). Results from fitting ANGSD-estimated spectra were similar to not modeling low depth,

219 suggesting that ANGSD is ineffective for these data (Table 1). As with simulated data, the likelihoods for  
 220 ANGSD at low depth were also low.

221 For the isolation-with-migration model fit to YRI and CEU data, the results were broadly similar (Ta-  
 222 ble S1.) For population sizes and the divergence time, inferences were more stable from GATK genotyping  
 223 and our low-pass model than from ANGSD-estimated AFS. By contrast, the inferred migration rate was  
 224 similar across analyses.

Table 1: One-population YRI model analysis results. Inferred demographic parameters in dadi using empirical GATK and ANGSD AFS. We analyzed GATK empirical spectra without (dadi) and with low-pass correction (low-pass).

parameter	AFS	model	depth			
			30×	10×	5×	3 ×
$\nu_{YRI}$	GATK	dadi	1.82	1.76	1.54	0.05
	GATK	low-pass	1.83	1.77	1.70	1.54
	ANGSD	dadi	1.81	1.80	1.67	0.08
$T$	GATK	dadi	0.43	0.51	0.88	0.001
	GATK	low-pass	0.42	0.49	0.51	0.43
	ANGSD	dadi	0.55	0.68	0.96	0.001
$\theta (\times 10^4)$	GATK	dadi	5.15	5.08	4.81	5.99
	GATK	low-pass	5.16	5.10	5.10	5.20
	ANGSD	dadi	5.52	5.32	5.03	6.78
log-likelihood	GATK	dadi	-297	-253	-533	-1312
	GATK	low-pass	-302	-259	-283	-339
	ANGSD	dadi	-475	-486	-1120	-5905

## 225 Discussion

226 We assessed the biases introduced by low-pass sequencing using GATK multi-sample genotype calling and  
 227 developed a model to mitigate these biases. In a simulated population undergoing growth, we found that  
 228 low-pass sequencing reduced the presence of low-frequency alleles (Fig. 1). Our model accounted for these  
 229 biases, contrasting with ANGSD, which created fluctuations in the AFS at low depth (Fig. 2). In scenar-  
 230 ios involving two populations, we observed similar biases, which our model effectively corrected, whereas  
 231 ANGSD introduced additional noise (Fig. S3 and S4). For demographic inference, using our model enabled  
 232 accurate parameter estimates even at low-pass depths, while neglecting low-depth biases resulted in substan-  
 233 tial inaccuracies (Fig. 3). ANGSD also yielded accurate estimates, but worse likelihoods. Empirical testing  
 234 using human data from the 1000 Genomes Project showcased the accuracy of our correction method in  
 235 improving demographic inference from low-pass data, outperforming both uncorrected analysis and ANGSD  
 236 results (Fig. 4 and Tables 1 and S1).

237 While ANGSD is recognized for its effectiveness in managing low-pass sequencing, our results showed  
 238 its difficulties in modeling medium-frequency alleles. This is reflected in lower likelihood scores, particu-  
 239 larly when comparing low-pass datasets to high-pass ones (Fig. 3). Despite their utility in incorporating  
 240 uncertainty related to low-pass sequencing (Nielsen et al. 2011; Fumagalli 2013; Korneliussen et al. 2014),  
 241 genotype likelihoods might not always accurately capture the entire range of allele frequencies. Despite  
 242 the AFS fluctuations, ANGSD yielded reliable parameter estimates for simulated data. But ANGSD was  
 243 unable to accurately estimate the demographic parameters of real datasets, as demonstrated in the analysis  
 244 of the 1000 Genomes Project data (Tables 1 and S1). This underscores the need for rigorous and critical  
 245 assessments of results by evaluating the likelihood of the model and conducting uncertainty analysis.

246 Variant discovery using GATK involves two main approaches: multi-sample (classic joint-calling) and  
 247 single-sample calling (Nielsen et al. 2011). We modeled multi-sample calling, which has higher statistical



248 power compared to single-sample calling (Nielsen et al. 2011; Poplin et al. 2018). But multi-sample calling  
249 can become computationally burdensome with larger sample sizes, leading to the development of incremental  
250 single-calling as a scalable alternative (McKenna et al. 2010; Auwera & O’Connor 2020). When our model  
251 was applied to incremental single-calling AFS from subsampled 1000 Genomes Project data, parameter  
252 inference was poor (Table S2). Therefore, our model should only be used with multi-sample calling, and a  
253 slightly different model may need to be developed for incremental single-calling.

254 We present a GATK multi-sample calling model designed to compensate for AFS biases introduced by  
255 low-pass sequencing. Although tailored for GATK, our model’s design allows for its extension to different  
256 pipelines with modifications to address the unique aspects of each calling algorithm. For example, our model  
257 currently assumes that a site is called when at least two reads supporting the alternative allele are found  
258 (Eq. 3 and 4), but this could be modified for other pipelines with different calling criteria. Our approach  
259 can thus be generalized to other calling pipelines, including those using short reads, long reads, and hybrid  
260 approaches (e.g. Bankevich et al. 2012; Poplin et al. 2018). Note that our mathematical model assumes a  
261 shared read depth distribution among all individuals, and some studies may vary depth among individuals.  
262 Simulations suggest, however, that our model remain accurate with uneven depths (Fig. S10).

263 Our approach can also be integrated into other AFS-based inference tools such as moments (Portik et al.  
264 2017; Leaché et al. 2019), fastsimcoal2 (Excoffier et al. 2013, 2021), GADMA (Noskova et al. 2020), and  
265 delimitR (Smith & Carstens 2020), because our approach modifies the model AFS, independent of how it  
266 is computed. Our approach may also be useful in Approximate Bayesian Computation (Beaumont 2010;  
267 Csilléry et al. 2012) and machine learning workflows (Pudlo et al. 2016; Smith & Carstens 2020), facilitating  
268 simulation of low-pass datasets. Note, however, that we model bias in the mean shape of the AFS under  
269 low-pass sequencing, not its full variance (Fig. S11). Furthermore, AFS-based analyses are used not only  
270 for demographic studies but also to examine natural selection, including inferring the distribution of fitness  
271 effects of new mutations (Eyre-Walker & Keightley 2007; Huang et al. 2021). Our approach can thus facilitate  
272 population genomics research across tools, approaches, and problem domains.

273 In conclusion, we have developed a robust correction for low-pass sequencing biases, significantly enhanc-  
274 ing the accuracy of demographic parameter estimation at various coverage depths. As the genetic research  
275 community continues to address challenges associated with low-pass data (Bryc et al. 2013; Korneliussen  
276 et al. 2014; Blischak et al. 2018; Meisner & Albrechtsen 2018), especially when constrained by economics or  
277 sample availability, our methodology provides enables more reliable genetic analysis.

## 278 Material and Methods

### 279 Simulating AFS under low-pass sequencing

280 We used msprime (Kelleher et al. 2016; Baumdicker et al. 2022) to generate SNP datasets via coalescent  
281 simulations. We simulated two demographic models. The demographic models were visualized using demes-  
282 draw (Gower et al. 2022). The first model, single-population exponential growth (Fig. S2A), involved two  
283 parameters: the relative population size  $\nu_1 = 10$  and time of past growth  $T = 0.1$  (in units of two times the  
284 effect population size generations). The second model, two-population isolation (Fig. S2B), involved three  
285 parameters: equal relative sizes of populations 1 and 2,  $\nu_1 = \nu_2 = 1$ , and divergence time in the past  $T = 0.1$ .  
286 For each model, we conducted 25 independent simulations. For the exponential growth model, we  
287 sampled 20 diploid individuals, whereas for the isolation model, we sampled 10 individuals per population.  
288 Both demographic scenarios used an ancestral effective population size  $N_e$  of 10,000, a sequence length of  
289  $10^7$  bp, a mutation rate of  $\mu = 10^{-7}$  per site per generation, and recombination rate of  $r = 10^{-7}$  per site  
290 per generation.

291 For simulations incorporating inbreeding, we used SLiM 4 (Messer 2013; Haller & Messer 2023). Datasets  
292 were generated under a bottleneck and growth model (Fig. S2C), with a population bottleneck of  $\nu_B = 0.25$ ,  
293 followed by a population expansion to  $\nu_F = 1.0$ . The time of the past bottleneck was set at  $T = 0.2$ , and  
294 the level of inbreeding was varied with  $F \in \{0.1, 0.5, 0.9\}$ . Inbreeding was introduced using the selfing rate,  
295 set to  $s = \frac{2F}{1+F}$ . Twenty-five independent simulations were conducted, with 20 individuals sampled for each

296 replicate. Simulation parameters were  $N_e = 1000$ ,  $L = 2 \times 10^6$  bp,  $\mu = 5 \times 10^{-6}$ , and  $r = 2.5 \times 10^{-6}$ , with  
297 a burn-in of 10,000 generations.

298 To create low-pass datasets, we used synthetic diploid genomes. For each simulation replicate, we gener-  
299 ated a random reference genome spanning 10 Mb with a GC content of 40%, resembling the human genome.  
300 Mutations were incorporated by altering single nucleotides at the positions observed in the SNP matrix  
301 generated during each simulation, assuming that all sites were biallelic. Diploid individual genomes were  
302 generated by randomly selecting two chromosomes from the population pool.

303 Using the synthetic individual genomes as templates, we simulated 126 bp paired-end short reads for each  
304 individual with InSilicoSeq v2.0.1. (Gourlé et al. 2019). We calculated the number of reads per scenario as  
305  $LC/R$ , where  $L$  is the genome length,  $C$  the coverage depth, and  $R$  the read length. Reads for each diploid  
306 chromosome were simulated with equal probability. Depth of coverage per individual was sampled from a  
307 normal distribution with means of 3, 5, 10, and 30 and corresponding standard deviations of 0.3, 0.5, 1,  
308 and 3 to explore coverage variability, which increased with coverage levels. These standard deviations were  
309 selected based on preliminary simulations that suggested they offer a realistic variance for each coverage  
310 level.

311 For each individual we aligned simulated reads to the reference genome using BWA v0.7.17 (Li et al. 2009).  
312 We then processed the aligned reads with SAMTools v1.10 (Li 2013) to perform sorting, indexing, and pileup  
313 generation. To generate GATK spectra, we used the GATK multi-sample approach via HaplotypeCaller  
314 v4.2 (McKenna et al. 2010; Auwera & O'Connor 2020). To minimize false positives, the identified variants  
315 underwent filtering based on GATK's Best Practices guidelines, with thresholds tailored to expected error  
316 rates and variant quality. These thresholds included depth-normalized variant confidence ( $QD < 2.0$ ),  
317 mapping quality ( $MQ < 40$ ), strand bias estimate ( $FS > 60.0$ ), and strand bias ( $SOR > 10.0$ ). The filtered  
318 SNP VCF files were subsequently used in demographic inference analyses to estimate population parameters  
319 based on the AFS of these variants. To generate ANGSD spectra, we used the BAM files containing  
320 information about each individual with reads aligned to the reference genome. Subsequently, realAFS was  
321 used to estimate a maximum-likelihood AFS through the Expectation-Maximization algorithm. ANGSD  
322 v0.94 analysis was executed with the following settings: `doSaf = 1`, `minMapQ = 1`, `minQ = 20`, and `GL = 2`.

### 323 Empirical subsampling of Human data

324 We used high-quality whole-genome sequencing data (30×) from the 1000 Genomes Project (1kGP), sourced  
325 from The International Genome Sample Resource data portal (<https://www.internationalgenome.org/>  
326 Fairley et al. 2020). The data comprised CRAM files aligned to the GRCh38 human reference genome.  
327 We focused on two sets of samples for our analysis: 10 randomly selected individuals from the Yoruba from  
328 Ibadan, Nigeria (YRI) samples and 10 from the Utah residents with Northern and Western European ancestry  
329 (CEU) samples. The specific individuals included for the YRI were NA18486, NA18499, NA18510, NA18853,  
330 NA18858, NA18867, NA18878, NA18909, NA18917, NA18924, and for the CEU NA07037, NA11829, NA11892,  
331 NA11918, NA11932, NA11994, NA12004, NA12144, NA12249, NA12273. Additionally, for a single-population  
332 demographic model, 20 YRI individuals were analyzed, which includes the initial 10 plus an additional 10  
333 samples: NA19092, NA19116, NA19117, NA19121, NA19138, NA19159, NA19171, NA19184, NA19204, and  
334 NA19223.

335 Initially, we converted the CRAM files to BAM format and indexed them using Picard tools (<https://broadinstitute.github.io/picard/>). We then isolated reads from chromosome 20 at the original 30×  
336 coverage, which we subsequently subsampled to 10×, 5×, and 3× coverage using samtools v1.10 (Li 2013)  
337 to emulate varying sequencing depths. Next, using GATK version 4.2.5 HaplotypeCaller (McKenna et al.  
338 2010; Auwera & O'Connor 2020), we called SNPs and indels from these varying coverage depths for each  
339 population. We employed multi-sample SNP calling, merging BAM files with identical coverage prior to  
340 processing with HaplotypeCaller. This approach yielded a raw output VCF file.

341 We also carried out a single-sample calling procedure. For this, individual BAM files were used directly  
342 as inputs for the GATK HaplotypeCaller with the `-ERC GVCF` flag to enable GVCF mode. Following this,  
343 we used GATK GenomicsDBImport to compile the individual variant calls into a cohesive data structure.  
344 This setup allowed us to conduct joint genotyping using GATK GenotypeGVCFs, ultimately producing a  
345

346 multi-sample VCF.

347 Following SNP calling, we employed GATK SelectVariants to filter out indels for both approaches, re-  
348 taining only SNPs. Quality filtering of SNPs was conducted using GATK VariantFiltration, applying criteria  
349 such as depth-normalized variant confidence ( $QD < 2.0$ ), mapping quality ( $MQ < 40$ ), strand bias estimate  
350 ( $FS > 60.0$ ), and overall strand bias ( $SOR > 10.0$ ). After quality filtering, the VCF files were annotated  
351 with ancestral allele information using the `vcftools fill-aa` module, based on data from the Ensembl  
352 Release 110 Database (Danecek et al. 2011).

353 Finally, we used ANGSD to generate an AFS by using BAM files as input. The sample allele frequencies  
354 were first estimated using ANGSD's `-doSaf` flag, using GATK genotype likelihoods. These likelihoods  
355 were then used to calculate the AFS via the Expectation-Maximization algorithm using ANGSD's `realAFS`  
356 program. In this way, we maintained the original sample sizes from the BAM files, resulting in AFS for 40  
357 chromosomes in the single-population analysis and 20 chromosomes per population in the two-population  
358 analysis.

### 359 Demographic inference using dadi

360 We used dadi (Gutenkunst et al. 2009) to fit demographic models to simulated and empirical datasets. For  
361 the GATK spectra, we used the VCF files as input and subsampled individuals to accommodate missing data.  
362 For the ANGSD spectra, we used them as input directly. Within dadi, we used three demographic models for  
363 the simulated datasets: (i) an exponential growth model: `dadi.Demographics1D.growth`; (ii) a divergence  
364 model with migration fixed to zero: `dadi.Demographics2D.split_mig`; (iii) an bottleneck then exponential  
365 growth model modified to incorporate inbreeding: `dadi.Demographics1D.bottlegrowth`. For the human  
366 datasets, we used two models: (i) a divergence with migration model: `dadi.Demographics2D.split_mig`  
367 and (ii) an instantaneous growth model: `dadi.Demographics1D.two_epoch`. The extrapolation grid points  
368 were set using the formula  $[max(ns) + 120, max(ns) + 130, max(ns) + 140]$ , where  $ns$  is the sample size of  
369 the AFS. Our low-coverage correction is also implemented in dadi-cli (Huang et al. 2023).

### 370 Acknowledgements

371 This work was supported by the National Institute of General Medical Sciences of the National Institutes of  
372 Health (R01GM127348 and R35GM149235 to R.N.G.).

### 373 References

- 374 Auwera GAVd, O'Connor BD (2020) *Genomics in the cloud: using Docker, GATK, and WDL in Terra*.  
375 O'Reilly, Beijing Boston Farnham Sebastopol Tokyo, first edition edition.
- 376 Balding DJ, Nichols RA (1995) A method for quantifying differentiation between populations at multi-allelic  
377 loci and its implications for investigating identity and paternity. *Genetica* 96:3.
- 378 Balding DJ, Nichols RA (1997) Significant genetic correlations among Caucasians at forensic DNA loci.  
379 *Heredity* 78:583.
- 380 Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM, Nikolenko SI, Pham S,  
381 Prjibelski AD, Pyshkin AV, Sirotkin AV, Vyahhi N, Tesler G, Alekseyev MA, Pevzner PA (2012) SPAdes:  
382 A new genome assembly algorithm and its applications to single-cell sequencing. *Journal of Computational*  
383 *Biology* 19:455.
- 384 Baumdicker F, Bisschop G, Goldstein D, Gower G, Ragsdale AP, Tsambos G, Zhu S, Eldon B, Ellerman  
385 EC, Galloway JG, Gladstein AL, Gorjanc G, Guo B, Jeffery B, Kretzschumar WW, Lohse K, Matschiner  
386 M, Nelson D, Pope NS, Quinto-Cortés CD, Rodrigues MF, Saunack K, Sellinger T, Thornton K, Van Ke-  
387 menade H, Wohns AW, Wong Y, Gravel S, Kern AD, Koskela J, Ralph PL, Kelleher J (2022) Efficient  
388 ancestry and mutation simulation with msprime 1.0. *Genetics* 220:iyab229.

- 389 Beaumont MA (2010) Approximate Bayesian computation in evolution and ecology. *Annual Review of*  
390 *Ecology, Evolution, and Systematics* 41:379.
- 391 Blischak PD, Barker MS, Gutenkunst RN (2020) Inferring the demographic history of inbred species from  
392 genome-wide SNP frequency data. *Molecular Biology and Evolution* 37:2124.
- 393 Blischak PD, Kubatko LS, Wolfe AD (2018) SNP genotyping and parameter estimation in polyploids using  
394 low-coverage sequencing data. *Bioinformatics* 34:407.
- 395 Bryc K, Patterson N, Reich D (2013) A novel approach to estimating heterozygosity from low-coverage  
396 genome sequence. *Genetics* 195:553.
- 397 Carstens BC, Smith ML, Duckett DJ, Fonseca EM, Thomé MTC (2022) Assessing model adequacy leads to  
398 more robust phylogeographic inference. *Trends in Ecology & Evolution* 37:402.
- 399 Csilléry K, François O, Blum MGB (2012) abc: an R package for approximate Bayesian computation (ABC).  
400 *Methods in Ecology and Evolution* 3:475.
- 401 Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE, Lunter G, Marth GT,  
402 Sherry ST, McVean G, Durbin R, 1000 Genomes Project Analysis Group (2011) The variant call format  
403 and VCFtools. *Bioinformatics* 27:2156.
- 404 Duckett DJ, Calder K, Sullivan J, Tank DC, Carstens BC (2023) Reduced representation approaches produce  
405 similar results to whole genome sequencing for some common phylogeographic analyses. *PLOS ONE*  
406 18:e0291941.
- 407 Duitama J, Kennedy J, Dinakar S, Hernández Y, Wu Y, Măndoiu II (2011) Linkage disequilibrium based  
408 genotype calling from low-coverage shotgun sequencing reads. *BMC Bioinformatics* 12:S53.
- 409 Excoffier L, Dupanloup I, Huerta-Sánchez E, Sousa VC, Foll M (2013) Robust demographic inference from  
410 genomic and SNP data. *PLoS Genetics* 9:e1003905.
- 411 Excoffier L, Marchi N, Marques DA, Matthey-Doret R, Gouy A, Sousa VC (2021) fastsimcoal2: demographic  
412 inference under complex evolutionary scenarios. *Bioinformatics* 37:4882.
- 413 Eyre-Walker A, Keightley PD (2007) The distribution of fitness effects of new mutations. *Nature Reviews*  
414 *Genetics* 8:610.
- 415 Fairley S, Lowy-Gallego E, Perry E, Flicek P (2020) The International Genome Sample Resource (IGSR)  
416 collection of open human genomic variation resources. *Nucleic Acids Research* 48:D941.
- 417 Fox EA, Wright AE, Fumagalli M, Vieira FG (2019) ngsLD: evaluating linkage disequilibrium using genotype  
418 likelihoods. *Bioinformatics* 35:3855.
- 419 Fumagalli M (2013) Assessing the effect of sequencing depth and sample size in population genetics inferences.  
420 *PLoS ONE* 8:e79667.
- 421 Gorjanc G, Cleveland MA, Houston RD, Hickey JM (2015) Potential of genotyping-by-sequencing for genomic  
422 selection in livestock populations. *Genetics Selection Evolution* 47:12.
- 423 Gourel H, Karlsson-Lindsjö O, Hayer J, Bongcam-Rudloff E (2019) Simulating Illumina metagenomic data  
424 with InSilicoSeq. *Bioinformatics* 35:521.
- 425 Gower G, Ragsdale AP, Bisschop G, Gutenkunst RN, Hartfield M, Noskova E, Schiffels S, Struck TJ, Kelleher  
426 J, Thornton KR (2022) Demes: a standard format for demographic models. *Genetics* 222:iyac131.
- 427 Gutenkunst RN, Hernandez RD, Williamson SH, Bustamante CD (2009) Inferring the joint demographic  
428 history of multiple populations from multidimensional SNP frequency data. *PLoS Genetics* 5:e1000695.

- 429 Haller BC, Messer PW (2023) SLiM 4: multispecies eco-evolutionary modeling. *The American Naturalist*  
430 201:E127.
- 431 Huang X, Fortier AL, Coffman AJ, Struck TJ, Irby MN, James JE, León-Burguete JE, Ragsdale AP,  
432 Gutenkunst RN (2021) Inferring genome-wide correlations of mutation fitness effects between popula-  
433 tions. *Molecular Biology and Evolution* 38:4588.
- 434 Huang X, Struck TJ, Davey SW, Gutenkunst RN (2023) dadi-cli: automated and distributed population  
435 genetic model inference from allele frequency spectra.
- 436 Kelleher J, Etheridge AM, McVean G (2016) Efficient coalescent simulation and genealogical analysis for  
437 large sample sizes. *PLOS Computational Biology* 12:e1004842.
- 438 Kim BY, Huber CD, Lohmueller KE (2017) Inference of the distribution of selection coefficients for new  
439 nonsynonymous mutations using large samples. *Genetics* 206:345.
- 440 Korneliussen TS, Albrechtsen A, Nielsen R (2014) ANGSD: Analysis of Next Generation Sequencing Data.  
441 *BMC Bioinformatics* 15:356.
- 442 Leaché AD, Portik DM, Rivera D, Rödel M, Penner J, Gvoždík V, Greenbaum E, Jongsma GFM, Ofori-  
443 Boateng C, Burger M, Eniang EA, Bell RC, Fujita MK (2019) Exploring rain forest diversification using  
444 demographic model testing in the African foam-nest treefrog *Chiromantis rufescens*. *Journal of Biogeog-*  
445 *raphy* 46:2706.
- 446 Li H (2013) Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.  
447 ArXiv:1303.3997 [q-bio].
- 448 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R, 1000  
449 Genome Project Data Processing Subgroup (2009) The Sequence Alignment/Map format and SAMtools.  
450 *Bioinformatics (Oxford, England)* 25:2078.
- 451 Li Y, Sidore C, Kang HM, Boehnke M, Abecasis GR (2011) Low-coverage sequencing: implications for design  
452 of complex trait association studies. *Genome Research* 21:940.
- 453 Lou RN, Jacobs A, Wilder AP, Therkildsen NO (2021) A beginner's guide to low-coverage whole genome  
454 sequencing for population genomics. *Molecular Ecology* 30:5966.
- 455 Maddison DR, Ruvolo M, Swofford DL (1992) Geographic origins of Human mitochondrial DNA: phyloge-  
456 netic evidence from control region sequences. *Systematic Biology* 41:111.
- 457 Marchi N, Winkelbach L, Schulz I, Brami M, Hofmanová Z, Blöcher J, Reyna-Blanco CS, Diekmann Y,  
458 Thiéry A, Kapopoulou A, Link V, Piuz V, Kreutzer S, Figarska SM, Ganiatsou E, Pukaj A, Struck  
459 TJ, Gutenkunst RN, Karul N, Gerritsen F, Pechtl J, Peters J, Zeeb-Lanz A, Lenneis E, Teschler-Nicola  
460 M, Triantaphyllou S, Stefanović S, Papageorgopoulou C, Wegmann D, Burger J, Excoffier L (2022) The  
461 genomic origins of the world's first farmers. *Cell* 185:1842.
- 462 Martin AR, Atkinson EG, Chapman SB, Stevenson A, Stroud RE, Abebe T, Akena D, Alemayehu M, Ashaba  
463 FK, Atwoli L, Bowers T, Chibnik LB, Daly MJ, DeSmet T, Dodge S, Fekadu A, Ferriera S, Gelaye B,  
464 Gichuru S, Injera WE, James R, Kariuki SM, Kigen G, Koenen KC, Kwobah E, Kyebuzibwa J, Majara  
465 L, Musinguzi H, Mwema RM, Neale BM, Newman CP, Newton CR, Pickrell JK, Ramesar R, Shiferaw W,  
466 Stein DJ, Teferra S, Van Der Merwe C, Zingela Z (2021) Low-coverage sequencing cost-effectively detects  
467 known and novel variation in underrepresented populations. *The American Journal of Human Genetics*  
468 108:656.
- 469 McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytzky A, Garimella K, Altshuler D, Gabriel  
470 S, Daly M, DePristo MA (2010) The Genome Analysis Toolkit: A MapReduce framework for analyzing  
471 next-generation DNA sequencing data. *Genome Research* 20:1297.



- 472 Meisner J, Albrechtsen A (2018) Inferring population structure and admixture proportions in low-depth  
473 NGS data. *Genetics* 210:719.
- 474 Messer PW (2013) SLiM: simulating evolution with selection and linkage. *Genetics* 194:1037.
- 475 Mota BS, Rubinacci S, Cruz Dávalos DI, G Amorim CE, Sikora M, Johannsen NN, Szmyt MH, Włodarczak  
476 P, Szczepanek A, Przybyła MM, Schroeder H, Allentoft ME, Willerslev E, Malaspina AS, Delaneau O  
477 (2023) Imputation of ancient human genomes. *Nature Communications* 14:3660.
- 478 Nielsen R, Korneliussen T, Albrechtsen A, Li Y, Wang J (2012) SNP Calling, genotype calling, and sample  
479 allele frequency estimation from new-generation sequencing data. *PLoS ONE* 7:e37558.
- 480 Nielsen R, Paul JS, Albrechtsen A, Song YS (2011) Genotype and SNP calling from next-generation se-  
481 quencing data. *Nature Reviews Genetics* 12:443.
- 482 Noskova E, Ulyantsev V, Koepfli KP, O'Brien SJ, Dobrynin P (2020) GADMA: Genetic algorithm for  
483 inferring demographic history of multiple populations from allele frequency spectrum data. *GigaScience*  
484 9:giaa005.
- 485 Poplin R, Chang PC, Alexander D, Schwartz S, Colthurst T, Ku A, Newburger D, Dijamco J, Nguyen N,  
486 Afshar PT, Gross SS, Dorfman L, McLean CY, DePristo MA (2018) A universal SNP and small-indel  
487 variant caller using deep neural networks. *Nature Biotechnology* 36:983.
- 488 Portik DM, Leaché AD, Rivera D, Barej MF, Burger M, Hirschfeld M, Rödel M, Blackburn DC, Fujita  
489 MK (2017) Evaluating mechanisms of diversification in a Guineo-Congolian tropical forest frog using  
490 demographic model selection. *Molecular Ecology* 26:5245.
- 491 Pudlo P, Marin JM, Estoup A, Cornuet JM, Gautier M, Robert CP (2016) Reliable ABC model choice via  
492 random forests. *Bioinformatics* 32:859.
- 493 Reid NM, Proestou DA, Clark BW, Warren WC, Colbourne JK, Shaw JR, Karchner SI, Hahn ME, Nacci D,  
494 Oleksiak MF, Crawford DL, Whitehead A (2016) The genomic landscape of rapid repeated evolutionary  
495 adaptation to toxic pollution in wild fish. *Science* 354:1305.
- 496 Sawyer SA, Hartl DL (1992) Population genetics of polymorphism and divergence. *Genetics* 132:1161.
- 497 Smith ML, Carstens BC (2020) Process-based species delimitation leads to identification of more biologically  
498 relevant species. *Evolution* 74:216.
- 499 Vieira FG, Fumagalli M, Albrechtsen A, Nielsen R (2013) Estimating inbreeding coefficients from NGS data:  
500 Impact on genotype calling and allele frequency estimation. *Genome Research* 23:1852.
- 501 Wakeley J (2009) *Coalescent theory: an introduction*. Roberts & Co. Publishers, Greenwood Village, Colo.



## Supplementary Material

### Modeling biases from low-pass genome sequencing to enable accurate population genetic inferences

Emanuel M. Fonseca<sup>\*</sup>, Linh N. Tran, Hannah Mendoza, Ryan N. Gutenkunst<sup>\*</sup>

Department of Molecular and Cellular Biology, University of Arizona, Tucson, AZ 85721, USA

**\*Corresponding author:** E-mail: [emanuelfonseca@arizona.edu](mailto:emanuelfonseca@arizona.edu); [rgutenk@arizona.edu](mailto:rgutenk@arizona.edu)

The correction for low-pass sequencing is performed using the publicly available dadi Python package, which can be accessed at <https://bitbucket.org/gutenkunstlab/dadi>. Additionally, the codebase for creating and analyzing both simulated and empirical datasets, ensuring reproducibility, is readily accessible on GitHub at <https://github.com/emanuelfonseca/low-coverage-sfs> and [https://github.com/lntan26/low-coverage-sfs/tree/main/empirical\\_analysis](https://github.com/lntan26/low-coverage-sfs/tree/main/empirical_analysis). Furthermore, we provide illustrative examples to assist users in implementing our methodology.

Table S1: Two-population model analysis results. Inferred demographic parameters in dadi using empirical GATK and ANGSD AFS. We analyzed GATK empirical spectra without (dadi) and with low-pass correction (low-pass).

parameter	AFS	model	depth			
			30×	10×	5×	3 ×
$\nu_{YRI}$	GATK	dadi	1.79	1.63	1.18	0.61
	GATK	low-pass	1.82	1.62	1.67	1.69
	ANGSD	dadi	1.69	1.58	1.26	0.87
$\nu_{CEU}$	GATK	dadi	0.38	0.37	0.31	0.17
	GATK	low-pass	0.38	0.38	0.36	0.34
	ANGSD	dadi	0.39	0.38	0.33	0.22
$T$	GATK	dadi	0.21	0.22	0.18	0.06
	GATK	low-pass	0.21	0.23	0.20	0.16
	ANGSD	dadi	0.21	0.22	0.20	0.07
$m$	GATK	dadi	1.80	2.00	2.24	1.68
	GATK	low-pass	1.80	2.00	1.89	1.66
	ANGSD	dadi	1.99	2.12	2.44	1.91
$\theta (\times 10^4)$	GATK	dadi	5.42	5.44	5.56	5.65
	GATK	low-pass	5.43	5.42	5.45	5.40
	ANGSD	dadi	6.04	6.01	6.03	6.25
log-likelihood	GATK	dadi	-2588	-2378	-2329	-2663
	GATK	low-pass	-2590	-2479	-2224	-1850
	ANGSD	dadi	-5518	-5595	-7074	-11029

Table S2: One-population model analysis results with single-sample calling using empirical GATK AFS. We analyzed GATK empirical single-sample call spectra without (dadi) and with low-pass correction (low-pass).

parameter	model	depth			
		30×	10×	5×	3 ×
$\nu_{YRI}$	dadi	1.85	1.87	1.82	1.56
	low-pass	1.86	1.93	2.73	3.60
$T$	dadi	0.43	0.45	0.51	0.48
	low-pass	0.42	0.40	0.24	0.24
$\theta (\times 10^3)$	dadi	5.13	5.05	4.62	4.31
	low-pass	5.14	5.10	4.96	4.49
log-likelihood	dadi	-284	-280	-457	-1755
	low-pass	-291	-317	-597	-1005

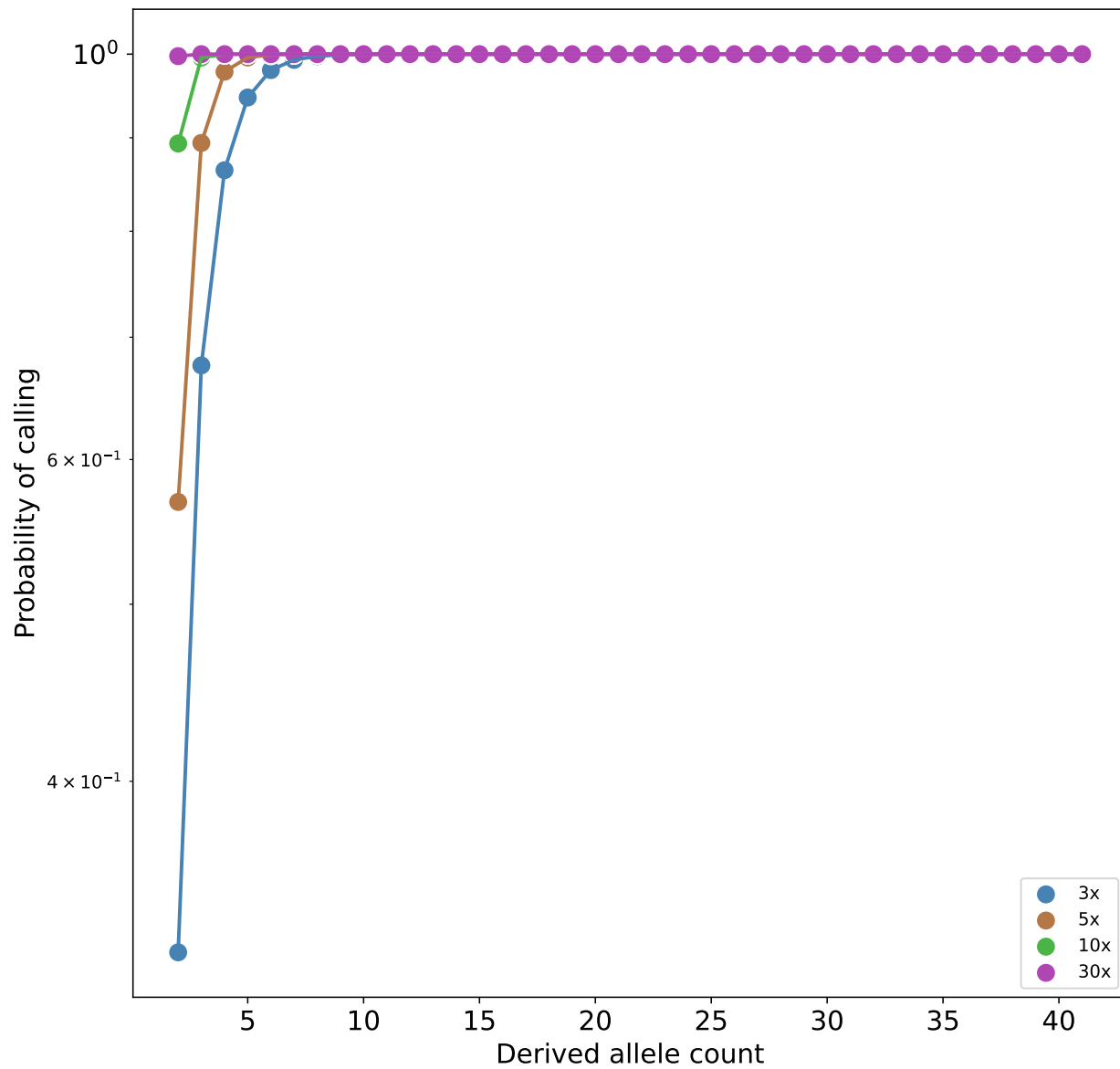


Figure S1: Probability of calling a variant site versus true allele frequency and coverage depth.

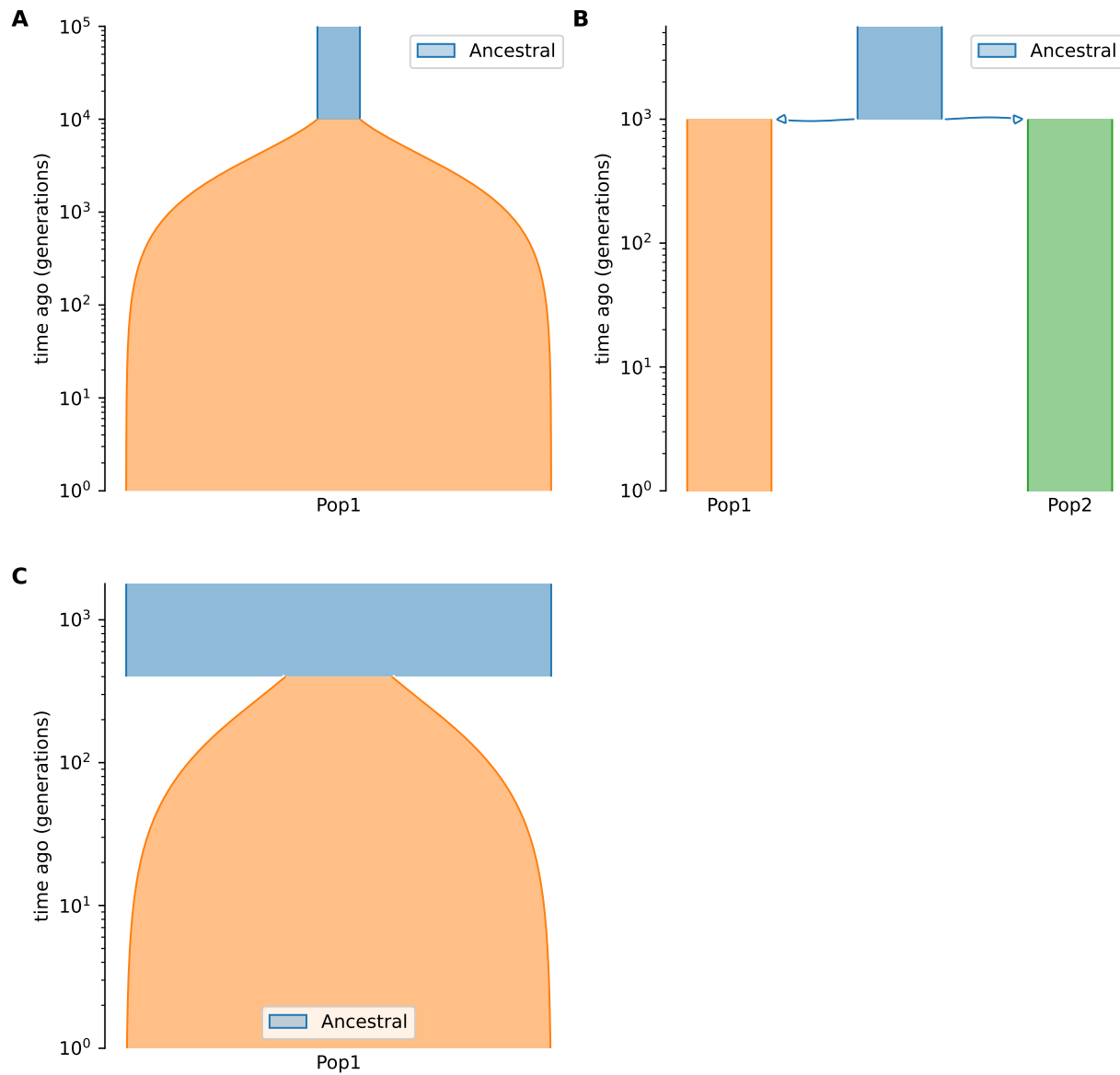


Figure S2: Representation of the demographic models used in the simulations: (A) single-population exponential growth model with parameters  $\nu_1 = 10$  and  $T = 0.1$ , (B) two-population isolation model with  $\nu_1 = \nu_2 = 1$  and  $T = 0.1$ , (C) single-population exponential growth model with inbreeding with parameters  $\nu_1 = 4$ ,  $T = 0.4$ , and  $F \in \{0.1, 0.5, 0.9\}$ .  $\nu$ ,  $T$ ,  $F$  represent relative population size, time in the past, and inbreeding coefficient, respectively. This plot was created with Demes (Gower et al. 2022)

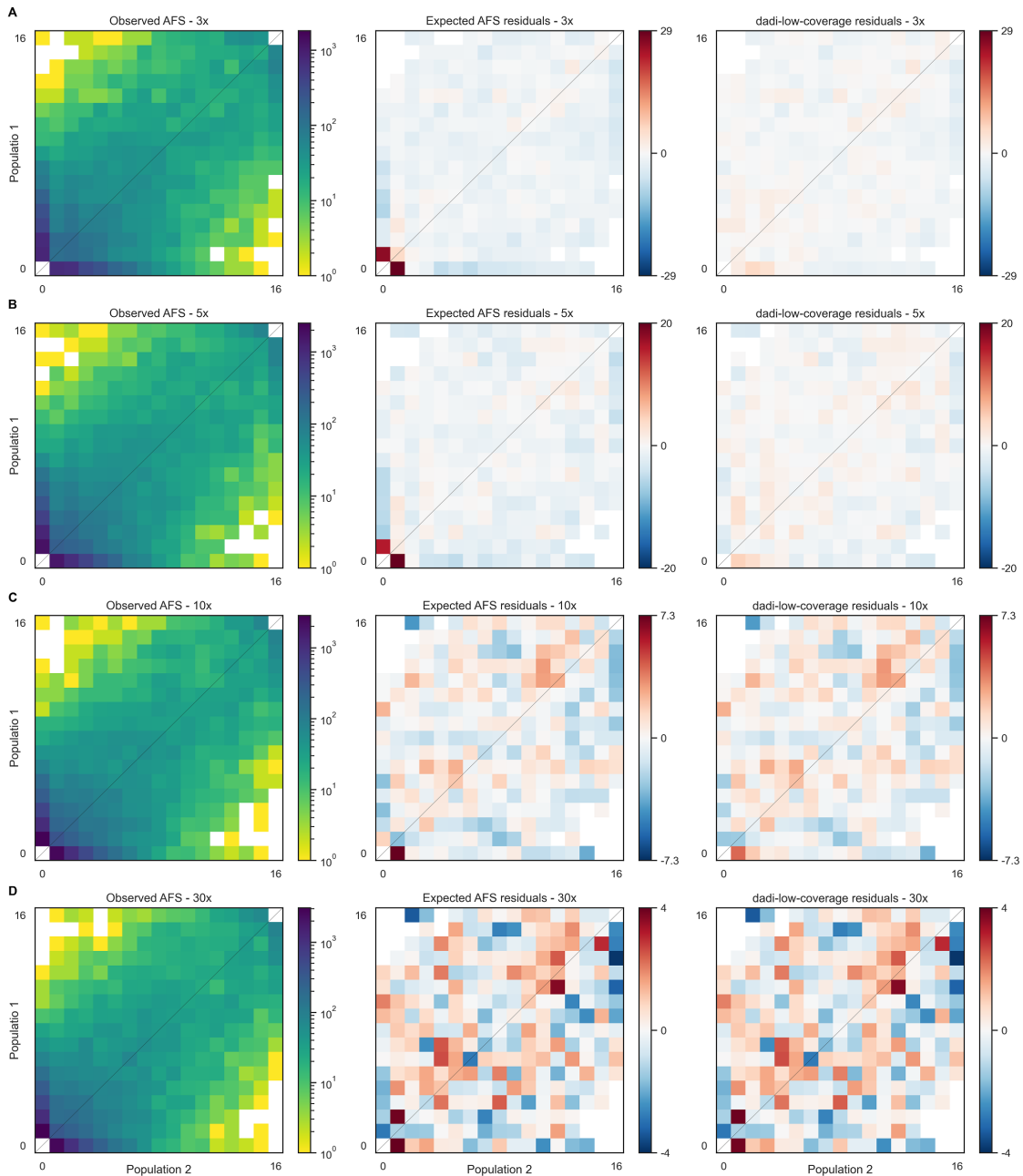


Figure S3: The observed 2D AFS is biased by low coverage. Deviation between the observed low-coverage AFS (first column) and the expected AFS (calculated by dadi) for the isolation demographic scenario is visualized through the residual plot (second column). Dark red residuals indicate that the observed low-coverage AFS is deficient in low-frequency alleles compared to the expectation. By contrast, the residuals between the observed AFS and the low-coverage model are much smaller. At 30 $\times$  coverage (D) the residuals become small and random, indicating agreement between all three spectra. Coverage depths compared are (A) 3 $\times$ , (B) 5 $\times$ , (C) 10 $\times$ , and (D) 30 $\times$ .

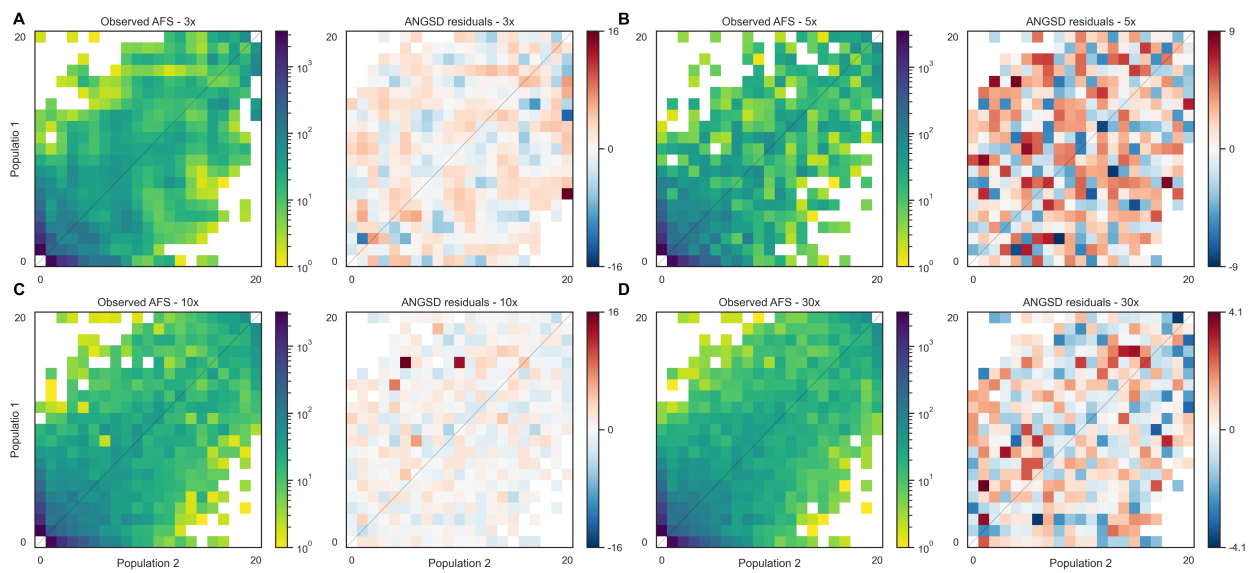


Figure S4: ANGSD creates fluctuations in the joint AFS. The joint AFS output by ANGSD exhibits sporadic very large residuals when compared with the true simulated AFS, similar to the oscillations seen in the single population AFS (Fig. 2). Coverage depths compared are (A) 3 $\times$ , (B) 5 $\times$ , (C) 10 $\times$ , and (D) 30 $\times$ .



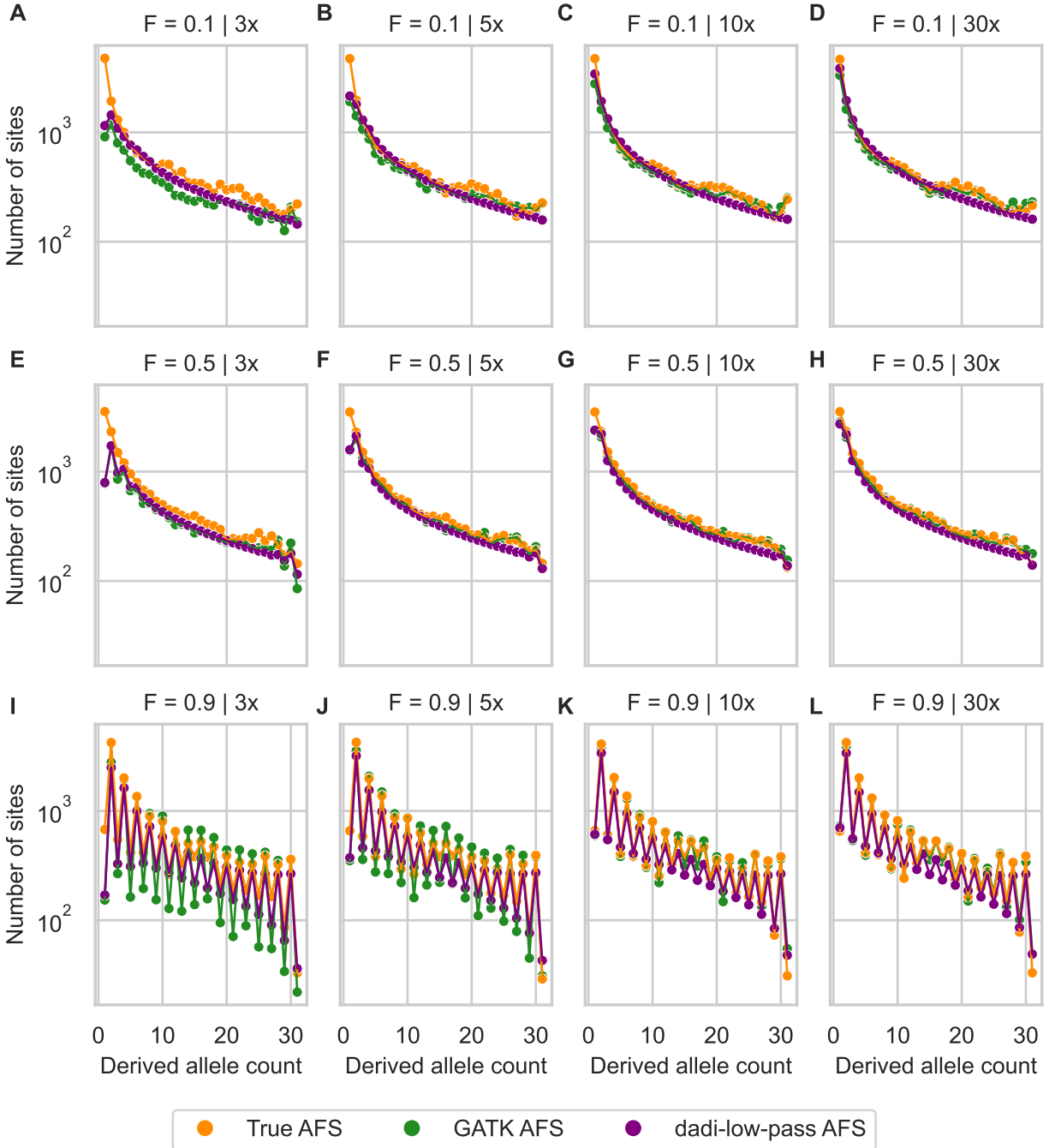


Figure S5: The observed AFS is impacted by low-pass sequencing (3 $\times$ , 5 $\times$ , 10 $\times$ , and 30 $\times$ ) and inbreeding ( $F \in \{0.1, 0.5, 0.9\}$ ). This figure presents a comparison of the observed AFS from low-pass variant calling with simulations in both the standard dadi and dadi-low-pass frameworks, using the true parameter values for a single-population model.

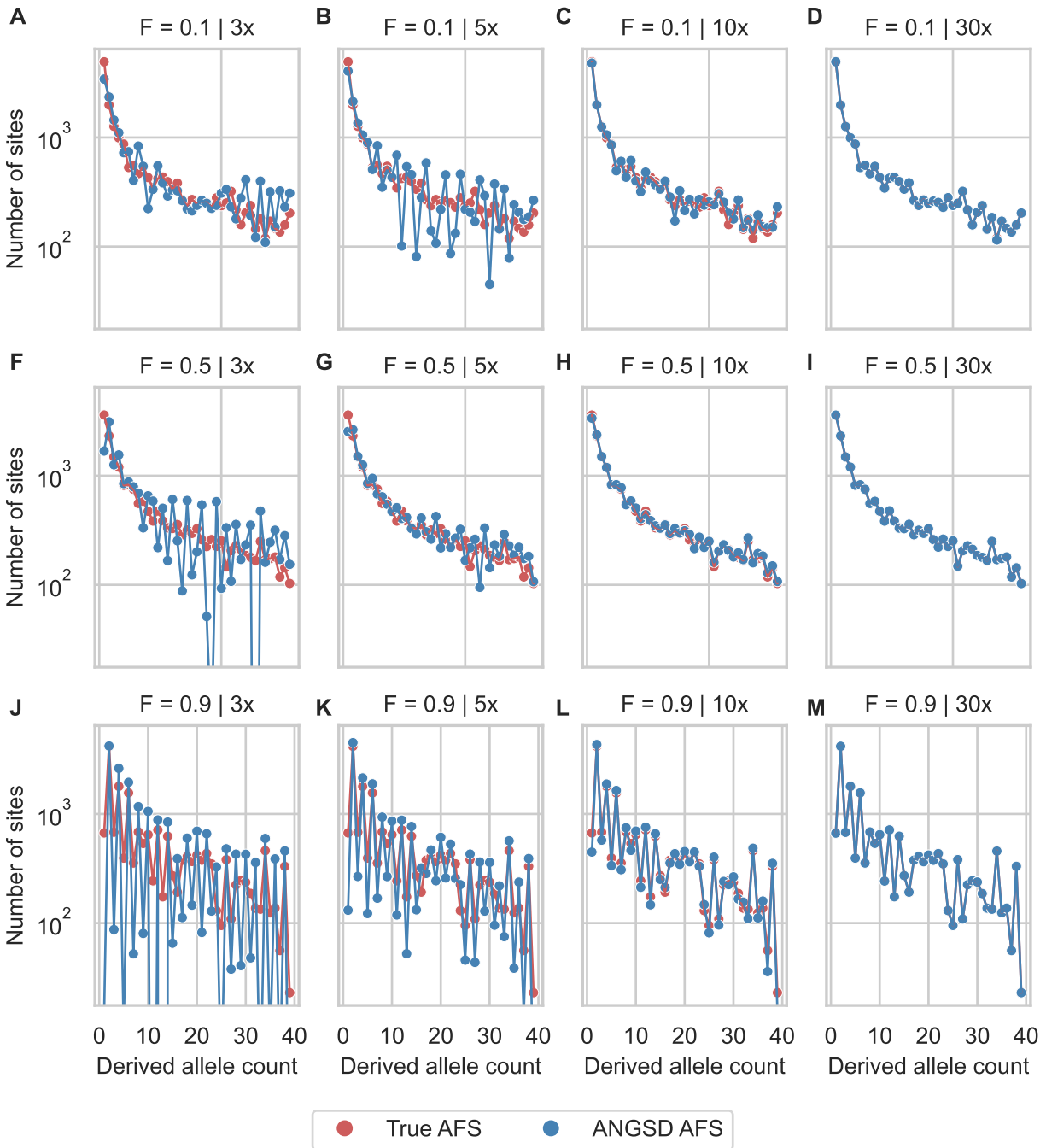


Figure S6: ANGSD corrects for the low-pass bias of the AFS, but it introduces fluctuations in inbreeding models. For the same simulations as Fig. S5, ANGSD (blue) was used to reconstruct the simulated AFS (red). Coverages were 3 $\times$ , 5 $\times$ , 10 $\times$ , and 30 $\times$ ) and inbreeding 0.1, 0.5, and 0.9.

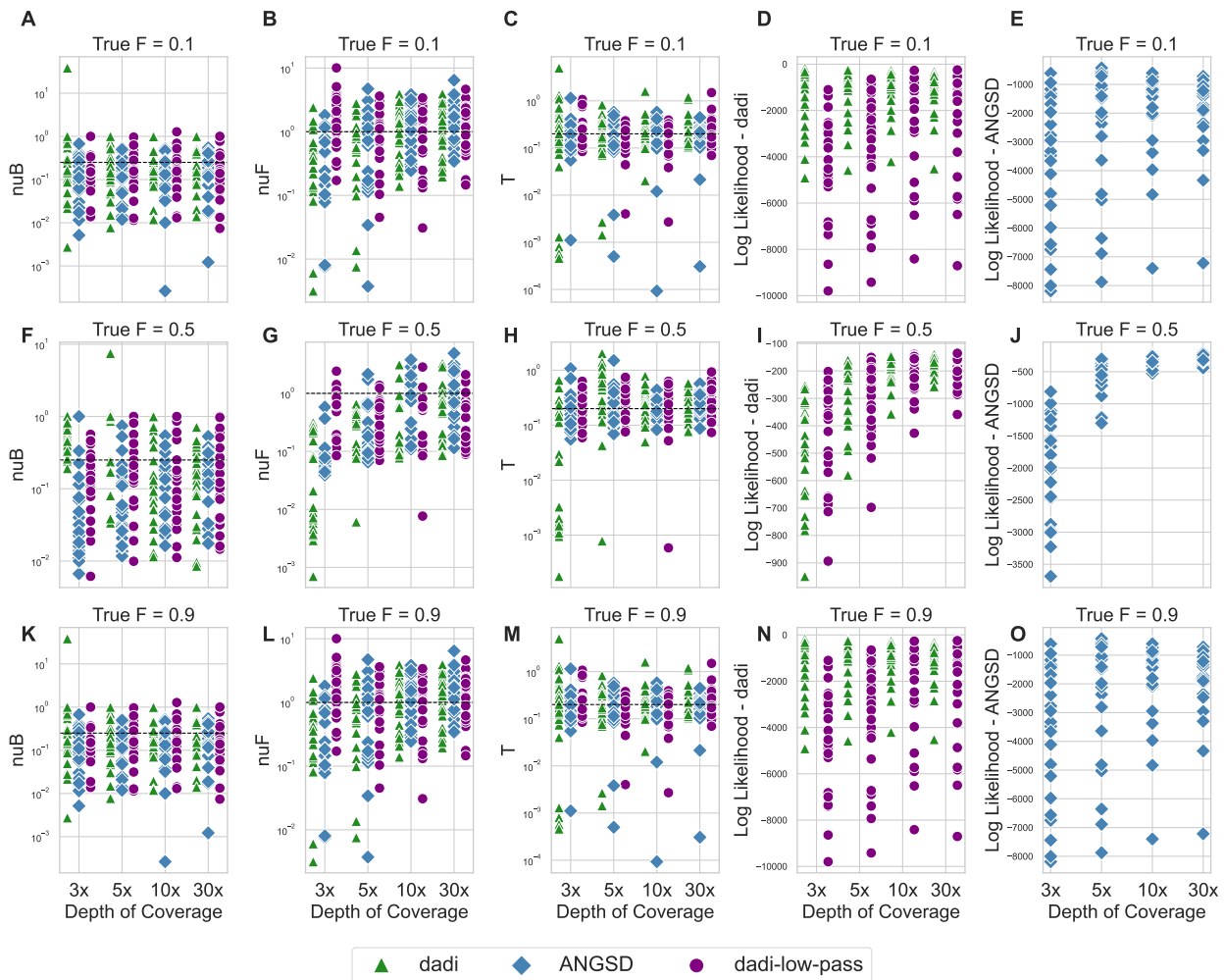


Figure S7: Graph showcasing the accuracy of parameter and likelihood estimations across various sequencing depths ( $3 \times$ ,  $5 \times$ ,  $10 \times$ , and  $30 \times$ ) and inbreeding ( $F \in \{0.1, 0.5, 0.9\}$ ) for a population bottleneck and growth model. The inbreeding parameters were kept fixed for both the low-pass calculation and the optimization process. Parameters were obtained through different methods, including dadi, both with and without corrections for low coverage, as well as ANGSD. Details of the graph include: (A), (F), (K) the estimated size after population bottleneck; (B), (G), (L) the estimated size after population expansion; (C), (H), (M) the time of population expansion; (D), (I), (N) log-likelihood calculations from dadi, highlighting the distinction between corrected and uncorrected model for low coverage; and (E), (J), (O) log-likelihood calculations from ANGSD. The black line present in the plots for (A), (B), (E), (F), (I), (J) and indicates the true value of the parameter, providing a standard for evaluating the accuracy of different approaches.

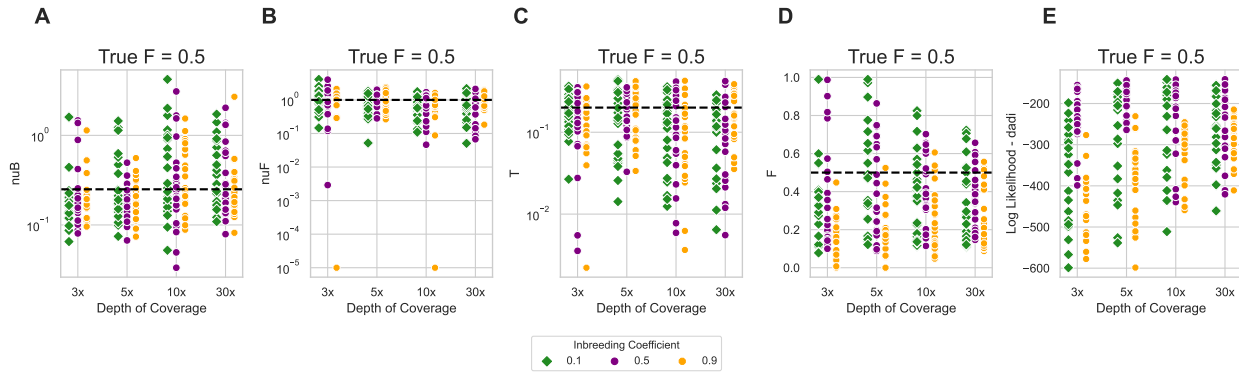


Figure S8: Graph showcasing the accuracy of parameter and likelihood estimations across various sequencing depths ( $3 \times$ ,  $5 \times$ ,  $10 \times$ , and  $30 \times$ ) and inbreeding ( $F \in \{0.1, 0.5, 0.9\}$ ) for a population expansion model under a true inbreeding value of 0.5. The inbreeding parameters used for the low-pass calculation were 0.1, 0.5, and 0.9. Parameters were obtained using dadi-low-pass. Details of the graph include: (A) the estimated size after population bottleneck; (B) the estimated size after population expansion; (C) the time of population expansion; (D) inferred inbreeding coefficient; (E) log-likelihood calculations from dadi-low-pass. The black line present in the plots for (A), (B), (C), and (D) indicates the true value of the parameter, providing a standard for evaluating the accuracy of different approaches.

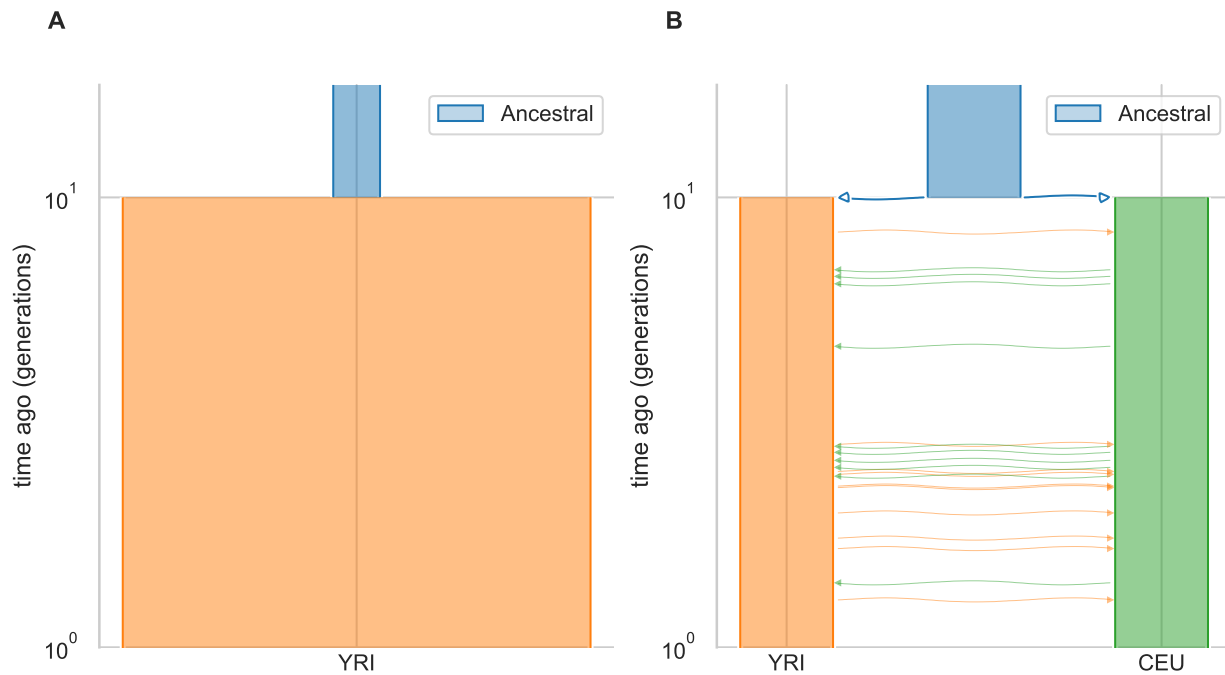


Figure S9: Representation of the demographic models used to analyse 1000 genomes datasets: (A) single-population two-epoch growth model with parameters, (B) two-population isolation with migration model. This plot was created with Demes (Gower et al. 2022)

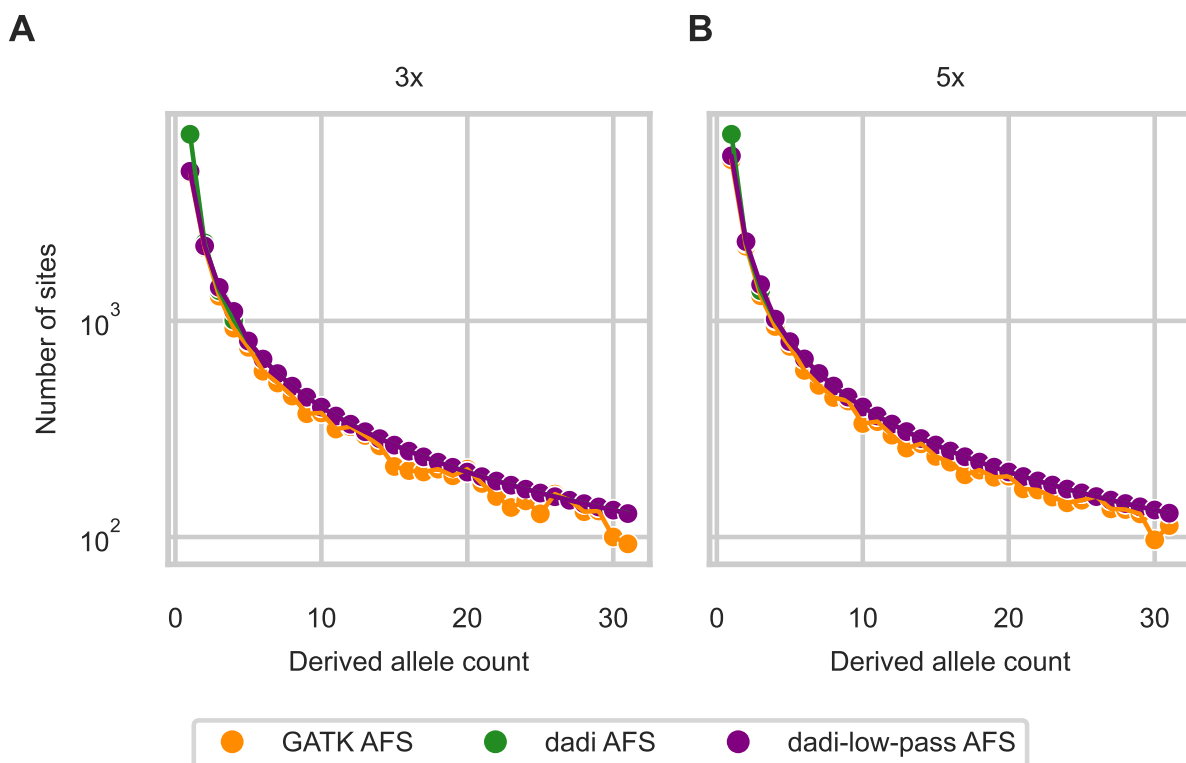


Figure S10: Unbalanced depth of coverage does not bias the dadi-low-pass model. Simulations were performed using 20 individuals, with half simulated under low-coverage conditions (A: 3× or B: 5×) and the other half under high-depth coverage (30×).

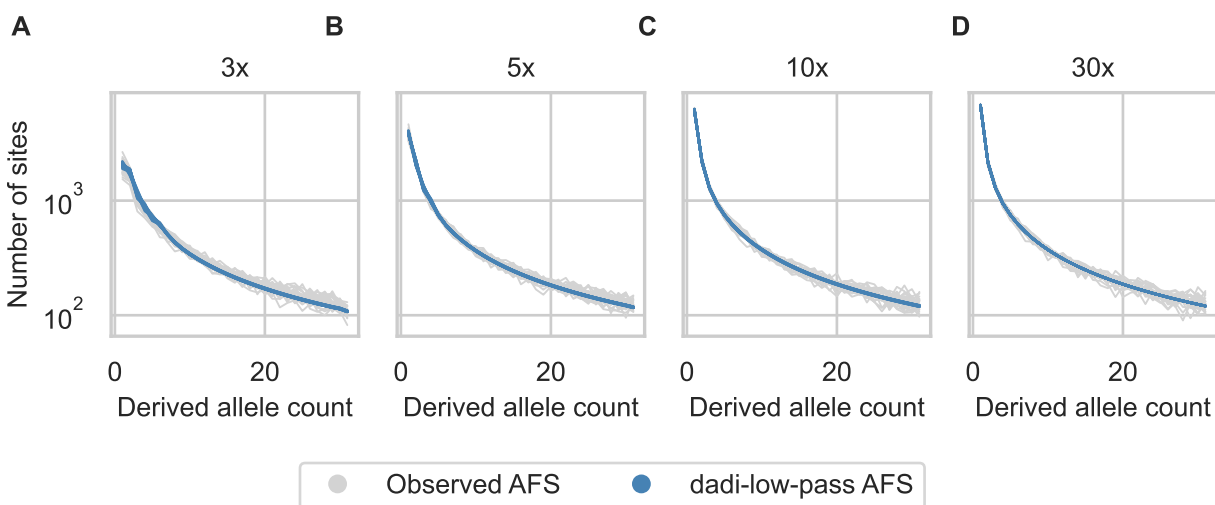


Figure S11: The simulated AFS under the low-pass model shows less variance compared to that observed in the simulated datasets. We generated 25 AFS for each condition.