

An empirical evaluation of genotype imputation of ancient DNA

Kristiina Ausmees ^{1,*} Federico Sanchez-Quinto,^{2,3} Mattias Jakobsson ³, Carl Nettelblad ¹

¹Department of Information Technology, Uppsala University, Uppsala 751 05, Sweden,

²Instituto Nacional de Medicina Genómica (INMEGEN), Mexico City 14610, Mexico,

³Human Evolution, Department of Organismal Biology, Uppsala University, Uppsala 752 36, Sweden

*Corresponding author: Department of Information Technology, Uppsala University, Box 337, Uppsala, 751 05, Sweden. Email: kristiina.ausmees@it.uu.se

Abstract

With capabilities of sequencing ancient DNA to high coverage often limited by sample quality or cost, imputation of missing genotypes presents a possibility to increase the power of inference as well as cost-effectiveness for the analysis of ancient data. However, the high degree of uncertainty often associated with ancient DNA poses several methodological challenges, and performance of imputation methods in this context has not been fully explored. To gain further insights, we performed a systematic evaluation of imputation of ancient data using Beagle v4.0 and reference data from phase 3 of the 1000 Genomes project, investigating the effects of coverage, phased reference, and study sample size. Making use of five ancient individuals with high-coverage data available, we evaluated imputed data for accuracy, reference bias, and genetic affinities as captured by principal component analysis. We obtained genotype concordance levels of over 99% for data with 1× coverage, and similar levels of accuracy and reference bias at levels as low as 0.75×. Our findings suggest that using imputed data can be a realistic option for various population genetic analyses even for data in coverage ranges below 1×. We also show that a large and varied phased reference panel as well as the inclusion of low- to moderate-coverage ancient individuals in the study sample can increase imputation performance, particularly for rare alleles. In-depth analysis of imputed data with respect to genetic variants and allele frequencies gave further insight into the nature of errors arising during imputation, and can provide practical guidelines for postprocessing and validation prior to downstream analysis.

Keywords: imputation; phasing; ancient DNA; low coverage; reference bias

Introduction

The possibility to sequence ancient DNA (aDNA) has increased capabilities to study archaeological remains and provided new insights into various aspects of human evolutionary history. Notable findings such as the detection of genetic introgression between anatomically modern humans and other hominins, confirmation of the African origin of modern humans, and an increased understanding of the spread of agriculture into Europe have been achieved through population genetic analyses of ancient and contemporary genomes (Nielsen *et al.* 2017).

Due to the age and varying preservation conditions that ancient samples may have been exposed to, aDNA has unique properties that pose methodological and computational challenges not present when working with data from present-day humans. Contamination of DNA from microbes and other nontarget sources can result in low proportions of endogenous DNA (Pääbo *et al.* 2004; Prüfer *et al.* 2010), leading to limitations in sample availability that can cause sequencing to high coverage depth to be impossible or prohibitively expensive. Contamination also leads to issues regarding data authenticity. In addition, the DNA molecule is subject to degradation over time, which can cause errors in the sequencing pipeline (Brotherton *et al.* 2007; Sánchez-Quinto *et al.* 2012).

Although the identification of patterns of damage unique to aDNA has allowed for methods of data authentication and improved the process of reconstructing ancient genomes (Stiller *et al.* 2006; Briggs *et al.* 2007; Krause *et al.* 2010; Sawyer *et al.* 2012), these characteristics nonetheless cause biases in sequencing and mapping that can impact variant calling and other forms of downstream analysis (Prüfer *et al.* 2010; Ginolhac *et al.* 2011; Parks and Lambert 2015). Consequently, studies of aDNA samples are often limited to low- to moderate-coverage data with higher degrees of uncertainty than modern samples typically exhibit.

Genotype imputation is a powerful tool that can increase the information content in a sample by inferring unobserved genotypes. Imputation has been widely applied in various scenarios analyzing modern data, e.g. to increase power of inference in genome-wide association studies and to conform samples from different studies for merged analysis (Zeggini *et al.* 2008; Spencer *et al.* 2009; Marchini and Howie 2010). For aDNA, the possibility to increase information content of sparse and noisy data can potentially improve the quality of results as well as expand the range of analyses that are possible to perform.

Many common imputation methods for unrelated samples rely on sequential probabilistic models in which missing genotypes are inferred based on similarity to other individuals. The

Received: February 28, 2022. Accepted: April 05, 2022

© The Author(s) 2022. Published by Oxford University Press on behalf of Genetics Society of America.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

estimation is founded on an assumption of the presence of short stretches of shared haplotypes that have been passed down from a distant common ancestor. Given the individual's genotypes, its haplotype phase is inferred, allowing missing variants to be predicted based on similarity to other samples. Many methods are able to leverage the information in the study sample as well as an additional panel of phased reference haplotypes when performing the phase estimation.

Many widely employed tools, such as MACH (Li *et al.* 2010), IMPUTE2 (Howie *et al.* 2009), and PHASE (Stephens *et al.* 2001; Stephens and Scheet 2005), are based on variants of the “the product of approximate conditionals” (PAC) framework (Li and Stephens 2003). This model represents the sample sequence as an imperfect mosaic of the reference haplotypes, generally considering all possible transitions over the state space, with explicit modeling of mutation and recombination. A haploid version of this framework is used in the software GLIMPSE (Rubinacci *et al.* 2021). A slightly different method is implemented in the software Beagle (Browning and Browning 2007), which is based on a model of local haplotype clusters based on similarity of the reference haplotypes at nearby markers. This results in a smaller state space with not all possible transitions considered at every position, reducing the computational burden. The effects of mutation and recombination are not explicitly modeled, but the change of cluster membership along the sequence can be seen as implicitly representing these biological processes.

The accuracy of genotype imputation depends on several factors, mainly related to the quality of the sample data and the properties of the phased reference panel. A larger sample size, as well as increased marker density and genotype accuracy generally results in better performance (Browning and Browning 2011). For data with high levels of uncertainty, imputation based on a probabilistic framework using genotype likelihoods rather than called genotypes may be beneficial (Browning and Yu 2009; Browning and Browning 2011; Nielsen *et al.* 2011), an option that is supported by some software tools, including Beagle v4.0. In Hui *et al.* (2020), a 2-step approach is introduced, in which genotype likelihood data is first used to obtain genotype probabilities based on a reference panel, and missing genotypes are subsequently imputed based on a subset of these which were able to be confidently called. They evaluate different pipelines on low coverage data below 1× and show that applying the 2-step method using Beagle v4.1 or GLIMPSE for calling genotypes and Beagle v5 for imputation gave similar overall accuracy as using the single-step methods Beagle v4.0 or GLIMPSE, but that the 2-step method gave more nuanced posterior genotype probabilities which allowed for a more informed postimputation filtering procedure.

Studies comparing different phased reference panels have yielded varying results, with some finding that highest performance is gained by using population-specific panels (Pistis *et al.* 2015; Mitt *et al.* 2017) and others indicating the benefits of a large reference with a high level of diversity, particularly for admixed populations with no clearly matching reference (Huang *et al.* 2009; Howie *et al.* 2011; Jostins *et al.* 2011). Using a phased genome reference panel from present-day individuals is currently the only option for ancient data, introducing a possible source of bias as it means that only variants that exist in the population today can be reproduced. Leveraging information in other ancient individuals by increasing study sample size may be a way to mitigate reference divergence, particularly as more sequenced ancient genomes become available, but benefits may be diminished in the context of sparse and uncertain data. The behavior of genotype likelihood-based imputation methods on data that exhibits

the characteristic properties of aDNA discussed above has not been fully explored, particularly in combination with low coverage levels below 1×.

Genotype imputation has previously been performed on ancient human data in e.g. Gamba *et al.* (2014), Jones *et al.* (2015), Martiniano *et al.* (2017), Antonio *et al.* (2019), and Cassidy *et al.* (2020). In these studies, imputation was mainly performed using Beagle v4.0 and used to maximize the information content of ancient samples and allow for analyses such as Runs of Homozygosity (RoH) that require dense, diploid genotypes. Performance evaluation was mainly done by comparison of genotypes imputed from masked data to corresponding high-coverage calls, and showed satisfactory concordance to motivate the use of imputed data for downstream population-genetic analyses. The goal of this study is to complement and extend previous work by performing a systematic evaluation of a commonly used genotype likelihood-based imputation pipeline on ancient data. We investigate how the particular issues of sample quality and reference divergence associated with aDNA affect imputation, focusing on practical considerations regarding methodology and performance evaluation.

Materials and methods

Data description and preprocessing

Ancient data

The ancient genome data used in this study consisted of five individuals for which high-coverage data between 19× and 57× was available (ans17, LBK, Loschbour, sf12, and ne1), as well as a set of 61 individuals with low- to moderate-coverage ranging from 0.1× to 16×. See [Supplementary Tables 1 and 2](#) for sample specifics and references to source publications.

The Genome Analysis Toolkit (GATK) v3.5.0 (McKenna *et al.* 2010) tool UnifiedGenotyper was used to generate genotype likelihoods from alignment data for each of the ancient samples individually. The allele callset used was that of the 1000 Genomes phase 3 panel (Auton *et al.* 2015), filtered to keep only autosomal, biallelic SNPs, resulting in a total of 77,818,182 markers. To avoid introducing a possible bias from nucleotide misincorporations due to postmortem damage, the generated VCF files were filtered to exclude all sites where the most likely genotype could have been inferred from a deaminated allele. For C → T deaminations, this was done by removing sites where the SNP was a C ↔ T transition and the most likely genotype contained a T allele. The corresponding treatment was performed for G → A deaminations. The software bcftools v1.6 (Li *et al.* 2009) was used for filtering.

As in previous studies considering aDNA, assessment of imputation performance was done by comparison of imputation results to corresponding high-quality (HQ) genotypes. For this, the five individuals for which high-coverage data were available were used. The HQ genotypes considered as gold standard were called from the original high-coverage data, following the same pipeline as described in their respective publications ([Supplementary Table 1](#)). The called genotypes were filtered to keep only biallelic SNPs with a minimum depth of 15 and a QUAL score of at least 50. Heterozygote sites were further filtered for both alleles having a minimum allele depth of 25% of the total depth. Low-coverage data for each evaluation individual was generated by downsampling reads using Picard version 2.0.1 (Broad Institute 2019), after which estimation of genotype likelihoods and filtering were performed in an identical manner to the method described above for the low- to moderate-coverage

individuals. The sparse data was used as input for imputation, and the resulting genotypes were compared to HQ data.

Reference panels

The reference material used for imputation was the 1000 Genomes phase 3 v5a panel of phased haplotypes and the GRCh37 genetic maps provided along with the Beagle software. The reference panel was filtered to only include biallelic SNPs, resulting in 27,904,756 markers over chromosomes 1–22. To evaluate the effects of the reference on imputation, two panels were considered: the entire data set of 2,504 samples, and a smaller one of 254 samples from European populations only. The two panels are denoted as “FULL” and “EUR,” respectively. The former was also used to estimate minor allele frequency (MAF) of SNPs when analysis of imputed data over the allele frequency spectrum was performed.

Imputation methodology

Imputation was performed using Beagle v4.0, with sample data split into segments of 50,000 markers with an overlap of 25,000, using genotype likelihoods as input. Version 4.0 was selected as the goal was to evaluate imputation based on probabilistic input in the context of extremely low coverages, and this is the latest version of the software that allows imputation to be performed based on genotype likelihoods. Further, as Beagle v4.0 uses a haplotype model that is constructed from both the phased reference panel as well as current estimates of the haplotypes of the study sample (Browning 2008), we also wanted to assess the effects of including other ancient individuals in the study sample, allowing the imputation to be informed by other low-coverage ancient data within the same probabilistic framework. The effects of including multiple ancient individuals in the imputation study sample were evaluated by performing imputation jointly as well as separately. In the first case, all 61 low-to-moderate coverage individuals as well as the five evaluation individuals were included in the study sample to impute. In the second case, imputation was performed separately per study individual, meaning only the phased reference haplotypes were used in the genotype estimation.

Performance evaluation metrics

Genotype concordance

The main metric employed for assessing imputation accuracy was genotype concordance/discordance, defined as the fraction of genotypes that were imputed correctly/incorrectly. This was measured separately for each of the five evaluation individuals by comparison to the HQ genotypes derived from dense data. For each of the five evaluation individuals, we extracted the imputed markers for which there were corresponding HQ genotypes available, and divided them into two disjoint sets denoted as “overlapping” and “nonoverlapping,” based on whether or not the corresponding downsampled individual data had overlapping reads for the site or not. Supplementary Table 3 shows the sizes of these sets for each evaluation individual and coverage level.

Reference affinity

To assess whether imputed genotypes show a systematic bias toward the reference panel, the degree to which a sample showed an affinity toward the reference was compared between imputed and HQ genotypes. Here, reference affinity was measured as the fraction of markers that have the same genotype as the most frequently occurring one in the reference panel.

Principal component analysis

Principal component analysis (PCA) is a method of projecting data onto a basis that maximizes the variance of the data, possibly revealing previously unseen patterns or features. PCA can be used to reduce the dimensionality of data, e.g. for visualization purposes, and in the field of aDNA it is commonly used to show ancient samples in the context of modern variation. Here, it was used as a means of illustrating the difference between imputed and corresponding high-coverage genotypes.

PCA was performed on diploid genotypes, with a modern panel consisting of 429 European samples from the Human Origins data set of Patterson et al. (2012), filtered to remove variants with MAF under 1% or missing call rates exceeding 10%. To handle the fact that the ancient samples did not have observed genotypes for all sites used in the PCA, we used the method of known data regression (KDR) (Arteaga and Ferrer 2002). A reference PCA based on the modern panel was initially defined. Estimation of scores for each ancient sample based on this model proceeded by considering the data of the reference samples corresponding to observed ancient genotypes, and fitting a linear regression model to their original PCA scores in the reference model. The software SMARTPCA from EIGENSOFT 7.2.1 was used to define the reference PCA, and the Python library scikit-learn was used for solving linear least-squares problems.

Results

Effects of reference panel and study sample size

Imputation performance was assessed for three combinations of reference panel and study sample size, denoted as imputation configurations and shown in Table 1. In the first configuration, imputation was performed individually per evaluation individual, using the EUR phased reference panel. For configurations 2 and 3, all ancient individuals were included in the imputation study sample, using the EUR and FULL reference panels, respectively. All results in this section are for imputation performed on sample data with 1× coverage, and to perform a comprehensive evaluation of the effect of imputation configuration, no posterior filter was imposed on the imputed genotypes. Results are shown for each of the five evaluation individuals, with results split into overlapping and nonoverlapping marker sets.

Figure 1 shows genotype concordance for the three imputation configurations. Overall, concordance rates were similar between individuals and reached 0.99 in all cases. The results indicate that the larger reference panel as well as the inclusion of ancient individuals in the imputation study sample improved performance. For overlapping markers, concordance rates increased slightly from just under 0.997 to somewhat above. Concordance was lower among the nonoverlapping markers in general, and it was also among these that the effects of varying imputation configurations were the most pronounced.

To investigate imputation performance over the allele frequency range, the imputed markers were binned according to MAF, and the average genotype discordance was assessed per bin. Because there is a high risk of chance agreement between

Table 1. Imputation configurations.

Configuration	Reference	Study sample
1	EUR	Single
2	EUR	All ancient
3	FULL	All ancient

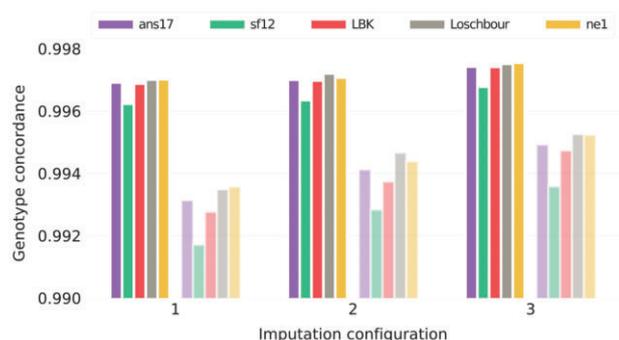


Fig. 1. Concordance of imputed genotypes for the three configurations in Table 1. Imputation was performed on data in which the five evaluation individuals were downsampled to 1× coverage, and the evaluation was based on unfiltered results, with fully colored bars showing results for sites at which the downsampled data had overlapping reads, and shaded bars for nonoverlapping markers, i.e. sites at which all overlapping reads were removed by the downsampling process.

homozygous genotypes and the reference majority in the case of low MAF markers, only sites at which the HQ data were heterozygous were considered, thus measuring the ability of the imputation to recover heterozygotes. Results for imputation configurations 1–3 are shown in Fig. 2, and again indicate that using a larger phased reference panel and study sample size increases performance. The effects were particularly visible among nonoverlapping markers and in the low MAF ranges, where heterozygotes are the most difficult to recover.

Effects of coverage

Next, we assessed the effects of coverage on imputation. For every level $C \in \{0.10, 0.25, 0.50, 0.75, 1.00, 1.25, 1.50, 1.75, 2.00\}$, an imputation run was performed using the five evaluation individuals downsampled to C ×. Imputation configuration 2 was used for all runs, and as in the previous section, no posterior filter was imposed on the resulting genotypes.

Figure 3 shows the concordance between imputed and HQ genotypes for all markers (A) and for heterozygotes (B), separated into overlapping and nonoverlapping marker sets for every evaluation sample and coverage level. A total concordance rate of 0.99 was reached around 0.25× for most individuals, with a plateau visible at 1× where levels around 0.9975 and 0.995 were obtained for overlapping and nonoverlapping sites, respectively. For heterozygote sites, concordance levels were below 0.975 throughout, reaching 0.95 at 1× and as low as 0.825 for the lowest coverage level of 0.1×.

To assess the level of systematic bias toward the variant that is most common in the reference, we compared the level of reference affinity of the imputed data to that of the corresponding HQ genotypes, considering these as a baseline for the similarity between the true genotypes and the reference. Figure 4 shows the difference in measured affinities between the HQ and imputed data for the nine coverage levels, over the allele frequency spectrum, averaged over the five evaluation individuals.

The negative values in the lower MAF ranges indicate that the imputed data shows a larger affinity toward the reference. While the extremely low coverages showed significantly higher levels of bias, the rates decreased and showed little variation for coverages 0.75× and higher. The results indicate a reduction in bias toward the reference with increasing MAF, showing little differences around MAF 0.3 for most coverage levels. A possible explanation for the positive difference at higher MAF values is that at these

markers, imputation errors do not tend toward the reference majority as strongly as in the lower MAF ranges. Overall concordance is also lower at high MAF due to higher frequency of heterozygotes.

In-depth performance analysis

In this section, we present further evaluation of imputed genotypes to assess properties relevant to downstream analysis. We considered results of imputation configuration 3 on data with 1× coverage, and imposed a filter of minimum posterior genotype probability of 0.99 on the imputed data. First, performance for different genotypes was evaluated. Figure 5 shows concordance of sites split according to genotype in the high-coverage data. Although performance remained poorer for heterozygote sites, the filtered data showed improved levels of over 0.99 throughout. The filtered results also showed little difference between overlapping and nonoverlapping marker sets at homozygote sites, while larger differences remained for heterozygotes. Inspection of performance for heterozygote sites across the allele frequency spectrum showed that discordance levels below 0.01 were reached around MAF 0.1, with over 85% of sites retained postfilter (Fig. 6).

Finally, PCA was used to visualize and compare genetic affinities of imputed and high-coverage data. Figure 7 shows that, within the variation represented by the first two principal components, scores of imputed genotypes map closely to those of the HQ data, particularly for the individuals ans17, LBK, and ne1. As illustrated in e.g. Günther and Jakobsson (2019), considering subsets of SNPs introduces noise in the PCA projection, resulting in less accurate projections with higher degree of uncertainty compared to using all available data. In Fig. 7, this is visualized by the distance between low-coverage and HQ data points for each individual. For each of the five individuals, relative distances to the HQ data are smaller for the imputed genotypes than the low-coverage data, indicating that information lost by the downsampling process has been retained by imputation.

Discussion

This study provides a systematic investigation of the application of imputation to human aDNA. We have corroborated results from similar experiments showing that overall genotype concordance levels of over 0.99 can be reached for data with 1× coverage, and provided an in-depth analysis of the qualities of imputed genotypes. Investigation of performance at coverage levels tending to the ultra-low indicated that the quality of imputed genotypes began to plateau around 0.75–1×, and that common variants were more robust to different reference panels as well as less prone to reference bias overall. The results suggest a MAF threshold of around 0.1 for minimizing genotype discordance of heterozygote calls, and that in cases where rare alleles are of interest, an increased diversity and size of the phased reference as well as the imputation study sample are particularly beneficial. The fact that information from ancient individuals can be leveraged to improve imputation, along with the genotype concordance levels shown, may in some cases motivate the sequencing of more ancient samples to lower coverage as a cost-effective alternative to that of fewer samples to higher coverage.

The presented work provides a framework of practical considerations for performing imputation, as well as a basis for further investigations. A systematic evaluation of different imputation methods and adaptation of the statistical models to the context of sparse and uncertain data may increase performance for

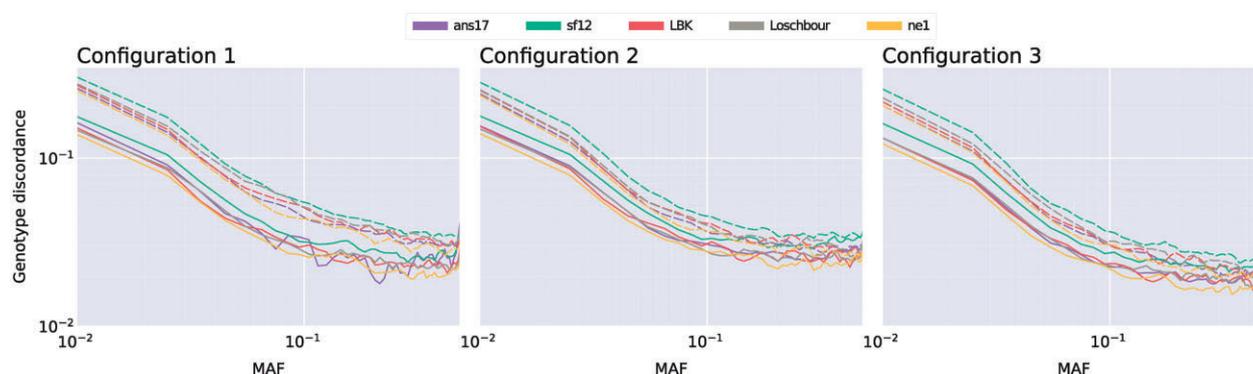


Fig. 2. Log-log plots displaying discordance of heterozygote genotypes for the five evaluation individuals. The subplots show results obtained using imputation configurations 1–3 described in Table 1, averaged over markers in 50 MAF bins. Results are shown for unfiltered imputed genotypes of data downsampled to $1\times$ coverage, with solid lines indicating overlapping markers, and dashed lines markers at which the downsampled data had no overlapping reads.

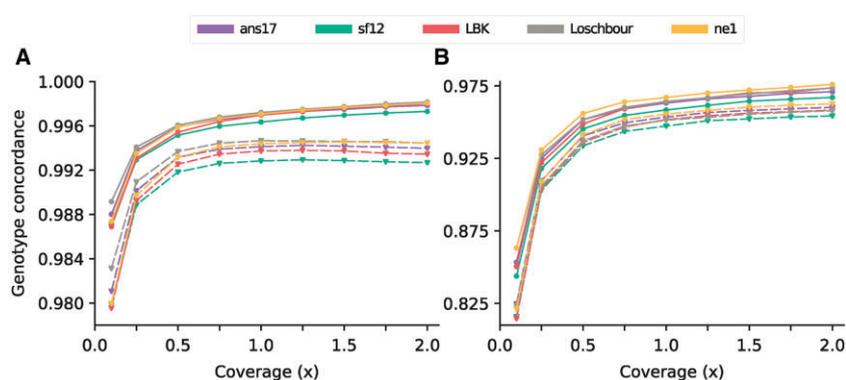


Fig. 3. Genotype concordance of imputed data for different levels of coverage of the input data, for (A) all markers and (B) markers at which the HQ genotype was heterozygote. Results are shown for the five evaluation individuals, with solid and dashed lines indicating SNPs with and without overlapping reads in the downsampled data, respectively. Imputation was performed using configuration 2 (Table 1), with no filter on posterior genotype probability imposed on the imputed genotypes.

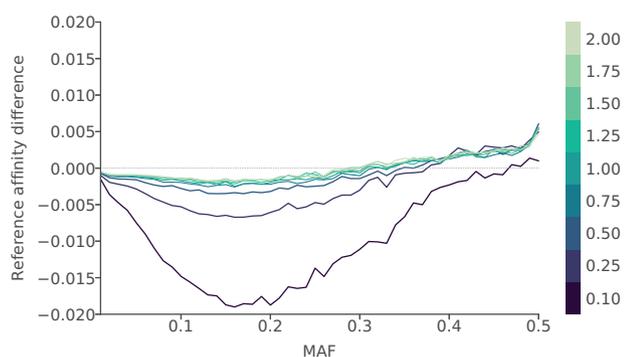


Fig. 4. Comparison of levels of reference bias for different coverage levels. Input data for imputation in which the evaluation individuals were downsampled to coverages ranging between 0.1 and $2.0\times$ was generated (displayed in different colors), and the difference in reference affinity between HQ and resulting imputed genotypes estimated, with negative values indicating increased reference affinity of the imputed genotypes compared to the HQ genotypes. Imputation was performed using configuration 2 (Table 1), with no filter on posterior genotype probability imposed on the imputed genotypes. Results are shown for all markers, aggregated into 50 MAF bins, and averaged over the evaluation individuals.

ancient samples. Further insights may also be gained by assessment of imputed data by means of different population genetic analyses, considering e.g. haplotype-based methods such as RoH as well as those based on allele frequencies.

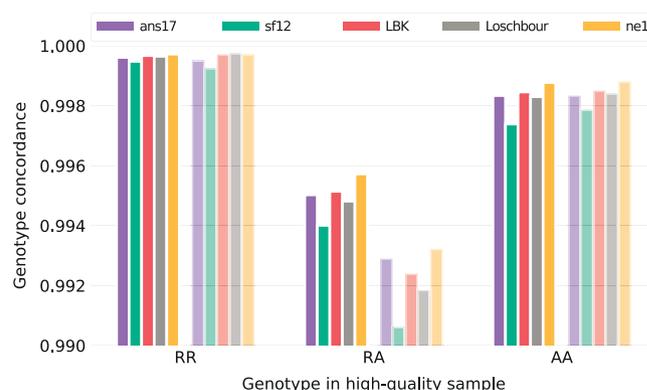


Fig. 5. Concordance of imputed genotypes, split according to genotype in the HQ data. Imputation was performed on data in which the five evaluation individuals were downsampled to $1\times$ coverage, using imputation configuration 3 (Table 1), and the resulting data was filtered for minimum genotype probability of 0.99. Fully colored bars indicate markers that had overlapping reads in the downsampled data, and shaded bars indicate sites that did not.

For the five evaluation individuals considered, in-depth performance analysis showed lower performance for the hunter-gatherer genomes sf12 and Loschbour, both in terms of genotype concordance and similarity in PCA-space to HQ data. A possible explanation is that the reference panel of present-day individuals

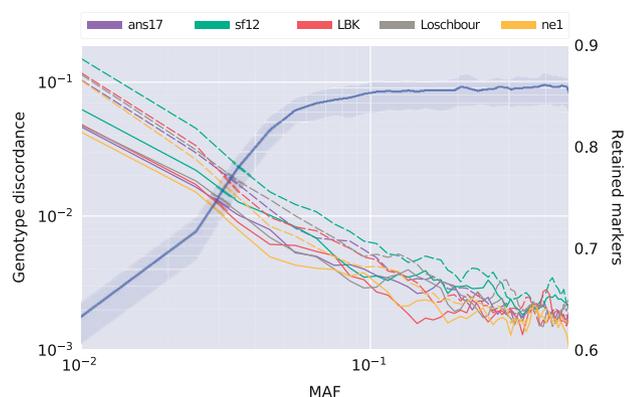


Fig. 6. Log-log plot of genotype discordance at heterozygote sites, averaged over 50 bins in the allele frequency spectrum and smoothed using a moving average using 3 points. Input data for imputation was downsampled to 1× and configuration 3 (Table 1) was used, after which a posterior filter of minimum genotype probability of 0.99 was applied. The fraction of markers retained after the filter, averaged over the five evaluation individuals, is shown in blue, with shaded regions indicating minimum and maximum. Solid and dashed lines indicate SNPs with and without overlapping reads in the downsampled data, respectively.

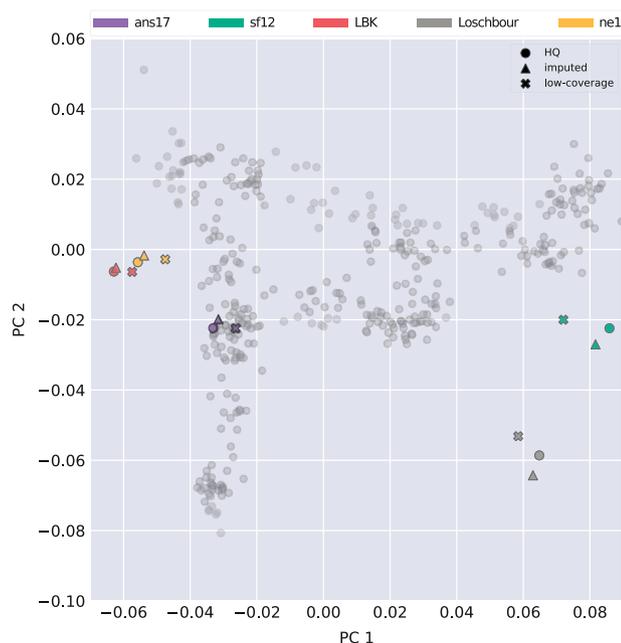


Fig. 7. PCA comparing HQ, imputed, and low-coverage data for the five evaluation individuals. A reference PCA was defined based on genotypes of modern European samples from the Human Origins (Patterson et al. 2012) data set, after which the KDR method was used to estimate scores of ancient individuals. Modern individuals are indicated by gray dots and ancient individuals colored according to the legend. Not all modern individuals are visible in the plot, see Supplementary Fig. 1 for a view of the full data set. The low-coverage data points correspond to the genotypes that have been downsampled to 1×, filtered, and used as input to imputation, which was performed using configuration 3 (Table 1), with a posterior filter of minimum genotype probability of 0.99 applied.

used for imputation does not contain individuals with a genetic composition that is similar to these individuals. As discussed in e.g. Skoglund et al. (2014) and Günther and Jakobsson (2016), hunter-gatherer individuals have shown a particular genetic profile that is not represented in the genetic variation of modern-day

people. The individuals LBK, ne1, and ans17, in contrast, are from early farming cultures, for which there are present-day European individuals who share a similar genetic make-up. Further investigations of the effects of sample ancestry and reference divergence will be required to customize imputation pipelines to individuals of varying genetic composition. As more sequenced ancient data becomes available, both imputation accuracy as well as the ability to assess various aspects of performance will be improved.

A possible source of bias in the evaluation of imputation performance is reference bias introduced by alignment and variant calling procedures on ancient sequence data. This can lead to alternate alleles being missed and heterozygous sites to be called as homozygous for the reference allele (Günther and Nettelblad 2019), causing ancient individuals to appear artificially similar to the modern reference genome. As imputation is expected to perform better for samples that are genetically close to the phased reference panel used, this could cause overestimated measures of accuracy. A possible means to mitigate this effect is the use of genotype callers that are tailored for aDNA such as snpAD (Prüfer 2018) and ATLAS (Link et al. 2017). As these callers also account for postmortem damage, the inclusion of transition-type polymorphisms in future imputation studies may also be motivated.

In this study, we have focused on a probabilistic imputation framework based on Beagle v4.0 that has been frequently used in the aDNA community. Other imputation pipelines have shown to give qualitatively similar results for imputing genotypes from ancient sequence data, and also indicated some performance trade-offs such as imputation accuracy at different parts of the allele frequency spectrum (Hui et al. 2020). Another relevant point is that more recent software for phasing and imputation such as GLIMPSE and Beagle v5 have increased focus on scalability to larger reference panels. Sample size, availability of computational resources and the type of downstream analysis intended are thus factors to consider in the selection of methodology for imputation of aDNA.

Data availability

The authors state that all data necessary for confirming the conclusions presented in the article are represented fully within the article and in its online [supplementary material](#).

[Supplementary material](#) is available at G3 online.

Acknowledgments

The authors acknowledge the use of computational resources provided by Swedish National Infrastructure for Computing (SNIC) at Uppsala Multidisciplinary Center for Advanced Computational Science (UPPMAX) under project SNIC 2017/7-162.

Funding

CN acknowledges funding by Formas (grant number 2020-00712) and MJ acknowledges the Knut and Alice Wallenberg foundation.

Conflicts of interest

None declared.

Literature cited

- Antonio ML, Gao Z, Moots HM, Lucci M, Candilio F, Sawyer S, Oberreiter V, Calderon D, Devitofranceschi K, Aikens RC, et al. Ancient Rome: a genetic crossroads of Europe and the Mediterranean. *Science*. 2019;366(6466):708–714.
- Arteaga F, Ferrer A. Dealing with missing data in MSPC: several methods, different interpretations, some examples. *J Chemometrics*. 2002;16(8–10):408–418.
- Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO, Marchini JL, McCarthy S, McVean GA, Abecasis GR, et al.; 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature*. 2015;526(7571):68–74.
- Briggs AW, Stenzel U, Johnson PLF, Green RE, Kelso J, Prüfer K, Meyer M, Krause J, Ronan MT, Lachmann M, et al. Patterns of damage in genomic DNA sequences from a neandertal. *Proc Natl Acad Sci USA*. 2007;104(37):14616–14621.
- Broad Institute. 2019. Version 2.0.1 Picard tools. <http://broadinstitute.github.io/picard/>.
- Brotherton P, Endicott P, Sanchez JJ, Beaumont M, Barnett R, Austin J, Cooper A. Novel high-resolution characterization of ancient DNA reveals c > u-type base modification events as the sole cause of post mortem miscoding lesions. *Nucleic Acids Res*. 2007;35(17):5717–5728.
- Browning BL, Yu Z. Simultaneous genotype calling and haplotype phasing improves genotype accuracy and reduces false-positive associations for genome-wide association studies. *Am J Hum Genet*. 2009;85(6):847–861.
- Browning SR. Missing data imputation and haplotype phase inference for genome-wide association studies. *Hum Genet*. 2008;124(5):439–450.
- Browning SR, Browning BL. Rapid and accurate haplotype phasing and missing-data inference for whole-genome association studies by use of localized haplotype clustering. *Am J Hum Genet*. 2007;81(5):1084–1097.
- Browning SR, Browning BL. Haplotype phasing: existing methods and new developments. *Nat Rev Genet*. 2011;12(10):703–714.
- Cassidy LM, Maoldúin RÓ, Kador T, Lynch A, Jones C, Woodman PC, Murphy E, Ramsey G, Dowd M, Noonan A, et al. A dynastic elite in monumental neolithic society. *Nature*. 2020;582(7812):384–388.
- Gamba C, Jones ER, Teasdale MD, McLaughlin RL, Gonzalez-Fortes G, Mattiangeli V, Domboróczki L, Kővári I, Pap I, Anders A, et al. Genome flux and stasis in a five millennium transect of European prehistory. *Nat Commun*. 2014;5:5257.
- Ginolhac A, Rasmussen M, Gilbert MTP, Willerslev E, Orlando L. mapdamage: testing for damage patterns in ancient DNA sequences. *Bioinformatics*. 2011;27(15):2153–2155.
- Günther T, Jakobsson M. Genes mirror migrations and cultures in prehistoric Europe—a population genomic perspective. *Curr Opin Genet Dev*. 2016;41:115–123. *Genetics of human origin*.
- Günther T, Jakobsson M. Population genomic analyses of DNA from ancient remains. In: DJ Balding, I Moltke, J Marioni, editors. *Handbook of Statistical Genomics*. Chapter 10. Hoboken (NJ): Wiley; 2019. p. 295–324.
- Günther T, Nettelblad C. The presence and impact of reference bias on population genomic studies of prehistoric human populations. *PLoS Genet*. 2019;15(7):e1008302.
- Howie B, Marchini J, Stephens M. Genotype imputation with thousands of genomes. *G3 (Bethesda)*. 2011;1(6):457–470.
- Howie BN, Donnelly P, Marchini J. A flexible and accurate genotype imputation method for the next generation of genome-wide association studies. *PLoS Genet*. 2009;5(6):e1000529.
- Huang L, Li Y, Singleton AB, Hardy JA, Abecasis G, Rosenberg NA, Scheet P. Genotype-imputation accuracy across worldwide human populations. *Am J Hum Genet*. 2009;84(2):235–250.
- Hui R, D’Atanasio E, Cassidy LM, Scheib CL, Kivisild T. Evaluating genotype imputation pipeline for ultra-low coverage ancient genomes. *Sci Rep*. 2020;10(1):18542.
- Jones ER, Gonzalez-Fortes G, Connell S, Siska V, Eriksson A, Martiniano R, McLaughlin RL, Gallego Llorente M, Cassidy LM, Gamba C, et al. Upper palaeolithic genomes reveal deep roots of modern Eurasians. *Nat Commun*. 2015;6:8912.
- Jostins L, Morley KI, Barrett JC. Imputation of low-frequency variants using the hapmap3 benefits from large, diverse reference sets. *Eur J Hum Genet*. 2011;19(6):662–666.
- Krause J, Briggs AW, Kircher M, Maricic T, Zwyns N, Derevianko A, Pääbo S. A complete mtDNA genome of an early modern human from Kostenki, Russia. *Curr Biol*. 2010;20(3):231–236.
- Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R; 1000 Genome Project Data Processing Subgroup. The sequence alignment/map format and samtools. *Bioinformatics*. 2009;25(16):2078–2079.
- Li N, Stephens M. Modeling linkage disequilibrium and identifying recombination hotspots using single-nucleotide polymorphism data. *Genetics*. 2003;165(4):2213–2233.
- Li Y, Willer CJ, Ding J, Scheet P, Abecasis GR. Mach: using sequence and genotype data to estimate haplotypes and unobserved genotypes. *Genet Epidemiol*. 2010;34(8):816–834.
- Link V, Kousathanas A, Veeramah K, Sell C, Scheu A, Wegmann D. Atlas: analysis tools for low-depth and ancient samples. 2017. bioRxiv.
- Marchini J, Howie B. Genotype imputation for genome-wide association studies. *Nat Rev Genet*. 2010;11(7):499–511.
- Martiniano R, Cassidy LM, ÓMaoldúin R, McLaughlin R, Silva NM, Manco L, Fidalgo D, Pereira T, Coelho MJ, Serra M, et al. The population genomics of archaeological transition in west Iberia: investigation of ancient substructure using imputation and haplotype-based methods. *PLoS Genet*. 2017;13(7):e1006852.
- McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D, Gabriel S, Daly M, et al. The genome analysis toolkit: a mapreduce framework for analyzing next-generation DNA sequencing data. *Genome Res*. 2010;20(9):1297–1303.
- Mitt M, Kals M, Pärn K, Gabriel SB, Lander ES, Palotie A, Ripatti S, Morris AP, Metspalu A, Esko T, et al. Improved imputation accuracy of rare and low-frequency variants using population-specific high-coverage WGS-based imputation reference panel. *Eur J Hum Genet*. 2017;25(7):869–876.
- Nielsen R, Akey JM, Jakobsson M, Pritchard JK, Tishkoff S, Willerslev E. Tracing the peopling of the world through genomics. *Nature*. 2017;541(7637):302–310.
- Nielsen R, Paul JS, Albrechtsen A, Song YS. Genotype and SNP calling from next-generation sequencing data. *Nat Rev Genet*. 2011;12(6):443–451.
- Parks M, Lambert D. Impacts of low coverage depths and post-mortem DNA damage on variant calling: a simulation study. *BMC Genomics*. 2015;16:19.
- Patterson N, Moorjani P, Luo Y, Mallick S, Rohland N, Zhan Y, Genschoreck T, Webster T, Reich D. Ancient admixture in human history. *Genetics*. 2012;192(3):1065–1093.
- Pistis G, Porcu E, Vrietze SI, Sidore C, Steri M, Danjou F, Busonero F, Mulas A, Zoledziewska M, Maschio A, et al. Rare variant genotype imputation with thousands of study-specific whole-genome sequences: implications for cost-effective study designs. *Eur J Hum Genet*. 2015;23(7):975–983.

- Prüfer K, Stenzel U, Hofreiter M, Pääbo S, Kelso J, Green RE. Computational challenges in the analysis of ancient DNA. *Genome Biol.* 2010;11(5):R47.
- Prüfer K. snpAD: an ancient DNA genotype caller. *Bioinformatics.* 2018;34(24):4165–4171.
- Pääbo S, Poinar H, Serre D, Jaenicke-Despres V, Hebler J, Rohland N, Kuch M, Krause J, Vigilant L, Hofreiter M, et al. Genetic analyses from ancient DNA. *Annu Rev Genet.* 2004;38:645–679.
- Rubinacci S, Ribeiro DM, Hofmeister RJ, Delaneau O. Efficient phasing and imputation of low-coverage sequencing data using large reference panels. *Nat Genet.* 2021;53(1):120–126.
- Sawyer S, Krause J, Guschanski K, Savolainen V, Pääbo S. Temporal patterns of nucleotide misincorporations and DNA fragmentation in ancient DNA. *PLoS One.* 2012;7(3):e34131.
- Skoglund P, Malmström H, Omrak A, Raghavan M, Valdiosera C, Günther T, Hall P, Tambets K, Parik J, Sjögren K-G, et al. Genomic diversity and admixture differs for stone-age Scandinavian foragers and farmers. *Science.* 2014;344(6185):747–750.
- Spencer CCA, Su Z, Donnelly P, Marchini J. Designing genome-wide association studies: sample size, power, imputation, and the choice of genotyping chip. *PLoS Genet.* 2009;5(5):e1000477.
- Stephens M, Scheet P. Accounting for decay of linkage disequilibrium in haplotype inference and missing-data imputation. *Am J Hum Genet.* 2005;76(3):449–462.
- Stephens M, Smith NJ, Donnelly P. A new statistical method for haplotype reconstruction from population data. *Am J Hum Genet.* 2001;68(4):978–989.
- Stiller M, Green RE, Ronan M, Simons JF, Du L, He W, Egholm M, Rothberg JM, Keates SG, Keats SG, et al. Patterns of nucleotide misincorporations during enzymatic amplification and direct large-scale sequencing of ancient DNA. *Proc Natl Acad Sci USA.* 2006;103(37):13578–13584.
- Sánchez-Quinto F, Schroeder H, Ramirez O, Avila-Arcos MC, Pybus M, Olalde I, Velazquez AMV, Marcos MEP, Encinas JMV, Bertranpetit J, et al. Genomic affinities of two 7,000-year-old Iberian hunter-gatherers. *Curr Biol.* 2012;22(16):1494–1499.
- Zeggini E, Scott LJ, Saxena R, Voight BF, Marchini JL, Hu T, de Bakker PIW, Abecasis GR, Almgren P, Andersen G, et al.; Wellcome Trust Case Control Consortium. Meta-analysis of genome-wide association data and large-scale replication identifies additional susceptibility loci for type 2 diabetes. *Nat Genet.* 2008;40(5):638–645.

Communicating editor: A. Sethuraman