

## An Ancient DNA Damage Toolkit

Christian Stentoft Michelsen<sup>1</sup>  , Mikkel Winther Pedersen<sup>2</sup>  , Antonio

<sup>4</sup> Fernandez-Guerra<sup>2</sup> , Lei Zhao<sup>2</sup>, Troels C. Petersen<sup>1</sup> , Thorfinn Sand

Korneliussen<sup>2</sup>  

 For correspondence:

[christianmichelsen@gmail.com](mailto:christianmichelsen@gmail.com)

(CM); [mwpedersen@sund.ku.dk](mailto:mwpedersen@sund.ku.dk)

(MW);

[tskorneliussen@sund.ku.dk](mailto:tskorneliussen@sund.ku.dk)

(TSK)

<sup>6</sup> <sup>1</sup> Niels Bohr Institute, University of Copenhagen; <sup>2</sup> Globe Institute, University of Copenhagen

<sup>8</sup> 

<sup>†</sup>Authors contributed equally.

### Abstract

**Present address:** Niels Bohr Institute, University of Copenhagen, Blegdamsvej 17, 2100 Copenhagen, Denmark

**Data availability:** Data is available on [Zenodo](#) or the [Github](#) repository.

**Funding:** This work was supported by Carlsberg Foundation Young Researcher Fellowship awarded by the Carlsberg Foundation [CF19-0712], and the Lundbeck Foundation Centre for Disease Evolution: [R302-2018-2155 to L.Z]. The funders had no role in the decision to publish.

**Competing interests:** The author declare no competing interests.

**1. Motivation** Under favourable conditions DNA molecules can persist for more than two million years (Kjaer et al in press). Such genetic remains make up invaluable resources to study past assemblages, populations and even the evolution of species. However, DNA is subjected to enzymatic, chemical and mechanical degradation, and hence over time decrease to ultra low concentrations which makes it highly prone to contamination by modern sources that are rich in DNA. Strict precautions and criteria(Llamas et al., 2017; Gilbert et al., 2005; Champlot et al., 2010) are therefore necessary to ensure that DNA from modern sources does not appear in the final data and to authenticate that the DNA is ancient. The most generally accepted and widely applied authenticity for ancient DNA studies is to test for elevated deaminated cytosine residues towards the termini of the molecules - DNA damage (Briggs et al., 2007; Dabney, Meyer, and Pääbo, 2013). To date this has primarily been used for single organisms (Jónsson et al., 2013)) and recently for read assemblies (Borry et al., 2021), however these methods have not been designed, nor are they computationally up-scalable for calculating DNA damage for ancient metagenome with tens-hundreds of thousands of species.

**2. Methods** Here we present metaDMG, a novel framework that takes advantage of the information already contained within standard alignment files to compute and statically

evaluate misincorporations due to DNA damage. It thus bypasses any need for initial  
28 classification, splitting reads by individual organisms, realigning these to the reference genome and lastly parse alignments to mapDamage2.0 (Jónsson et al., 2013). We  
29 furthermore, implemented a Bayesian approach that combines a geometric damage profile with a beta-binomial model to fit the entire model to the individual misincorporations at all  
30 taxonomic levels. metaDMG were hereafter benchmarked using simulated data of single  
31 genomes, metagenomes and tested on published datasets, before comparing its  
32 performance with pydamage.

3. **Results** we find metaDMG to be a factor 10 faster than previous methods and more accurate  
33 even for complex metagenomes with tens of thousands of species. Our simulation show  
34 that metaDMG can estimate DNA damage at taxonomic levels with less than 100 reads and  
35 that uncertainties decreases with increased number of reads, but also that our estimate is  
36 more significant with increased C-T misincorporations.  
37 we also show that sequencing errors does not influence our ability to estimate DNA  
38 damage.  
39 improve the damage estimates in precision and time compared to previous methods.

4. **Conclusion** metaDMG is a state-of-the-art programme suite for computing DNA damage  
40 estimation, nucleotide misincorporation, gc-content, DNA fragmentation of simple and  
41 complex ancient genomic datasets. Additionally it includes PMDtool statistics (Skoglund  
42 et al., 2014) that allow for extraction of reads with damage, making it a complete tools  
43 package for ancient DNA damage authentication.

44 **keywords:** ancient DNA, DNA damage estimation, DNA damage, metaDMG, metagenomics, .

---

## 50 1 | INTRODUCTION

Throughout an organism life it contaminates its environment with DNA, cells or tissue and hereby  
51 leave genetic traces behind. As the cell leaves its host, DNA repair mechanisms stops and the DNA  
52 are subjected to intra and extra-cellular enzymatic, chemical and mechanical degradation, resulting  
53 in fragmentation and molecule alterations that over time lead to the characteristics of ancient  
54 DNA (Briggs et al., 2007; Dabney, Meyer, and Pääbo, 2013). Ancient DNA (aDNA) has been shown to

56 persist in a diverse variety of environmental contexts, e.g. within fossils such as bones, soft-tissue  
57 and hair, as well as in geological sediments, archaeological layers, ice-cores, permafrost soil for  
58 thousands and even millions of years (Kjaer et al in review), (Cappellini et al., 2018). Common for  
59 all is that they have an accumulated amount of deaminated cytosines towards the termini of the  
60 DNA strand, which when amplified for sequencing results in mis-incorporations of thymines on the  
61 cytosines (Dabney, Meyer, and Pääbo, 2013; Ginolhac et al., 2011). This postmortem damage with  
62 regards to DNA is characterized by the four Briggs parameters (Briggs et al., 2007). First, ancient  
63 DNA molecule tend to be short, and is likely to have a single stranded overhang due to hydrolytic  
64 cleavage of the  $\beta$ -N-glycosidic bond which results in a nucleotide loss and eventual breakage of the  
65 DNA. On this overhang the cytosine is especially exposed to deamination which results in the ob-  
66 served high proportion of  $C \rightarrow T$  substitutions at the single stranded overhang  $\epsilon_{ss}$ , and a somewhat  
67 higher  $C \rightarrow T$  at the double stranded part  $\epsilon_{ds}$ . The length of the single stranded part (*overhang*)  
68 follows a geometric distribution  $\lambda$ , and finally there might be breaks at the backbone in the double  
69 stranded part  $v$ . It is possible to estimate these four Briggs parameters Jónsson et al., 2013 but  
70 these four parameters are rarely used directly for asserting "ancientness", and researchers work-  
71 ing with ancient DNA tend to simply use the empirical  $C \rightarrow T$  on the first position of the fragment  
72 together with other supporting summary statistic of the experiment. Estimating mis-incorporation  
73 due to DNA damage, molecule fragmentation, and nick frequencies have become standard for sin-  
74 gle individual sources like hair, bones, teeth and also applied on small subsets of species in ancient  
75 environmental metagenomes (Mikkel W. Pedersen et al., 2016; Murchie et al., 2021; Zavala et al.,  
76 2021; Wang, Mikkel Winther Pedersen, et al., 2021). While this is a relatively fast process for single  
77 individuals it becomes increasingly demanding, iterative and time consuming as the samples and  
78 the diversity within increases, as in the case for metagenomes from ancient soil, sediments, dental  
79 calculus, coprolites and other ancient environmental sources. It has therefore been practice to es-  
80 timate damage for only the key taxa of interest in a metagenome, as a metagenomic sample easily  
81 includes tens of thousands of different taxonomic entities, that would make a complete estimate  
82 an impossible task.

To overcome this limitation, we designed a program suite metaDMG with novel test statistics that  
83 takes into account all relevant information provided alignments to both single or multiple refer-  
84 ence genomes. We find metaDMG, to estimate DNA damage both faster, more accurate and able to  
85 process complex metagenomes within hours. metaDMG is designed for and up-scales equally with

the increasingly large datasets that are currently generated in the field of ancient environmental  
88 DNA. However, it also outperforms standard tools that estimate DNA damage in single genomes  
and samples with low complexity. Furthermore, it can even compute an global damage estimate  
90 for a metagenome as a whole. Importantly, metaDMG is compatible with the NCBI taxonomy and  
use ngsLCA (add Wang et al. 2021 ngsLCA paper) to perform a naïve last common ancestor classi-  
92 fication of the aligned reads to get precise damage estimates at all taxonomic nodes. Lastly, it is  
also designed to accommodate custom taxonomies and hence metagenomic assembled genomes  
94 (MAGs) as references.

To test metaDMG specificity, sensitivity and performance we use multiple sets of simulations and  
96 test parameters which all show that metaDMG not only outperforms existing methods in the case  
of single-genome damage estimation but also enables fast and accurate estimation in ancient  
98 metagenomes. Lastly, we apply our metaDMG on a representative mix of nine different metagenomes  
that span the variety of environments hitherto published, including lakes sediments, cave sedi-  
100 ments and ancient chewed birch tar.

## 2 | METHODS & MATERIALS

102 Perhaps the most basic bioinformatic analyses is the difference between two nucleotide sequences.  
This assumes that we have a haploid representation of our target organisms and larger differences  
104 can be interpreted as larger genetic differences. Obtaining a haploid representation is none trivial,  
firstly our target organism might not be haploid and we need to construct a consensus genome,  
106 secondly data from modern day sequencers are essentially a sampling with replacement process  
and we need to infer the relative location of each of the possible millions or even billions of short  
108 DNA fragments, this is the process which is called mapping or alignment. Thirdly, and the focus for  
this manuscript, is the quantification of the presence of postmortem damage (PMD) in DNA. PMD  
110 mainly manifests as an excess of cytosine to thymine substitutions at the termini of fragments that  
has been prepared for sequencing. A priori we can not directly observe these actual biochemical  
112 changes but we can align each fragment and consider the difference between reference and read  
as possible PMD, and it is even possible to use the excess of C to T at the single fragment level to  
114 separate modern from ancient (data with PMD) (Skoglund et al., 2014). Expanding from the single  
read all reads for a sequencing experiment and genome to tabulate the overall substitution or

116 mismatch rates to obtain a statistic of the damage (Borry et al., 2021) or even estimate the four  
118 Briggs parameters that is traditionally used to characterize the damage signal (Jónsson et al., 2013).

120 We have devised a general ancient DNA damage toolkit with a special emphasis in a metagenomic setting which implements and expands existing relevant methods but also expands with  
122 several state of the art novel methodologies. At the most basic level we have reimplemented the approach given in (Skoglund et al., 2014) which allows for the extracting and separation of highly  
124 damaged DNA reads. Secondly under the assumption of vast amounts of data we have defined a full multinomial regression model building on the method in (Cabanski et al., 2012), we show that  
126 this will give superior and stable results if it is possible to obtain high depth and coverage data.  
128 However, in standard ancient DNA context it is generally not possible to obtain vast amounts of data and we propose two novel tests statistics that is especially suited for this scenario. To our knowledge there are no currently available methods that is geared towards damage analysis in a  
130 metagenomic setting and existing approaches are essentially based on remapping against the single target organism and does not take into account any possible issues with regards to reads being  
132 well assigned or specified. Our solution called `metaDMG` (pronounced metadamage), estimates the damage patterns in metagenomic samples in a three step approach. First, the lowest common  
134 ancestor (LCA) for each read (mapped to a multi-species reference database) is computed and the mismatch matrix for each leaf node (e.g. taxonomic ID or contig, depending on the database  
136 used) is computed based on the mapped reads. Second, `metaDMG` fits a damage model to each leaf node to compute the ancient damage estimates. Finally, the results are visualized in the `metaDMG` dashboard, which is a state of the art graphical user interface that allows for fast and user-friendly interaction with the results for further downstream analysis and visualization.

## 138 2.1 | Lowest Common Ancestor and Mismatch matrices

For environmental DNA (eDNA) studies we routinely apply a competitive alignment approach where  
140 we consider all possible alignments for a given read. Each read is mapped against a multi species reference databases (e.g. nucleotide or RefSeq from NCBI or custom downloaded). A single read  
142 might map to a highly conserved gene that is shared across higher taxonomic ranks such as class or even domains. This read will not provide relevant information due to the generality, whereas a  
144 read that maps solely to a single species or species from a genus would be indicative of the read being well classified. We seek to obtain the pattern or signal of damage which is done by the tab-

146 ulation of the cycle specific mismatch rates between our reference and observed sequence for all  
well classified reads.

148 In details we compute the lowest common ancestor (lca) for all alignments for each read,  
this is done using (Wang, T. S. Korneliussen, et al., 2022) and if a read is well classified or properly  
150 assigned based on a user defined threshold (specics, genus or family) we tabulate the mismatches  
for each cycle, if a read is not well assigned it is discarded. Pending on the run mode we allow  
152 for the construction of these mismatch tables on three different lavel. Either we obtain a basic  
single global mismatch matrix, which could be relevant in a standard single genome aDNA study  
154 and similar to the tabulation used in (Jónsson et al., 2013). Secondly we can obtain per reference  
counts or if a taxonomy database has been supplied we allow for the aggregation from leaf nodes  
156 to the internal taxonomic ranks towards the root.

To suit as many users as possible, metaDMG takes as input an alignment file (.bam, .sam, or  
158 .sam.gz), where Each read is hereafter allowed an equal chance to map against the multiple refer-  
ences. One read can therefore attract multiple alignments, and we thus first seek to find the lowest  
160 common ancestor (LCA) among the alignments based on the tree structure from the databases and  
a user defined read-reference similarity interval (Wang, T. S. Korneliussen, et al., 2022). Note that  
162 metaDMG is not limited to the NCBI database and allow for custom databases as well.

Regardless of runmode or weighing scheme used in the possible aggregation we obtain the  
164 nucleotide substitution frequencies across reads which provides us with the position dependent  
mismatch matrices,  $\underline{M}(x)$ , with  $x$  denoting the position in the read, starting from 1. At a specific  
166 position,  $M_{\text{ref} \rightarrow \text{obs}}(x)$  describes the number of nucleotides that was mapped to a reference base  $B_{\text{ref}}$   
but observed to be  $B_{\text{obs}}$ , where  $B \in \{A, C, G, T\}$ . The number of C to T transitions, e.g., is denoted  
168 as  $M_{C \rightarrow T}(x)$ .

When calculating the mismatch matrix, two different approaches can be taken. Either all align-  
170 ments of the read will be counted, which we will refer to as weight-type 0, or the counts will be  
normalized by the number of alignments of each read; weight-type 1 (default).

## 172 2.2 | Damage Estimation

The damage pattern observed in aDNA has several features which are well characterized. By mod-  
174 elling these, one can construct observables sensitive to aDNA signal. We model the damage pat-  
terns seen in ancient DNA by looking exclusively at the  $C \rightarrow T$  transitions in the forward direction

<sup>176</sup> (5') and the  $G \rightarrow A$  transitions in the reverse direction (3'). For each LCA, we denote the number of transitions  $k(x)$  as:

$$\begin{aligned} \text{178} \quad k(x) = & \begin{cases} M_{C \rightarrow T}(x) & \text{for } x > 0 \quad (\text{forward}) \\ M_{G \rightarrow A}(x) & \text{for } x < 0 \quad (\text{reverse}) \end{cases} \end{aligned} \quad (1)$$

<sup>180</sup> and the number of the reference counts  $N(x)$ :

$$\begin{aligned} \text{182} \quad N(x) = & \begin{cases} \sum_{i \in B} M_{C \rightarrow i}(x) & \text{for } x > 0 \quad (\text{forward}) \\ \sum_{i \in B} M_{G \rightarrow i}(x) & \text{for } x < 0 \quad (\text{reverse}), \end{cases} \end{aligned} \quad (2)$$

<sup>184</sup> where the sum is over all four bases. The damage frequency is thus  $f(x) = k(x)/N(x)$ . A natural choice of likelihood model would be the binomial distribution. However, we found that a binomial likelihood lacks the flexibility needed to deal with the large amount of variance (overdispersion) <sup>186</sup> we found in the data due to bad references and misalignments.

To accommodate overdispersion, we instead apply a beta-binomial distribution,  $\mathcal{P}_{\text{BetaBinomial}}$ , which <sup>188</sup> treats the probability,  $p$ , as a random variable following a beta distribution<sup>1</sup> with mean  $\mu$  and concentration  $\phi$ :  $p \sim \text{Beta}(\mu, \phi)$ . The beta-binomial distribution has the the following probability density function:

$$\mathcal{P}_{\text{BetaBinomial}}(k | N, \mu, \phi) = \binom{N}{k} \frac{B(k + \mu\phi, N - k + \phi(1 - \mu))}{B(\mu\phi, \phi(1 - \mu))}, \quad (3)$$

where  $B$  is defined as the beta function:

$$B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}, \quad (4)$$

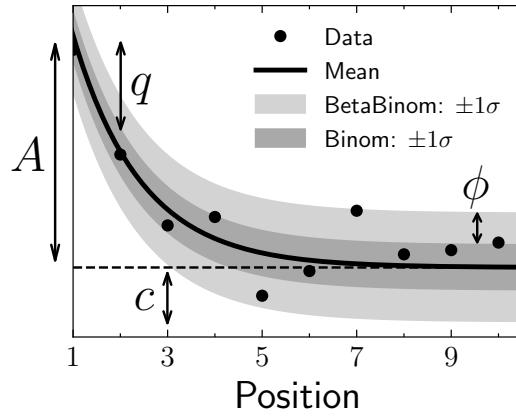
<sup>1</sup> Note that we do not parameterize the beta distribution in terms of the common  $(\alpha, \beta)$  parameterization, but instead using the more intuitive  $(\mu, \phi)$  parameterization. One can re-parameterize  $(\alpha, \beta) \rightarrow (\mu, \phi)$  using the following two equations:  $\mu = \frac{\alpha}{\alpha+\beta}$  and  $\phi = \alpha + \beta$  (Cepeda-Cuervo and Cifuentes-Amado, 2017). <sup>192</sup> <sup>194</sup> with  $\Gamma(\cdot)$  being the gamma function (Cepeda-Cuervo and Cifuentes-Amado, 2017).

The close resemblance to a binomial model is most easily seen by comparing the mean and variance of a random variable  $k$  following a beta-binomial distribution,  $k \sim \mathcal{P}_{\text{BetaBinomial}}(N, \mu, \phi)$ :

$$\begin{aligned} \mathbb{E}[k] &= N\mu \\ \mathbb{V}[k] &= N\mu(1 - \mu) \frac{\phi + N}{\phi + 1}. \end{aligned} \quad (5)$$

<sup>200</sup> The expected value of  $k$  is similar to that of a binomial distribution and the variance of the beta-binomial distribution reduces to a binomial distribution as  $\phi \rightarrow \infty$ . The beta-binomial distribution <sup>202</sup> can thus be seen as a generalization of the binomial distribution.

Note that both equation (3) and (5) relates to damage at a specific base position, i.e. for a single <sup>204</sup>  $k$  and  $N$ . To estimate the overall damage in the entire read using the position dependent counts,



**Figure 1.** Illustration of the damage model. The figure shows data points as circles and the damage,  $f(x)$ , as a solid line. The amplitude of the damage is  $A$ , the offset is  $c$ , and the relative decrease in damage pr. position is given by  $q$ . The damage uncertainty for a binomial model is shown in dark grey and the uncertainty for a beta-binomial model in light grey. The additional uncertainty of the beta-binomial model, compared to the binomial model, is related to  $\phi$ , see equation (5).

*k(x)* and  $N(x)$ , we model  $\mu$  as position dependent,  $\mu(x)$ , and assume a position-independent concentration,  $\phi$ . We model the damage frequency with a modified geometric sequence, i.e. exponential decreasing for discrete values of  $x$ :

$$208 \quad \tilde{f}(x; A, q, c) = A(1 - q)^{|x|-1} + c. \quad (6)$$

*210* Here  $A$  is the amplitude of the damage and  $q$  is the relative decrease of damage pr. position. A  
*211* background,  $c$ , was added to reflect the fact that the mismatch between the read and reference  
*212* might be due to other factors than just ancient damage. As such, we allow for a non-zero amount  
*213* of damage, even as  $x \rightarrow \infty$ . This is visualized in Fig. 1 along with a comparison between the classical  
*214* binomial model and the beta-binomial model.

To estimate the fit parameters,  $A$ ,  $q$ ,  $c$ , and  $\phi$ , we apply Bayesian inference to utilize domain  
<sup>2</sup> Parameterized as  $(\mu, \phi)$   
*216* specific knowledge in the form of priors. We assume weakly informative beta-priors<sup>2</sup> for both  $A$ ,  $q$ ,  
*217* and  $c$ . In addition to this, we assume an exponential prior on  $\phi$  with the requirement of  $\phi > 2$  to

218 avoid too much focus on 0-or-1 probabilities (McElreath, 2020). The final model is thus:

$$\begin{aligned} \text{[A prior]} & A \sim \text{Beta}(0.1, 10) \\ \text{[q prior]} & q \sim \text{Beta}(0.2, 5) \\ \text{[c prior]} & c \sim \text{Beta}(0.1, 10) \\ \text{[phi prior]} & \phi \sim 2 + \text{Exponential}(1000) \\ \text{[likelihood]} & k_i \sim \mathcal{P}_{\text{BetaBinomial}}(N_i, \tilde{f}(x_i; A, q, c), \phi), \end{aligned} \tag{7}$$

where  $i$  is an index running over all positions.

226 We define the damage due to deamination,  $D$ , as the background-subtracted damage frequency at the first position:  $D \equiv \tilde{f}(|x| = 1) - c$ . As such,  $D$  is the damage related to ancientness. Using the 228 properties of the beta-binomial distribution, eq. (5), we find the mean and variance of the damage:

$$\begin{aligned} \mathbb{E}[D] & \equiv \bar{D} = A \\ \mathbb{V}[D] & \equiv \sigma_D^2 = \frac{A(1-A)}{N} \frac{\phi + N}{\phi + 1}. \end{aligned} \tag{8}$$

230 Since  $D$  estimates the overexpression of damage due to ancientness, not only the mean is relevant but also the certainty of  $D > 0$ . We quantify this through the significance  $Z = \bar{D}/\sigma_D$  232 which is thus the number of standard deviations ("sigmas") away from zero. Assuming a Gaussian distribution of  $D$ ,  $Z > 2$  would indicate a probability of  $D$  being larger than zero, i.e. containing 234 ancient damage, with more than 97.7% probability. These two values allows us to not only quantify the amount of ancient damage (ie.  $\bar{D}$ ) but also the certainty of this damage ( $Z$ ) without even having 236 to run multiple models and comparing these.

We perform the Bayesian inference of the parameters models using Hamiltonian Monte Carlo 238 (HMC) sampling which is a particular of Monte Carlo Markov Chain (MCMC) algorithm (Betancourt, 2018). Specifically, we use the NUTS implementation in NumPyro (Phan, Pradhan, and Jankowiak, 240 2019), a Python package which uses JAX (Bradbury et al., 2018) under the hood for automatic differentiation and JIT compilation. We treat each leaf node of the LCA as being independent and 242 generate 1000 MCMC samples after an initial 500 samples as warm up.

Since running the full Bayesian model is computationally expensive, we also allow for a faster, 244 approximate method by just fitting the maximum a posteriori probability (MAP) estimate. We use iMinuit (Dembinski et al., 2021) for the MAP optimization with Numba acceleration (Lam, Pitrou, 246 and Seibert, 2015) for even faster run times. On a Macbook M1 Pro model from 2021, the timings

for running the full Bayesian model is  $1.41 \pm 0.04$  s/fit and for the MAP it is  $4.34 \pm 0.07$  ms/fit, showing  
more than a 2 order increase in performance (around 300x) for the approximate model. Both  
models allow for easy parallelisation to decrease the computation time.

### 250    2.3 | Visualisation

We provide an interactive dashboard to properly visualise the results from the modelling phase,  
252    see <https://metadmg.onrender.com/> for an example. The dashboard allows for filtering, styling and  
variable selection, visualizing the mismatch matrix related to a specific leaf node, and exporting of  
254    both fit results and plots. By filtering, we include both filtering by sample, by specific cuts in the fit  
results (e.g. requiring  $D$  to be above a certain threshold), and even by taxonomic level (e.g. only  
256    looking tax IDs that are part of the Mammalia class). We greatly believe that a visual overview of  
the fit results increase understanding of the data at hand. The dashboard is implemented with  
258    Plotly plots and incorporated into a Dash dashboard (Plotly, 2015).

## 3 | SIMULATION STUDY

260    To ascertain the performance of our test statistic and implementation we performed various rig-  
orous simulation studies to quantify possible issues with bias and accuracy in a synthetic setting  
262    that should mimick the various issues and complications that exist with real world data. We con-  
ducted two sets of simulations, one to gauge the performance of the damage model itself and one  
264    to determine the performance of the full metaDMG pipeline, i.e. both LCA and damage model.

### 3.1 | Single-genome Simulations

266    The first set of simulations was performed by taking a single, representative genome and adding  
post mortem damage together with sequencing noise. This was followed by a standard mapping  
268    step and finally damage estimation using metaDMG. The deamination was applied using NGSNGS  
(Henriksen, Zhao, and T. Korneliussen, 2022) which is a recent implementation of the original  
270    Briggs model similar to Gargammel (Neukamm, Peltzer, and Nieselt, 2021). In this step we vary  
the simulated amount of damage added (in particular the single-stranded DNA deamination,  $\delta_{ss}$ )  
272    in the original Briggs model (Briggs et al., 2007)), the number of reads, and the fragment length  
distribution.

274       ./ngsngs [something]

```
bowtie2 [something]
```

276 We chose five different, representative genomes, in each of these varying the three simulation  
parameters. These genomes where the homo sapiens, the betula, and three microbial organisms  
278 with respectively low, median, and high amount of GC-content. For each of these simulations, we  
performed 100 independent replicates to measure the variability of the parameter estimation and  
280 quantify the robustness of the estimates. We simulated eight different sets of damage (approximately  
0%, 1%, 2%, 5%, 10%, 15%, 20%, 30%), 13 sets of different number of reads (10, 25, 50, 100, 250,  
282 500, 1.000, 2.500, 5.000, 10.000, 25.000, 50.000, 100.000), three sets of different fragment length distri-  
butions (samples from a *log-normal* distribution with mean 35, 60, and 90, each with a standard  
284 deviation of 10), and five different genomes, each simulation set repeated 100 times.

In addition to this, we also create 1000 repetitions of the non-damaged simulations for Homo  
286 Sapiens to be able to gauge the risk of finding false positives. Finally, to show that the damage esti-  
mates that metaDMG provides are independent of the contig size, we artificially create three different  
288 genomes by sampling 1.000, 10.000 or 100.000 different basepairs from a uniform categorical dis-  
tribution of {A, C, G, T}.

290 To be able to compare our estimates to a known value, we generate 1.000.000 reads using  
NGSNGS without any added sequencing noise for each of other sets of simulation parameters.  
292 The difference in damage frequency at position 1 and 15 is then the value to compare to:

$$D_{\text{known}} = \frac{f(x=1) - f(x=15)}{2} + \frac{f(x=-1) - f(x=-15)}{2}, \quad (9)$$

294 where we take the average of the C to T damage frequency difference and the G to A damage  
frequency difference.

296 The fastq files were simulated with NGSNGS using the above mentioned simulation parameters,  
all with the same quality scores profiles as used in ART (Huang et al., 2012), based on the Illumina  
298 HiSeq 2500 (150 bp). The mapping was performed using Bowtie-2 with the -no-unal flag (Langmead  
and Salzberg, 2012).

### 300 3.2 | Metagenomic Simulations

While the previously mentioned simulation study is perfectly aimed at quantifying the performance  
302 of the damage model in the case of single-reference genomics it does lack the complexity related  
to metagenomic samples. Therefore, we also conduct a more advanced simulation study to deter-

304 mine the accuracy of the full `metaDMG` pipeline.

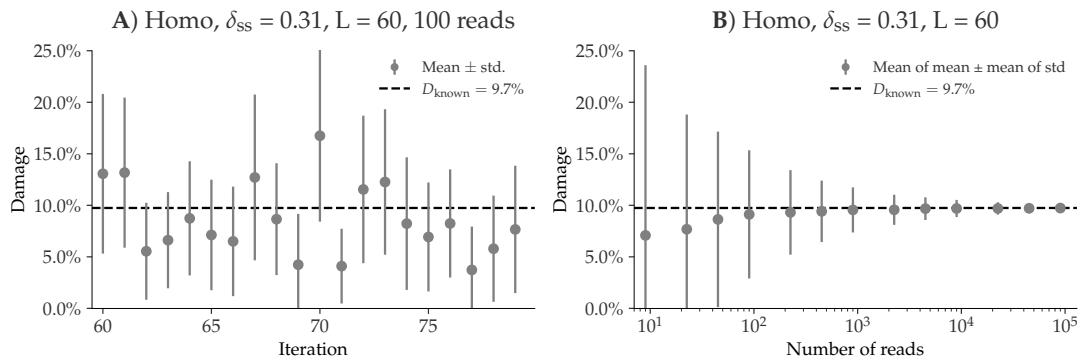
The previously mentioned simulation study quantifies the damage model's performance for  
306 single-reference genomics, but it does not address the complexity of metagenomic studies. Therefore,  
308 we also conducted a more advanced simulation study to determine the performance of the  
`metaDMG` pipeline in a standard eDNA setting. Based on an ancient metagenome scenario, we created  
310 a synthetic dataset that mimics the composition, fragment length distribution, and damage  
patterns for each genome. We selected X metagenomes (Supp table XXX) covering several environmental  
conditions and ages. First, we mapped the reads of each metagenome with `bowtie2` against  
312 a database that contained the GTDB r202 (Parks et al., 2018) species cluster reference sequences,  
all organelles from NCBI RefSeq (NCBI Resource Coordinators, 2018), and the reference sequences  
314 from CheckV (Nayfach et al., 2021). We used `bam-filter 1.0.11` with the flag `--read-length-freqs` to  
get the mapped read length distribution for each genome and their abundance. The genomes with  
316 an observed-to-expected coverage ratio greater than 0.75 were kept. The filtered BAM files were  
processed by `metaDMG` to obtain the misincorporation matrices. The abundance tables, fragment  
318 length distribution, and misincorporation matrices were used in aMGSIM-smk v0.0.1 (Fernandez-  
Guerra, 2022), a Snakemake workflow (Mölder et al., 2021) that facilitates the generation of many  
320 synthetic ancient metagenomes. The data used and generated by the workflow can be obtained  
from Figshare link (XXX). We then performed taxonomic profiling using the same parameters used  
322 for the synthetic reads generated by aMGSIM-smk.

## 4 | RESULTS

324 The accuracy of all methods in `metaDMG` was tested in various simulation scenarios. In general we  
find that `metaDMG` yields accurate, precise damage estimates even in extreme low-coverage data.

### 4.1 | Single-genome Simulations

The results of the single-genome simulations can be seen in Figure 2. The left part of the figure  
328 shows `metaDMG` damage estimates based on the *homo sapiens* genome with the Briggs parameter  
 $\delta_{SS} = 0.31$  and a fragment length distribution with mean 60, each of the simulations generated with  
330 100 simulated reads for 10 representative simulations. When the damage estimates are low, the  
distribution of  $D$  is highly skewed (restricted to positive values) leading to errorbars sometimes  
332 going into negative damage, which of course represents un-physical values. The right hand side of



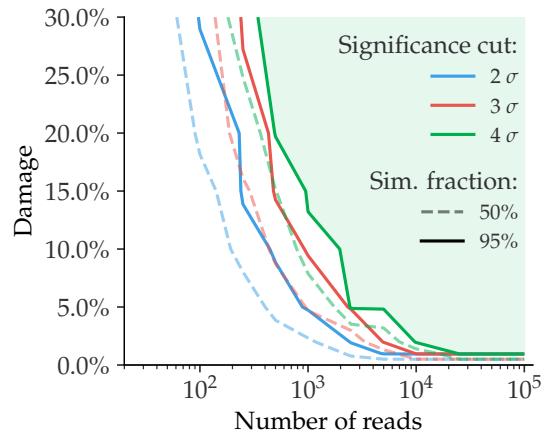
**Figure 2.** Overview of the single-genome simulations based on the homo sapiens genome with the Briggs parameter  $\delta_{SS} = 0.065$  and a fragment length distribution with mean 60. **A)** This plot shows the estimated damage ( $D$ ) of 10 simulations with 100 simulated reads. The grey points show the mean damage (with its standard deviation as errorbars). The known damage ( $D_{known}$ ) is shown as a dashed line, see eq. (9). **B)** This plot shows the average damage as a function of the number of reads. The grey points show the average of the individual means (with the average of the standard deviations as errors).

the figure visualizes the average amount of damage across a varying number of reads. This shows  
 334 that the damage estimates converge to the known value with more data, and that one needs more than 100 reads to even get strictly positive damage estimates (when including uncertainties).

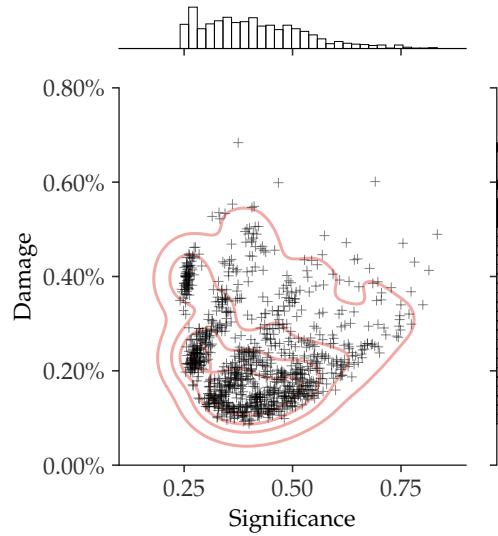
336 Across more than 5 different species, 3 different fragment length distributions, and 3 different contig length distributions, each with 100 simulations for 104 different sets of simulation parameters, the only difference we note in the damage estimates is between species with low, median, and high GC-levels. In general, species with higher GC-levels exhibit lower variations in their damage  
 340 estimates compared to species with lower GC-levels, leading to high-GC species requiring fewer reads to establish damage estimates.

342 Based on the single-genome simulations, we can compute the relationship between the amount of damage in a species and the number of reads required to correctly infer that the given species  
 344 is damaged, see Figure 3. If we want to find damage with a significance of more than 2 (solid blue line) in a sample with around 5% expected damage, it requires about 1000 reads to be 95% certain  
 346 that we will find results this good. Said in other words: given 100 different fits, each with 1000 reads and around 5% damage, one would expect to find damage (with a  $Z > 2$ ) in 95 of the total  
 348 100 samples, on average. If we lose the requirement such that it is okay to only find it in every second fit, it would be enough with only around 250 reads in each fit (dashed blue line).

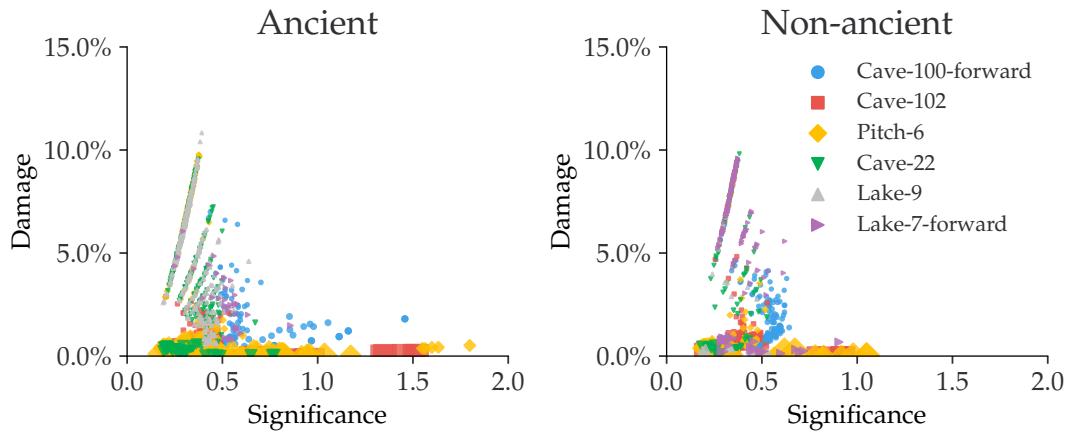
350 Finally, to quantify the risk of incorrectly assigning damage to a non-damaged species, we cre-



**Figure 3.** Given a specific significance cut, the solid contour line shows the relationship between the amount of damage and the number of reads required to be able to correctly infer damage in 95% of the species. The dashed line shows the similar value for a simulation fraction of 50%. The green part of the figure shows the “good” region of number of reads and estimated damage, given than one wants to be more than 95% certain of correctly identifying damage with more than  $4\sigma$  confidence.



**Figure 4.** This figure shows the inferred damage estimates of 1000 independent simulations, each with 1000 reads and no artificial ancient damage applied, with the inferred damage shown on the y-axis and the significance on the x-axis. Each simulation is shown as a single cross and the red lines show the kernel density estimate (KDE) of the damage estimates. The marginal distributions are shown as histograms next to the scatter plot.

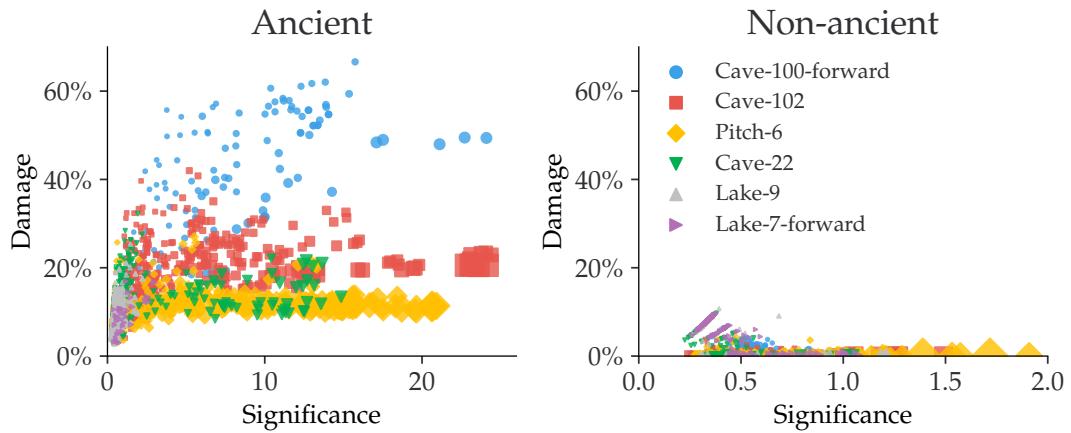


**Figure 5.** Estimated amount of damage as a function of significance using the fragSim data. The left figure shows the damage of the species that we simulated to be ancient (however with no deamination added yet) and the right figure shows the same for the species that are not going to have deamination added.

ated 1000 independent simulations for a varying number of reads, where none of them had any ar-  
 352 tificial ancient damage applied, only sequencing noise. Figure 4 shows the damage ( $D$ ) as a function  
 of the significance ( $Z$ ) for the case of 1000 simulated reads. Even though the estimated damage is  
 354 larger than zero, the damage is non-significant since the significance is less than one. When looking  
 at all the figures across the different number of reads, see Figure bayesian\_zero\_damage\_plots.pdf,  
 356 we note that a loose cut requiring that  $D > 1\%$  and  $Z > 2$  would filter out all of non-damaged  
 points. Overall the conclusion being that our devised test statistic is conservative and has low  
 358 false positive rate.

## 4.2 | Metagenomic Simulations

360 With the full metagenomic simulation pipeline we can further probe the performance of metaDMG.  
 By looking at the six different metagenomic scenarios at different steps in the pipeline we are able  
 362 to show that metaDMG provides relevant, accurate damage estimates. First of all, we run metaDMG on  
 the six samples after fragmentation with FragSim. Since no deamination has yet been added at  
 364 this step in the pipeline, this is also a test of the risk of getting false positives. The results can be  
 seen in Figure 5 where we see the damage estimates for both the species that we simulate to be  
 366 ancient and the species that we do not add deamination to. We see that the damage estimates are  
 quite similar, as expected, and that our previously established loose cut of  $D > 1\%$  and  $Z > 2$  still  
 368 filters out all of non-damaged points.



**Figure 6.** Estimated amount of damage as a function of significance using the ART data. The left figure shows the damage of the species that we simulated to be ancient and the right figure shows the same for the species that have not had deamination added.

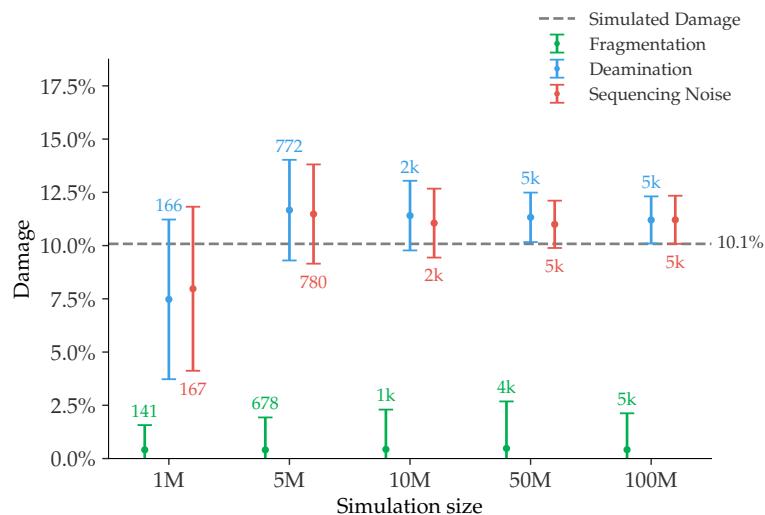
In comparison we can look at Figure 6 which shows the same plot, but after the deamination 370 (deamSim) and sequencing errors (ART) has been added. Here we see a clear difference between the ancient and the non-ancient ones, as expected. The non-ancient species would still not pass 372 the loose cut, however, we note that a large number of the ancient samples would. By looking at Figure 6 we see that not all of the samples show similar amount of damage. These observations 374 are summarised in Table 1 where we see that Cave-100-forward, Cave-102, Pitch-6 all have more than 60% of their ancient species labelled as damaged according to the loose cut, Cave-22 (18%) 376 and Lake-7-forward (12%) a bit lower, while Lake-9 (0.5%) does not show any clear signs of damage. However, once we condition on the requirement of having more than 100 reads, the fraction of 378 ancient species correctly identified as ancient increases to more than 90% for most the samples.

To better understand the damage estimates, we can look a them individually. Figure 7 shows 380 the *Stenotrophomonas maltophilia* species from the Pitch-6 sample. We see that none of the fragmentation-only files were estimated to have damage and that most of the deamination and 382 final files including sequencing errors have damage – at a simulation size of 1 million, the significance of both are  $Z \approx 1.9$ , so this one of the few fits with more than 100 reads that does not 384 pass the loose cut. Furthermore, we notice that the error bars decrease with simulation size, as expected.

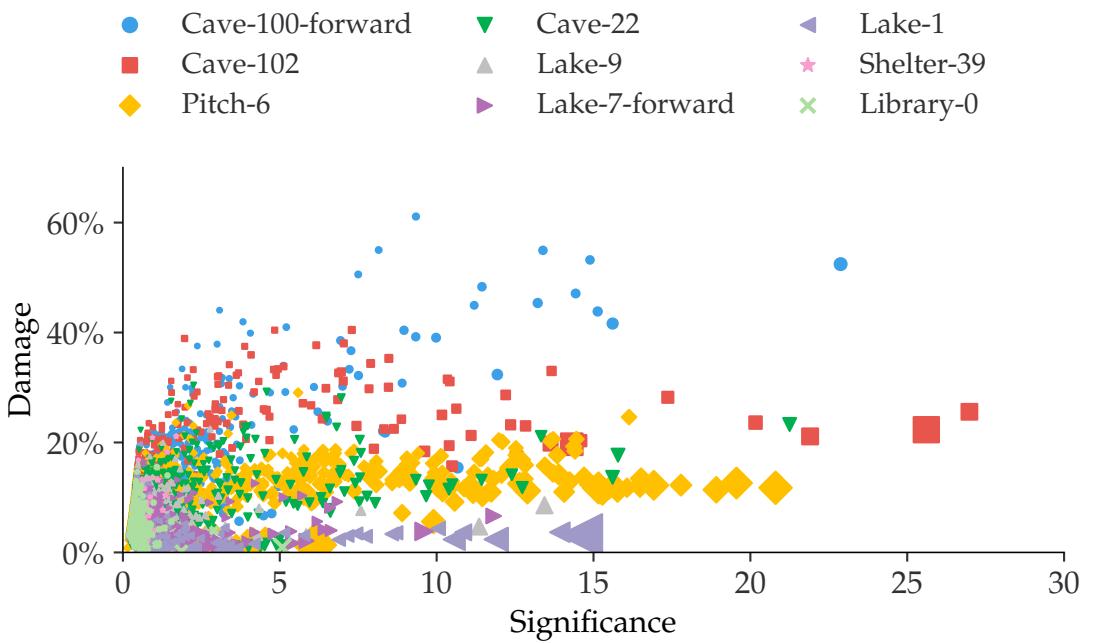
386 The rest of the metagenomic simulation results are shown in Figure XXX.

**Table 1.** Number of ancient species for each of the six simulated samples. The first column is the total number of species, the second column is the total number of species that would pass the loose cut of  $D > 1\%$  and  $Z > 2$ , the third column is the number of species with more than 100 reads, and the final column is the number of species with more than 100 reads that also do pass the cut.

Sample	Total	Pass	+100 Reads	+100 Reads and Pass		
Cave-100-forward	135	107	79.3%	88	87	98.9%
Cave-102	500	326	65.2%	309	285	92.2%
Pitch-6	415	260	62.7%	274	260	94.9%
Cave-22	393	71	18.1%	73	69	94.5%
Lake-9	410	2	0.5%	8	0	0%
Lake-7-forward	32	4	12.5%	6	4	66.7%



**Figure 7.** Damage estimates of the *Stenotrophomonas maltophilia* species from the Pitch-6 sample. Damage is shown as a function of simulation size, with the fragmentation files in green, the deamination files in blue and the final files including sequencing errors in red. All errors are  $1\sigma$  error bars (standard deviation). The number of reads for each fit is shown as text and since this was a species simulated to have ancient damage, the simulated amount of damage is shown as a dashed grey line.



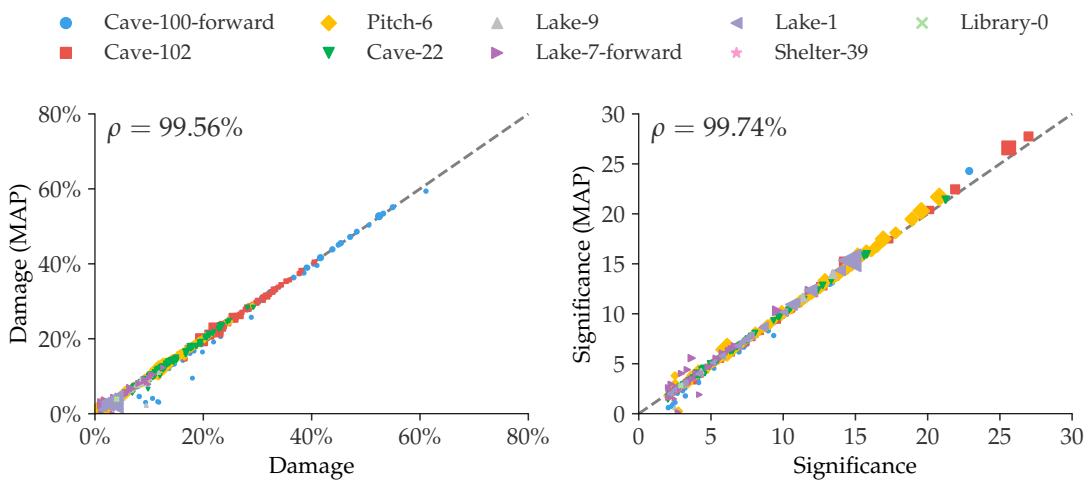
**Figure 8.** Estimated amount of damage as a function of significance using the real data.

### 4.3 | Real Data

388 The results from running the full metaDMG pipeline on real data can be seen in Figure 8. The figures  
 shows Blablabla, real life data here. We find that the loose cut ( $D > 1\%$ ,  $Z > 2$ ) accepts only one of  
 390 the fits from the control test Library-0, which would not have been accepted by more conservative  
 cut ( $D > 2\%$ ,  $Z > 3$ , more than 100 reads).

### 392 4.4 | Bayesian vs. MAP

Due to increased computational burden of running the full Bayesian model compared to faster,  
 394 approximate MAP model, in samples with several thousand species, the MAP model is often the  
 most realistic model to use due to time constraints. In this case, it is of course important to know  
 396 that the damage estimates are indeed trustworthy. Figure 9 compares the estimated damage  
 between the Bayesian model and the MAP model and the estimated significances for species with  
 398  $D > 1\%$ ,  $Z > 2$  and more than 100 reads. The figure shows that the vast majority of species map  
 400 1:1 between the Bayesian and the MAP model. One should note, though, that the few species with  
 the highest mismatch, all are based on forward-only fits, i.e. with no information from the reverse  
 strand, which thus leads to less data to base the fits on. For the comparison with no cuts, see  
 402 Figure 1 in appendix.



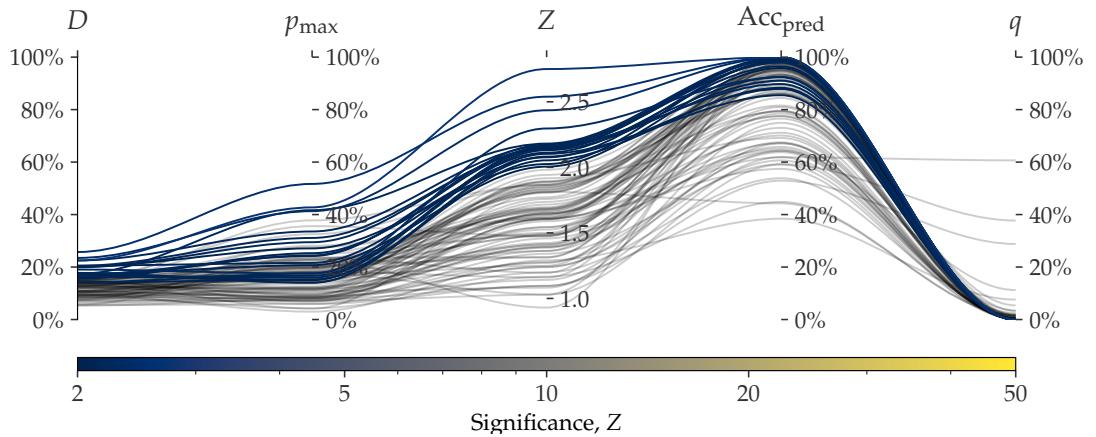
**Figure 9.** Comparison between the full Bayesian model and the fast, approximate, MAP model for the estimated damage and significance. The figure shows data after a loose cut of  $D > 1\%$ ,  $Z > 2$  and more than 100 reads. The dashed, grey line shows the 1:1 ratio and the correlation,  $\rho$ , is shown in the upper right corner.

## 4.5 | Existing Methods

404 We have also compared `metaDMG` to existing methods such as PyDamage (Borry et al., 2021). Since  
 405 PyDamage does not include the LCA step, this comparison is based on the non-LCA mode (local-  
 406 mode) of `metaDMG`. This mode iterates through the different assigned species for all mapped reads  
 407 and estimates the damage for each. In general, we find that `metaDMG` is more conservative, accurate  
 408 and precise in its damage estimates.

On example of this is can be found in Figure 10, which shows both the `metaDMG` and PyDamage  
 410 results of the 100 Homo Sapiens single-genome simulations with 100 reads and 15% added artificial  
 411 damage (and a fragment length distribution with mean 60).

412 To compare the computational performance, we use the Pitch-6 sample which has 11.433  
 413 unique taxa. When using only a single core, PyDamage took 1105 s to compute all fits, while `metaDMG`  
 414 took 88 s, a factor of 12.6x faster. Out of the 88 s, `metaDMG` spent 53 s on the actual fits, the rest was for  
 415 loading and reading the alignment file and computing the mismatch matrices. This makes `metaDMG`  
 416 more than 20x faster than PyDamage for the fit computation. For the rest of the timings, see Ta-  
 417 ble 2. PyDamage requires the alignment file to be sorted by chromosome position and be supplied  
 418 with an index file, allowing it to iterate fast through the alignment file, at the expense of compu-  
 419 tational load before running the actual damage estimation. `metaDMG` on the other hand requires the  
 420 reads to be sorted by name to minimize the time it takes to run the LCA, which however, is not



**Figure 10.** Parallel Coordinates plot comparing metaDMG and PyDamage for the *Homo Sapiens* single-genome simulation with 100 reads and 15% added artificial damage. The different axis shows the five different variables: metaDMG-damage ( $D$ , by metaDMG), PyDamage-damage ( $p_{\max}$ , by PyDamage), significance ( $Z$ , by metaDMG), predicted accuracy ( $\text{Acc}_{\text{pred}}$ , by PyDamage), and the p-value ( $q$ , by PyDamage). Each of the 100 simulations are plotted as single lines showing the values of the different dimensions. Simulations with  $D > 1\%$  and  $Z > 2$ , i.e. damaged according to the loose metaDMG cut, are shown in color proportional to their significance. Non-damaged simulations are shown in semi-transparent black lines.

tested in this comparison.

## 5 | DISCUSSION

Preliminary work indicates that the computational performance of the models can be even further optimized by using Julia (Bezanson et al., 2017), which shows around 7x optimization for the Bayesian model (~ 0.2 s/fit) and 4x for the MAP model (~ 1.1 ms/fit).

- 426     • contig length
- simulation setup
- 428     • more covariates, readlength? other substitutions
- no linkage
- 430     • taxa independence,
- weight
- 432     • improved fishing (pmdtools)
- relevant testcases?
- 434     • basis for environmental deme differences?

**Table 2.** Computational performance of PyDamage and metaDMG. The table contains the times it takes to run either PyDamage or metaDMG on the full Pitch-6 sample containing 11.433 species. The timings are shown for both single-processing case (1 core) and multi-processing (2 and 4 cores). The timings were performed on a Macbook M1 Pro model from 2021. “12.6x” means that metaDMG was 12.6 times faster than PyDamage for that particular test.

Cores	Pitch-6		Pydamage		metaDMG	
	Total	Fits	Total	Fits		
1	1105 s	1102 s	88 s	12.6x	53 s	20.8x
2	592 s	590 s	66 s	9.0x	25 s	23.6x
4	398 s	397 s	54 s	7.4x	14s	28.4x

## 5.1 | Acknowledgment

<sup>436</sup> Acknowledgements here

## 5.2 | Data availability

<sup>438</sup> Source code is hosted at GitHub: <https://github.com/metaDMG-dev>. Sequencing data can be found at: <https://somewhere.com> XXX.

## 440 REFERENCES

- Betancourt, Michael (2018). "A Conceptual Introduction to Hamiltonian Monte Carlo". In: *arXiv:1701.02434* 442 [*stat*]. arXiv: 1701.02434.
- Bezanson, Jeff et al. (2017). "Julia: A fresh approach to numerical computing". In: *SIAM review* 59.1. 444 Publisher: SIAM, pp. 65–98. URL: <https://julialang.org/>.
- Borry, Maxime et al. (2021). "PyDamage: automated ancient damage identification and estimation 446 for contigs in ancient DNA de novo assembly". en. In: *PeerJ* 9. Publisher: PeerJ Inc., e11845. ISSN: 2167-8359. DOI: [10.7717/peerj.11845](https://doi.org/10.7717/peerj.11845). URL: <https://peerj.com/articles/11845> (visited on 2022).
- 448 Bradbury, James et al. (2018). *JAX: composable transformations of Python NumPy programs*. URL: <http://github.com/google/jax>.
- 450 Briggs, Adrian W. et al. (2007). "Patterns of damage in genomic DNA sequences from a Neandertal". en. In: *Proceedings of the National Academy of Sciences* 104.37. Publisher: National Academy of 452 Sciences Section: Biological Sciences, pp. 14616–14621. ISSN: 0027-8424, 1091-6490. DOI: [10.1073/pnas.0704665104](https://doi.org/10.1073/pnas.0704665104). URL: <https://www.pnas.org/content/104/37/14616>.
- 454 Cabanski, Christopher R. et al. (2012). "ReQON: a Bioconductor package for recalibrating quality 456 scores from next-generation sequencing data". In: *BMC Bioinformatics* 13.1, p. 221. ISSN: 1471-2105. DOI: [10.1186/1471-2105-13-221](https://doi.org/10.1186/1471-2105-13-221). URL: <https://doi.org/10.1186/1471-2105-13-221> (visited on 2022).
- 458 Cappellini, Enrico et al. (2018). "Ancient Biomolecules and Evolutionary Inference". In: *Annual Review 460 of Biochemistry* 87.1. \_eprint: <https://doi.org/10.1146/annurev-biochem-062917-012002>, pp. 1029-1060. DOI: [10.1146/annurev-biochem-062917-012002](https://doi.org/10.1146/annurev-biochem-062917-012002). URL: <https://doi.org/10.1146/annurev-biochem-062917-012002> (visited on 2022).
- 462 Cepeda-Cuervo, Edilberto and MARÍA VICTORIA Cifuentes-Amado (2017). "Double Generalized Beta- 464 Binomial and Negative Binomial Regression Models". en. In: *Revista Colombiana de Estadística* 40.1. Publisher: Universidad Nacional de Colombia., pp. 141–163. ISSN: 0120-1751. DOI: [10.15446/rce.v40n1.61779](https://doi.org/10.15446/rce.v40n1.61779).
- 466 Champlot, Sophie et al. (2010). "An efficient multistrategy DNA decontamination procedure of PCR 468 reagents for hypersensitive PCR applications". eng. In: *PloS One* 5.9, e13042. ISSN: 1932-6203. DOI: [10.1371/journal.pone.0013042](https://doi.org/10.1371/journal.pone.0013042).

- Dabney, Jesse, Matthias Meyer, and Svante Pääbo (2013). "Ancient DNA Damage". In: *Cold Spring Harbor Perspectives in Biology* 5.7, a012567. ISSN: 1943-0264. DOI: [10.1101/cshperspect.a012567](https://doi.org/10.1101/cshperspect.a012567).  
470 URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3685887/> (visited on 2022).
- 472 Dembinski, Hans et al. (2021). *scikit-hep/iminuit*: v2.8.2. DOI: [10.5281/ZENODO.3949207](https://doi.org/10.5281/ZENODO.3949207). (Visited on 2021).
- 474 Fernandez-Guerra, Antonio (2022). *genomewalker/aMGSIM-smk*: v0.0.1. DOI: [10.5281/zenodo.7298422](https://doi.org/10.5281/zenodo.7298422).  
URL: <https://doi.org/10.5281/zenodo.7298422>.
- 476 Gilbert, M. Thomas P. et al. (2005). "Assessing ancient DNA studies". en. In: *Trends in Ecology & Evolution* 20.10, pp. 541–544. ISSN: 0169-5347. DOI: [10.1016/j.tree.2005.07.005](https://doi.org/10.1016/j.tree.2005.07.005). URL: <https://www.sciencedirect.com/science/article/pii/S0169534705002260> (visited on 2022).
- 478 Ginolhac, Aurelien et al. (2011). "mapDamage: testing for damage patterns in ancient DNA sequences". In: *Bioinformatics* 27.15, pp. 2153–2155. ISSN: 1367-4803. DOI: [10.1093/bioinformatics/btr347](https://doi.org/10.1093/bioinformatics/btr347). URL: <https://doi.org/10.1093/bioinformatics/btr347> (visited on 2022).
- 480 Henriksen, Ramus, Lei Zhao, and Thorfinn Korneliussen (2022). *NGSNGS*: v0.5.0. DOI: [10.5281/zenodo.7326212](https://doi.org/10.5281/zenodo.7326212). URL: <https://doi.org/10.5281/zenodo.7326212>.
- 482 484 Huang, Weichun et al. (2012). "ART: a next-generation sequencing read simulator". eng. In: *Bioinformatics (Oxford, England)* 28.4, pp. 593–594. ISSN: 1367-4811. DOI: [10.1093/bioinformatics/btr708](https://doi.org/10.1093/bioinformatics/btr708).
- 486 Jónsson, Hákon et al. (2013). "mapDamage2.0: fast approximate Bayesian estimates of ancient DNA damage parameters". en. In: *Bioinformatics* 29.13, pp. 1682–1684. ISSN: 1367-4803, 1460-2059.  
DOI: [10.1093/bioinformatics/btt193](https://doi.org/10.1093/bioinformatics/btt193).
- 488 Lam, Siu Kwan, Antoine Pitrou, and Stanley Seibert (2015). "Numba: a LLVM-based Python JIT compiler". In: *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC*. LLVM '15. New York, NY, USA: Association for Computing Machinery, pp. 1–6. ISBN: 978-1-4503-4005-2.  
DOI: [10.1145/2833157.2833162](https://doi.org/10.1145/2833157.2833162). URL: <https://github.com/numba/numba>.
- 490 494 Langmead, Ben and Steven L. Salzberg (2012). "Fast gapped-read alignment with Bowtie 2". en. In: *Nature Methods* 9.4. Number: 4 Publisher: Nature Publishing Group, pp. 357–359. ISSN: 1548-7105. DOI: [10.1038/nmeth.1923](https://doi.org/10.1038/nmeth.1923). URL: <https://www.nature.com/articles/nmeth.1923> (visited on 2022).
- 496 Llamas, Bastien et al. (2017). "From the field to the laboratory: Controlling DNA contamination in human ancient DNA research in the high-throughput sequencing era". en. In: *STAR: Science & Technology of Archaeological Research* 3.1, pp. 1–14. ISSN: 2054-8923. DOI: [10.1080/20548923](https://doi.org/10.1080/20548923).

- 500 2016.1258824. URL: <https://www.tandfonline.com/doi/full/10.1080/20548923.2016.1258824> (visited on 2022).
- 502 McElreath, Richard (2020). *Statistical rethinking: a Bayesian course with examples in R and Stan*. 2nd ed. CRC texts in statistical science. Boca Raton: Taylor and Francis, CRC Press. ISBN: 978-0-367-13991-9.
- 504 Mölder, Felix et al. (2021). *Sustainable data analysis with Snakemake*. en. Tech. rep. 10:33. Type: article. F1000Research. DOI: [10.12688/f1000research.29032.2](https://doi.org/10.12688/f1000research.29032.2). URL: <https://f1000research.com/articles/10-33> (visited on 2022).
- 508 Murchie, Tyler J. et al. (2021). "Collapse of the mammoth-steppe in central Yukon as revealed by ancient environmental DNA". en. In: *Nature Communications* 12.1. Number: 1 Publisher: Nature Publishing Group, p. 7120. ISSN: 2041-1723. DOI: [10.1038/s41467-021-27439-6](https://doi.org/10.1038/s41467-021-27439-6). URL: <https://www.nature.com/articles/s41467-021-27439-6> (visited on 2022).
- 512 Nayfach, Stephen et al. (2021). "CheckV assesses the quality and completeness of metagenome-assembled viral genomes". en. In: *Nature Biotechnology* 39.5. Number: 5 Publisher: Nature Publishing Group, pp. 578–585. ISSN: 1546-1696. DOI: [10.1038/s41587-020-00774-7](https://doi.org/10.1038/s41587-020-00774-7). URL: <https://www.nature.com/articles/s41587-020-00774-7> (visited on 2022).
- 516 NCBI Resource Coordinators (2018). "Database resources of the National Center for Biotechnology Information". In: *Nucleic Acids Research* 46.D1, pp. D8–D13. ISSN: 0305-1048. DOI: [10.1093/nar/gkx1095](https://doi.org/10.1093/nar/gkx1095). URL: <https://doi.org/10.1093/nar/gkx1095> (visited on 2022).
- 520 Neukamm, Judith, Alexander Peltzer, and Kay Nieselt (2021). "DamageProfiler: fast damage pattern calculation for ancient DNA". In: *Bioinformatics* 37.20, pp. 3652–3653. ISSN: 1367-4803. DOI: [10.1093/bioinformatics/btab190](https://doi.org/10.1093/bioinformatics/btab190). URL: <https://doi.org/10.1093/bioinformatics/btab190> (visited on 2022).
- 524 Parks, Donovan H. et al. (2018). "A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree of life". en. In: *Nature Biotechnology* 36.10. Number: 10 Publisher: Nature Publishing Group, pp. 996–1004. ISSN: 1546-1696. DOI: [10.1038/nbt.4229](https://doi.org/10.1038/nbt.4229). URL: <https://www.nature.com/articles/nbt.4229> (visited on 2022).
- 526 Pedersen, Mikkel W. et al. (2016). "Postglacial viability and colonization in North America's ice-free corridor". en. In: *Nature* 537.7618. Number: 7618 Publisher: Nature Publishing Group, pp. 45–49. ISSN: 1476-4687. DOI: [10.1038/nature19085](https://doi.org/10.1038/nature19085). URL: <https://www.nature.com/articles/nature19085> (visited on 2022).

- 530 Phan, Du, Neeraj Pradhan, and Martin Jankowiak (2019). "Composable Effects for Flexible and Accelerated Probabilistic Programming in NumPyro". In: *arXiv:1912.11554 [cs, stat]*. arXiv: 1912.11554.
- 532 Plotly, Technologies (2015). *Collaborative data science*. Place: Montreal, QC Publisher: Plotly Technologies Inc. URL: <https://plot.ly>.
- 534 Skoglund, Pontus et al. (2014). "Separating endogenous ancient DNA from modern day contamination in a Siberian Neandertal". In: *Proceedings of the National Academy of Sciences* 111.6. Publisher: Proceedings of the National Academy of Sciences, pp. 2229–2234. DOI: [10.1073/pnas.1318934111](https://doi.org/10.1073/pnas.1318934111). URL: <https://www.pnas.org/doi/10.1073/pnas.1318934111> (visited on 2022).
- 538 Wang, Yucheng, Thorfinn Sand Korneliussen, et al. (2022). "ngsLCA—A toolkit for fast and flexible lowest common ancestor inference and taxonomic profiling of metagenomic data". en. In: *Methods in Ecology and Evolution* n/a.n/a(). \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/2041-210X.14006>. ISSN: 2041-210X. DOI: [10.1111/2041-210X.14006](https://doi.org/10.1111/2041-210X.14006). URL: <https://onlinelibrary.wiley.com/doi/abs/10.1111/2041-210X.14006> (visited on 2022).
- 544 Wang, Yucheng, Mikkel Winther Pedersen, et al. (2021). "Late Quaternary dynamics of Arctic biota from ancient environmental genomics". en. In: *Nature* 600.7887. Number: 7887 Publisher: Nature Publishing Group, pp. 86–92. ISSN: 1476-4687. DOI: [10.1038/s41586-021-04016-x](https://doi.org/10.1038/s41586-021-04016-x). URL: <https://www.nature.com/articles/s41586-021-04016-x> (visited on 2022).
- 548 Zavala, Elena I. et al. (2021). "Pleistocene sediment DNA reveals hominin and faunal turnovers at Denisova Cave". en. In: *Nature* 595.7867. Number: 7867 Publisher: Nature Publishing Group, pp. 399–403. ISSN: 1476-4687. DOI: [10.1038/s41586-021-03675-0](https://doi.org/10.1038/s41586-021-03675-0). URL: <https://www.nature.com/articles/s41586-021-03675-0> (visited on 2022).
- 550

## Appendix 1

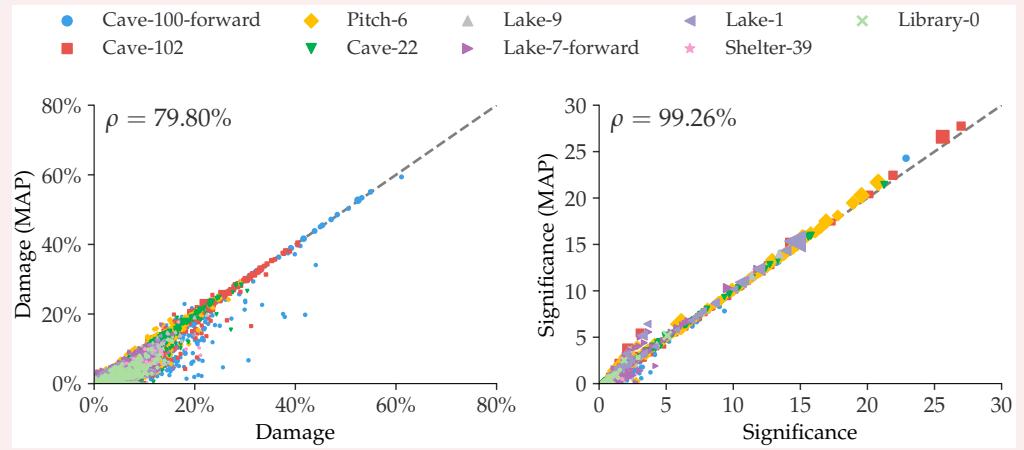
552

### A | EXAMPLE FIGURE

This is an example of including a figure in the appendix.

554

556



**Appendix 1—figure 1.** Comparison between the full Bayesian model and the fast, approximate, MAP model for the estimated damage and significance. The dashed, grey line shows the 1:1 ratio.

## B | EXAMPLE TABLE

560 This is an example of including a table in the appendix.

562 **Appendix 2—table 1.** An example table.

Speed (mph)	Driver	Car	Engine	Date
407.447	Craig Breedlove	Spirit of America	GE J47	8/5/63
413.199	Tom Green	Wingfoot Express	WE J46	10/2/64
434.22	Art Arfons	Green Monster	GE J79	10/5/64
468.719	Craig Breedlove	Spirit of America	GE J79	10/13/64
526.277	Craig Breedlove	Spirit of America	GE J79	10/15/65
536.712	Art Arfons	Green Monster	GE J79	10/27/65
555.127	Craig Breedlove	Spirit of America, Sonic 1	GE J79	11/2/65
576.553	Art Arfons	Green Monster	GE J79	11/7/65
600.601	Craig Breedlove	Spirit of America, Sonic 1	GE J79	11/15/65
622.407	Gary Gabelich	Blue Flame	Rocket	10/23/70
633.468	Richard Noble	Thrust 2	RR RG 146	10/4/83
763.035	Andy Green	Thrust SSC	RR Spey	10/15/97

## C | MULTINOMIAL LOGISTIC REGRESSIONS

### 566 Full Multinomial Logistic Regression models

Postmortem damages will have impacts on the NGS (next generation sequencing) reads. 566  
A common phenomenon is the calling error rates increases from nucleotide C to T due to 568  
the cytosine deamination process. Unawareness of this will lead to inaccurate inferences. 570  
Evidences show that the magnitude of such changes are related to the positions the site is 572  
within a read (the fraction of the ancient DNA). Here we present 3 slightly different ways to 574  
unveil the relationship between the calling error rates and the mismatching reference/read  
pairs as well as the site positions within a read. The methods are based on the multinomial  
logistic regressions.

### Data Description

We perform the regression based on the summary statistic of the mismatch matrix, i.e.,  $\underline{M}(x)$ ,  
576 which is a table which contains the counts of reads of different reference/read categories  
578 (in total 16) and positions on the forward/reversed strand (15 positions on each direction).

Table 1 and Table 2 give an example of the data format we use for the inference.

Ref.	Read Counts								
	A				C				
	Read	A	C	G	T	A	C	G	T
1	12794053	8325	28769	16073	10404	8045811	8020	2092619	
2	13480290	6812	21107	12102	9151	8260185	6531	1145605	
3	12760253	6131	18859	10327	7772	8385423	5899	914709	
4	12995572	5240	17671	8940	7880	8345892	5252	767237	
5	12930102	4601	17021	8188	8374	8474964	5161	703283	
6	12879355	4684	16435	7536	8726	8571141	4811	643607	
7	12684349	4557	15298	7394	8835	8727254	4762	586674	
8	12585563	4454	15497	7236	8898	8888173	5058	527691	
9	12468622	4309	14704	6942	8948	9076851	4673	481170	
10	12491183	4437	14567	6912	9103	9237982	4702	443329	
11	12430899	4296	14083	6515	9313	9364121	4609	404431	
12	12419506	4226	13985	6503	9342	9357468	4367	371475	
13	12469412	4147	13851	6375	9586	9386737	4588	345390	
14	12549936	4045	13650	6246	9673	9324488	4628	322294	
15	12566555	4174	13499	6213	9735	9305820	4518	301360	
-1	11599167	8800	16164	14851	90888	9613102	10843	19810	
-2	11985637	8769	14044	12040	28799	9561124	7184	18424	
-3	12941743	7805	13861	12001	24988	9400151	6368	15466	
-4	12808985	7141	12885	9889	23067	9509723	5421	14901	
-5	12869585	6954	12100	9428	22349	9464831	5789	13987	
-6	12784911	6440	12080	8735	20556	9566794	6544	14021	
-7	12878349	5946	12311	8225	19480	9566359	6478	16419	
-8	12719722	9521	12156	8131	19226	9725468	6709	23434	
-9	12652860	5634	11940	7671	18035	9762224	6321	31667	
-10	12566817	5448	11850	7178	17353	9701382	6306	37831	
-11	12702498	5309	12092	7568	16121	9526031	6035	43215	
-12	12731940	5207	11933	6856	15637	9533858	5557	47650	
-13	12697647	4989	12199	7153	15072	9508117	5434	51614	
-14	12689924	4944	11891	6816	15050	9525285	5237	bioRxiv preprint doi: https://doi.org/10.1101/555982; this version posted April 29, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a CC-BY-NC-ND 4.0 International license.	
-15	12660634	4746	11753	6732	14815	9561359	5184	59633	

582

584

**Appendix 3—table 1.** The read counts per position given the reference nucleotides are A or C of an ancient human data. The negative position indices are the position on the reversed strand. In the manuscript, the elements (the values of a specific nucleotide read counts per position given the reference nucleotide is A or C) in this table are denoted as  $M_{A \rightarrow i}(x)$  or  $M_{C \rightarrow i}(x)$ .

Ref.	Read Counts								
	G					T			
	Read	A	C	G	T	A	C	G	T
1	16389	8976	9639767	86584	11733	15878	8351	11718463	
2	17614	6483	9510149	26655	10761	13958	7011	11974947	
3	15164	5949	9488917	23374	9509	13767	6046	12839015	
4	14844	5186	9566468	21960	8170	12509	5585	12721790	
5	14005	5612	9497118	20468	7186	11991	5233	12795244	
6	13671	6195	9622572	19096	6948	11683	4790	12686645	
7	16648	6394	9609855	18594	6203	12122	4780	12794172	
8	23659	6405	9768666	17341	6131	11847	4758	12626614	
9	31680	6139	9785449	17034	5998	12040	4469	12579260	
10	38484	5982	9700857	16235	5487	11546	4175	12513653	
11	44665	5722	9536341	15284	5651	12044	4176	12646627	
12	48949	5371	9547134	14569	5449	11663	4060	12684645	
13	53076	5234	9543953	14090	5262	11785	4046	12631297	
14	57343	5186	9551477	13855	5257	11768	4006	12624840	
15	61236	5137	9583481	13667	5122	11733	3947	12612416	
-1	2078554	7947	8096447	11847	15732	28461	8551	12890628	
-2	1138478	6656	8232666	10760	12299	20759	6999	13446882	
-3	921712	5970	8399013	8643	10514	18226	6564	12718084	
-4	775038	5720	8319235	8416	9415	17800	5388	12977322	
-5	710955	5499	8462058	8926	8526	17088	4911	12886576	
-6	647761	5052	8545455	9193	7640	16351	4879	12852322	
-7	593854	4872	8693834	9318	7600	15523	5048	12664576	
-8	535542	7828	8889921	9399	7163	18704	4718	12510123	
-9	486549	4696	9075263	9522	7109	14547	4611	12409220	
-10	448895	4622	9226758	9432	6816	14567	4668	12438344	
-11	409027	4654	9352528	9544	6575	14019	4611	12388650	
-12	376069	4637	9344701	9419	6511	13874	4486	12390148	
-13	350609	4655	9384853	9885	6197	13877	4327	12432024	
-14	326760	4595	9337266	9889	5986	13928	4403	12490990	
-15	305014	4541	9310617	10065	5919	13442	4232	12529684	

**Appendix 3—table 2.** The read counts per position given the reference nucleotides are G or T of the same human data as in Table 1. The negative position indices are the position on the reversed strand. In the manuscript, the elements (the values of a specific nucleotide read counts per position given the reference nucleotide is G or T) in this table are denoted as  $M_{G \rightarrow i}(x)$  or  $M_{T \rightarrow i}(x)$ .

The terminology used here might not be standard. The term full regression here is to distinguish itself from the folded regression discussed later, which simply means inferring the coefficients of forward strand and reversed strand separately. Full regression includes both unconditional regression and conditional regression. The unconditional regression's objective is to infer the probability of observing a read of nucleotide  $j$  and its reference is  $i$  at position  $x$ , i.e.,  $P_{i \rightarrow j}(x)$  while the conditional regression's target is to estimate the probability of observing a read of nucleotide  $j$  given its reference is  $i$  at position  $x$ , i.e.,  $P_{j|i}(x)$ . Their relationship is as follows:

$$P_{j|i}(x) = \frac{P_{i \rightarrow j}(x)}{\sum_{j \in B} P_{i \rightarrow j}(x)}.$$

So in fact, unconditional regression can give us more detailed inferred results (extra information the nucleotide composition per position of the reference, which may be related to the prepared libraries).

### Unconditional Regression Likelihood

The unconditional regression's log-likelihood function is defined as follows,

$$\begin{aligned} l_1 &= \sum_x \sum_{i,j \in B} M_{i \rightarrow j}(x) \log P_{i \rightarrow j}(x) \\ &= \sum_x \left[ M(x) \log P_{T \rightarrow T}(x) + \sum_{(i,j) \neq (T,T)} M_{i \rightarrow j}(x) \log \frac{P_{i \rightarrow j}(x)}{P_{T \rightarrow T}(x)} \right], \end{aligned} \quad (10)$$

where  $M(x) = \sum_{i,j \in B} M_{i \rightarrow j}(x)$ . According to the multinomial logistic regression, we assume,

$$\log \frac{P_{i \rightarrow j}(x)}{P_{T \rightarrow T}(x)} = \sum_{n=0}^{\text{order}} \alpha_{i,j,x,n} x^n \quad (11)$$

Applying Equation 11 to Equation 10, we have

$$l_1 = \sum_x \left\{ -M(x) \log \left[ 1 + \sum_{(i,j) \neq (T,T)} \exp \left( \sum_{n=0}^{\text{order}} \alpha_{i,j,x,n} x^n \right) \right] + \sum_{(i,j) \neq (T,T)} M_{i \rightarrow j}(x) \sum_{n=0}^{\text{order}} \alpha_{i,j,x,n} x^n \right\} \quad (12)$$

620 The number of inferred parameters ( $\alpha_{i,j,x,n}$ ), for the full conditional regression is  $30 \times (\text{order} + 1)$ .  
 And the relevant derivatives of the unconditional regression likelihood are as follows,

622

$$\frac{\partial l_1}{\partial \alpha_{i,j,x,n}} = -M(x) \frac{x^n \exp\left(\sum_{n=0}^{\text{order}} \alpha_{i,j,x,n} x^n\right)}{1 + \sum_{(i,j) \neq (T,T)} \exp\left(\sum_{n=0}^{\text{order}} \alpha_{i,j,x,n} x^n\right)} + M_{i \rightarrow j}(x) x^n. \quad (13)$$

626  
624

### Conditional Regression Likelihood

628 Viewed as the sum of log-likelihoods given the reference nucleotide  $i \in \mathcal{B}$ , the conditional regression's log-likelihood function is,

630

$$l_2 = \sum_{i \in \mathcal{B}} \sum_x \sum_{j \in \mathcal{B}} M_{i \rightarrow j}(x) \log P_{j|i}(x)$$

632

$$= \sum_{i \in \mathcal{B}} \sum_x \left[ M_i(x) \log P_{T|i}(x) + \sum_{j \neq T} M_{i \rightarrow j}(x) \log \frac{P_{j|i}(x)}{P_{T|i}(x)} \right], \quad (14)$$

where  $M_i(x) = \sum_{j \in \mathcal{B}} M_{i \rightarrow j}(x)$ . Furthermore, if we assume,

634

$$\log \frac{P_{j|i}(x)}{P_{T|i}(x)} = \sum_{n=0}^{\text{order}} \beta_{i,j,x,n} x^n \quad (15)$$

By applying Equation 15 to Equation 14, we can obtain,

636

$$l_2 = \sum_{i \in \mathcal{B}} \sum_x \left\{ -M_i(x) \log \left[ 1 + \sum_{j \neq T} \exp\left(\sum_{n=0}^{\text{order}} \beta_{i,j,x,n} x^n\right) \right] + \sum_{j \neq T} M_{i \rightarrow j}(x) \sum_{n=0}^{\text{order}} \beta_{i,j,x,n} x^n \right\} \quad (16)$$

638

The number of inferred parameters ( $\beta_{i,j,x,n}$ ) for the full unconditional regression is  $24 \times (\text{order} + 1)$ . And the relevant derivatives of the conditional likelihood are as follows,

640

$$\frac{\partial l_2}{\partial \beta_{i,j,x,n}} = -M_i(x) \frac{x^n \exp\left(\sum_{n=0}^{\text{order}} \beta_{i,j,x,n} x^n\right)}{1 + \sum_{j \neq T} \exp\left(\sum_{n=0}^{\text{order}} \beta_{i,j,x,n} x^n\right)} + M_{i \rightarrow j}(x) x^n. \quad (17)$$

642

### Folded Regression

The folded regressions use the same log-likelihood functions as the full regression (i.e., Equation 12 and 16) but are conducted based on a presumable symmetric PMD pattern, i.e., the probability of  $C \rightarrow T$  at the position  $x$  of an random chosen ancient DNA strand is assumed to equal to the probability of  $G \rightarrow A$  at the position  $-x$ . Such an theoretical assumption go match the current ancient library preparation process [Meyer's paper and

Rasmus H's paper].

$$\alpha_{i,j,x,n} = \alpha_{c(i),c(j),-x,n}, \quad (18)$$

$$\beta_{i,j,x,n} = \beta_{c(i),c(j),-x,n}, \quad (19)$$

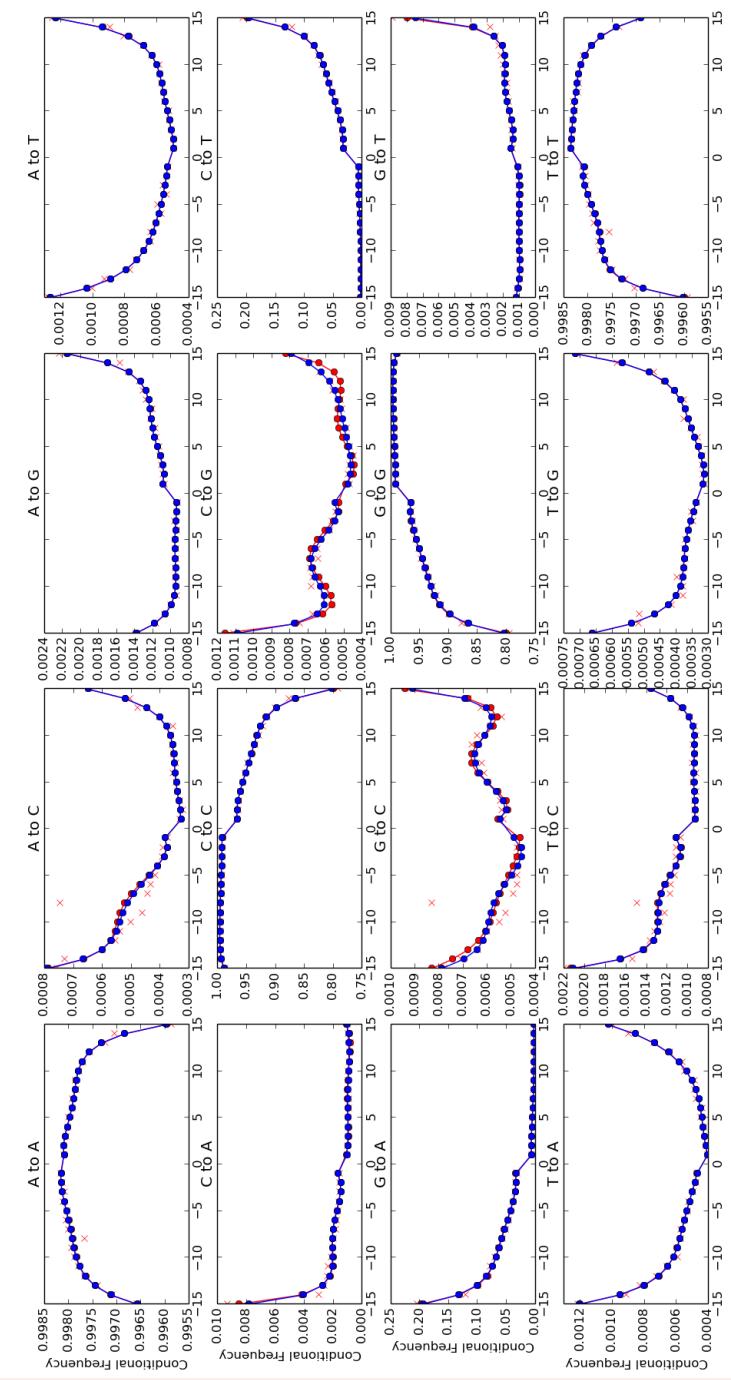
where  $c(i)$  means the complimentary nucleotide of the nucleotide  $i$ , e.g.,  $c(A) = T$  and  $c(G) = C$ .

By doing the folded regression, we halve the number of inferred parameters ( $\alpha_{i,j,x,n}$  or  $\beta_{i,j,x,n}$ ). Hence The number of inferred parameters for the folded unconditional regression is  $15 \times (\text{order} + 1)$ , and that of folded conditional regression is  $12 \times (\text{order} + 1)$ .

## Results for multinomial logistic regression

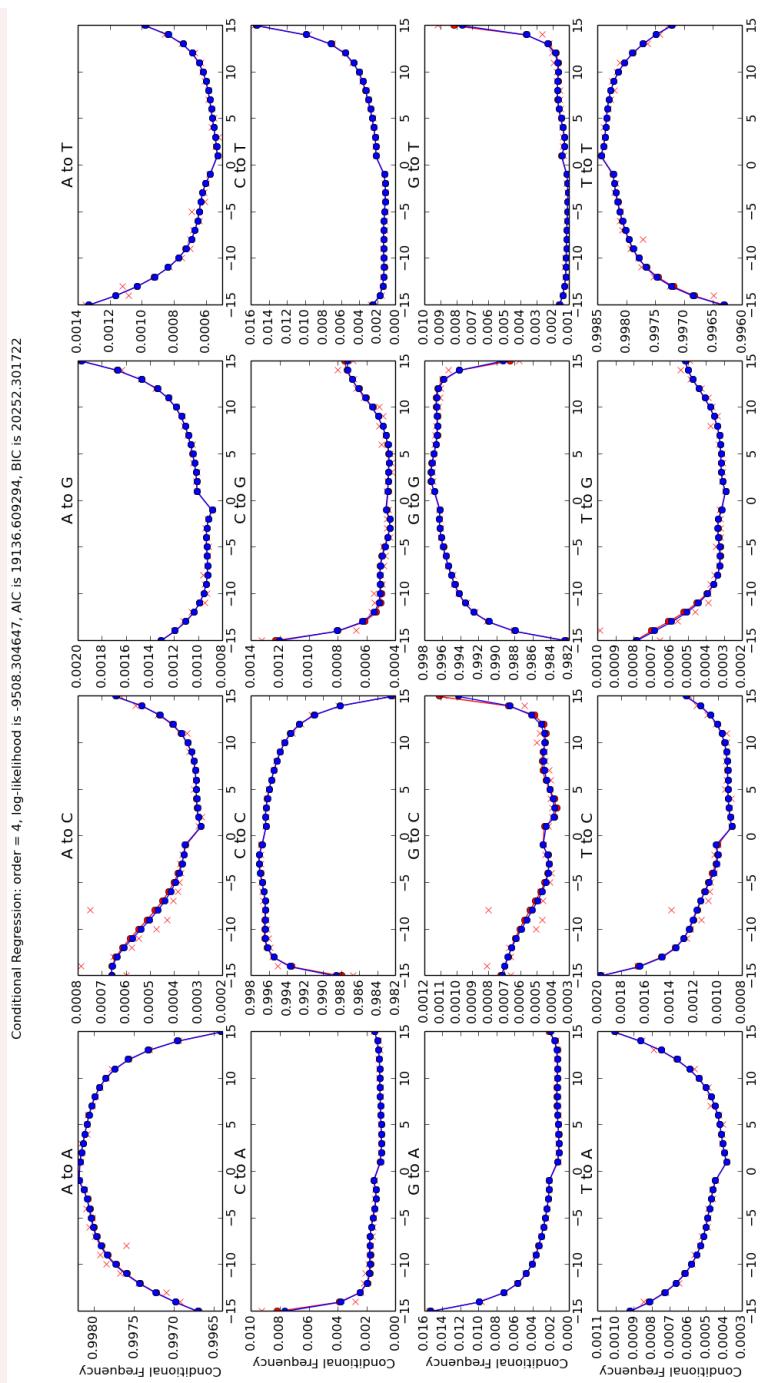
Currently, the optimization of the likelihood functions are based on the C++ library of gsl and use the function `gsl_multimin_fminimizer_nmsimplex2`. with the initial searching point is set to be the results of logistic regression. We here present here 4 figures pertaining to showcase the performance of our model. The regression methods are based on the summary statistic of the counts of mismatches and the optimization is therefore in the scale of miliseconds. Fig. 1 and Fig. 2 are the conditional regression results of the ancient and control human data correspondingly. And Fig. 3 and Fig. 4 are the folded conditional regression results of the same data as above. Our codes can also do the unconditional regression, but I have not generated the results for now.

Conditional Regression: order = 4, log-likelihood is -34526.568889, AIC is 69173.137778, BIC is 70313.881805



**Appendix 3—figure 1.** Conditional regression results with the order 4 of the ancient human data.

Each panel of figure represents a specific reference/read pair and plots its frequency across different positions. The positions from left to right are  $-1$  to  $-15$  and  $15$  to  $1$ .

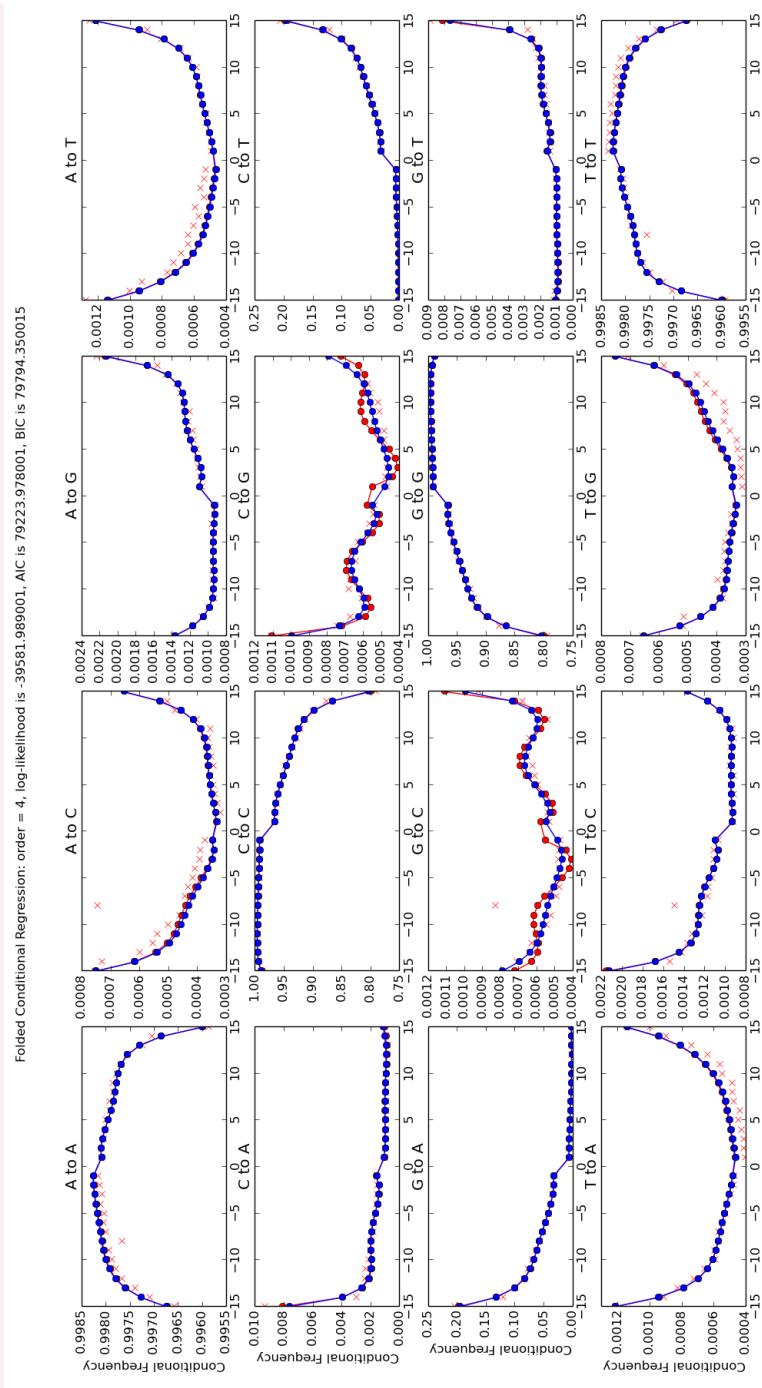


674

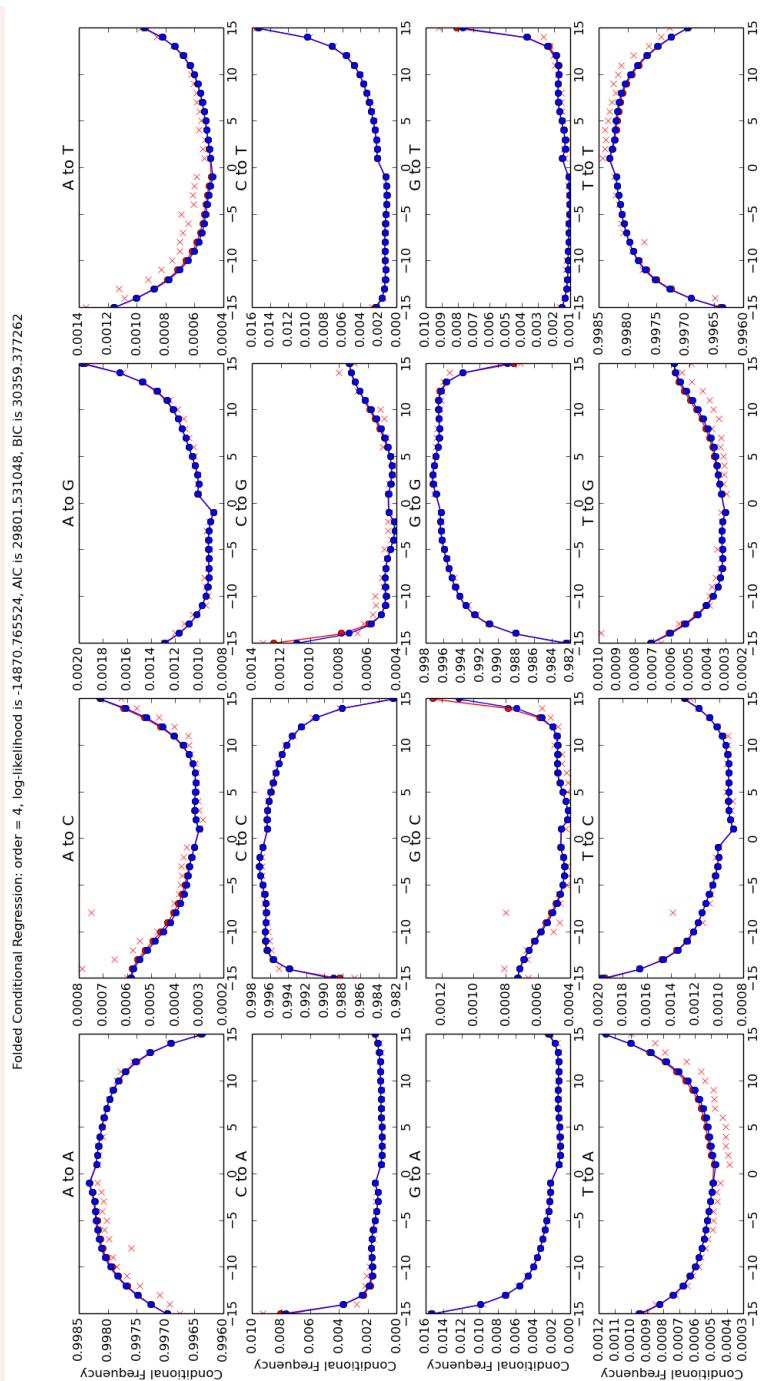
**Appendix 3—figure 2.** Conditional regression results with the order 4 of the control human data.

676

Each panel of figure represents a specific reference/read pair and plots its frequency across different positions. The positions from left to right are  $-1$  to  $-15$  and  $1$  to  $15$ .



**Appendix 3—figure 3.** Folded conditional regression results with the order 4 of the ancient human data. Each panel of figure represents a specific reference/read pair and plots its frequency across different positions. The positions from left to right are  $-1$  to  $-15$  and  $15$  to  $1$ .



684

**Appendix 3—figure 4.** Folded conditional regression results with the order 4 of the control human data. Each panel of figure represents a specific reference/read pair and plots its frequency across different positions. The positions from left to right are  $-1$  to  $-15$  and  $15$  to  $1$ .

686

Generally speaking, either the full multinomial regression or conditional regression, though describing a much more detailed PMD pattern, could suffer from an overfitting issue when the data is limited, while the simpler regression model in the main text shows an accept-

690

able statistic power even with extremely small amount of data [A figure to cite?], we thus

recommend the readers to use the simpler regression model when less data is applied.

692

## D | PMDTOOLS

We use a way introduced by (Skoglund et al., 2014) to fish out the ancient strands with intensive PMD patterns from samples.

According to (Skoglund et al., 2014), three nonmutually exclusive events can lead to an observation of  $C \rightarrow T$  or  $G \rightarrow A$ , namely (i) a true biological polymorphism (occurring at rate  $\pi$ ), (ii) a sequence error (rate  $\epsilon$ , can be extracted from the base quality scores of the site on the sampled strand), and (iii) in the case of damaged DNA, the damaged nucleotide frequencies are assumed to be only related to its position from either termini of the ancient fragment ( $C \rightarrow T$  from 5' end, and  $G \rightarrow A$  from 3' end),

$$D_x = C + p(1-p)^{|x|}, \quad (20)$$

704

where  $C = 0.01$  and  $p = 0.3$  are both constants.

The observation "Match" is defined as the case when we observe a  $C$  at a position whose reference is also a  $C$  or a  $G$  at a position whose reference is also a  $G$ . And the observation "Mismatch" represents the situation when we get a  $T$  or an  $A$  at a position whose reference nucleotide is a  $C$  or a  $G$ , respectively. The likelihoods of whether or not a specific fragment is damaged given the observation are calculated in the subsequent subsections.

706

708

710

### Model with PMD

If a strand is damaged, the probability that we observe a "Match" event at position  $x$  of this strand can be viewed as the sum of probabilities of three mutually exclusive events: (i) no biological difference between the reference and the sampled nucleotide, no damage and no sequencing error, (ii) no biological difference, damaged but the sequencing error lead to a "Match" observation, and (iii) no damage, and both the sequencing error and the biological

718

divergence contribute to a "Match" observation,

720

$$P(\text{Match} | x, \text{PMD}) = (1 - \pi)(1 - \epsilon)(1 - D_x) + (1 - \pi)\epsilon D_x + \pi\epsilon(1 - D_x), \quad (21)$$

722

The likelihood that the focal strand is damaged given the observation at position  $x$  is  $S_x$  can then be calculated as follows,

724

$$L(\text{PMD} | S_x) = \chi_{S_x=\text{Match}} P(\text{Match} | x, \text{PMD}) + \chi_{S_x=\text{Mismatch}} P(\text{Mismatch} | x, \text{PMD}), \quad (23)$$

726

where  $\chi_{S_x}$  is an indicator function.

## Model without PMD

Similarly, the "Match" event at position  $x$  in the case without PMD (the NULL model) can be decomposed as two exclusive events: (i) no biological divergence and no sequencing error, or (ii) both biological divergence and sequencing error contribute to a "Match" observation. And we have the following equations,

730

$$P(\text{Match} | x, \text{NULL}) = (1 - \pi)(1 - \epsilon) + \pi\epsilon \quad (24)$$

732

$$P(\text{Mismatch} | x, \text{NULL}) = 1 - P(\text{Match} | x, \text{NULL}) \quad (25)$$

734

$$L(\text{NULL} | S_x) = \chi_{S_x=\text{Match}} P(\text{Match} | x, \text{NULL}) + \chi_{S_x=\text{Mismatch}} P(\text{Mismatch} | x, \text{NULL}) \quad (26)$$

736

Skoglund et al., 2014 defines the likelihood ratio of a strand between the PMD model and the NULL model as its postmortem damage score (PMDS),

738

$$\text{PMDS} = \log \left\{ \frac{\prod_x L(\text{PMD} | S_x)}{\prod_x L(\text{NULL} | S_x)} \right\}, \quad (27)$$

740

The strands with the PMDS exceeding a empirical p-value threshold (???) will be fished out as intensively damaged fragments.