

# metaDMG – An Ancient DNA Damage

## 2 Toolkit

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

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

✉ For correspondence:

christianmichelsen@gmail.com

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

(MW);

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

(TSK)

Korneliussen<sup>2</sup>  

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

8

---

<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 at 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.

- 10 1. **Motivation** Under favourable conditions DNA molecules can persist for more than two million year (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 that the taxa is authenticated as 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 estimating DNA damage for ancient metagenomes with tens and even hundreds of thousands of species.
- 24 2. **Methods** 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 classification,  
28 splitting reads by individual organisms, realigning these to the reference genome and lastly  
29 parse alignments to mapDamage2.0 (Jónsson et al., 2013). We have implemented a  
30 Bayesian approach that combines a modified geometric damage profile with a  
31 beta-binomial model to fit the entire model to the individual misincorporations at all  
32 taxonomic levels. metaDMG was hereafter benchmarked using sets of simulated data of single  
33 genomes and metagenomes. Lastly, it was tested on published datasets and its  
34 performance compared to existing methods.

3. **Results** We find metaDMG to be a factor of 10 faster than previous methods and more  
35 accurate – even for complex metagenomes with tens of thousands of species. Our  
36 simulations show that metaDMG can estimate DNA damage at taxonomic levels down to 100  
37 reads, that the estimated uncertainties decrease with increased number of reads and that  
38 the estimates are more significant with increased number of C to T misincorporations.

40 4. **Conclusion** metaDMG is a state-of-the-art program for ancient DNA damage estimation and  
41 further allows for the computation of nucleotide misincorporation, GC-content, and DNA  
42 fragmentation for both simple and complex ancient genomic datasets. Additionally it  
43 includes the PMDtool statistics (Skoglund et al., 2014) that allow for the extraction of  
44 individual reads with ancient damage, making it a complete package for ancient DNA  
45 damage authentication.

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

---

## 48 1 | INTRODUCTION

Throughout the life of an organism it contaminates its environment with DNA, cells, or tissue, thus  
50 leaving genetic traces behind. As the cell leaves its host, DNA repair mechanisms stops and the DNA  
51 is subjected to intra and extra cellular enzymatic, chemical, and mechanical degradation, resulting  
52 in fragmentation and molecular alterations that over time lead to the characteristics of ancient  
53 DNA (Briggs et al., 2007; Dabney, Meyer, and Pääbo, 2013). Ancient DNA (aDNA) has been shown  
54 to persist in a diverse variety of environmental contexts, e.g. within fossils such as bones, soft-  
tissue, and hair, as well as in geological sediments, archaeological layers, ice-cores, permafrost soil

56 for thousands and even millions of years (Kjaer et al in review), (Cappellini et al., 2018). Common  
for all is that they have an accumulated amount of deaminated cytosines towards the termini of  
58 the DNA strand, which, when amplified, results in misincorporations of thymines on the cytosines  
(Dabney, Meyer, and Pääbo, 2013; Ginolhac et al., 2011).

60 Even though postmortem DNA damage is characterized by the four Briggs parameters, see  
(Briggs et al., 2007), they are rarely used directly for asserting "ancientness". Researchers working  
62 with ancient DNA tend to simply use the empirical C → T on the first position of the fragment  
together with other supporting summary statistic of the experiment (Jónsson et al., 2013).

64 Estimating misincorporation due to DNA damage, molecule fragmentation, and nick frequencies  
have become standard for single individual sources like hair, bones, teeth and also applied on  
66 smaller subsets of species in ancient environmental metagenomes (Pedersen et al., 2016; Murchie  
et al., 2021; Zavala et al., 2021; Yucheng Wang, Pedersen, et al., 2021). While this is a relatively fast  
68 process for single individuals it becomes increasingly demanding, iterative, and time consuming as  
the samples and the diversity within increases, as in the case for metagenomes from ancient soil,  
70 sediments, dental calculus, coprolites, and other ancient environmental sources. It has therefore  
been practice to estimate damage for only the key taxa of interest in a metagenome, as metage-  
72 nomic samples easily includes tens of thousands of different taxonomic entities, which would make  
a complete estimate an impossible task. To overcome this limitation, we designed a program called  
74 metaDMG (pronounced metadamage) which includes test statistics that takes all relevant information  
provided alignments to both single or multiple reference genomes into account.

76 Our research shows that metaDMG is both faster at ancient DNA damage estimation, provides  
more accurate damage estimates is, and able to process complex metagenomes within hours in-  
78 stead of days as metaDMG is designed with the increasingly large datasets, that are currently gener-  
ated in the field of ancient environmental DNA, in mind. At the same time, it outperforms standard  
80 tools that estimate DNA damage for single genomes and samples with low complexity. Further-  
more, it can even compute a global damage estimate for a metagenome as a whole. Lastly, metaDMG  
82 is compatible with the NCBI taxonomy and use ngsLCA (Yucheng Wang, T. S. Korneliussen, et al.,  
2022) to perform a last common ancestor (LCA) classification of the aligned reads to get precise  
84 damage estimates at all taxonomic nodes. It also allows for custom taxonomies and thus also the  
use of metagenomic assembled genomes (MAGs) as references.

86 This paper is organized as follows. First we present the XXX, then we YYY. Finally we ZZZ.

**Table 1.** Metagenomic samples, Mikkel, XXX. “Name” is the name used throughout this paper. “Site” is the type of metagenomic site. “Type” is the type of XXX. “Age” is the approximate age of the sample in kyr Bp. “Sediment” is the name type of sediment. “Instrument” is the Illumina model. “Library” is the XXX, where D.S. means double stranded and S.S. means single stranded. “Reads” is the number of reads (in millions) after filtering and trimming. “Source” is the source of the data.

Name	Site	Type	Age (kyr)	Sediment	Instrument	Library	Reads (M)	Source
Library-0	Control	Control	0	Reagents	HiSeq4000	D.S.	1.86	(Ardelean et al., 2020)
Pitch-6	Syltholmen pitch	Chewed organic material	5.7	Organic material	HiSeq2500	D.S.	95.59	(Jensen et al., 2019)
Lake-1	Spring Lake	Lake gyttja/sediment	1.4	Organic material	HiSeq 100	D.S.	16.89	(Pedersen et al., 2016)
Lake-7	Lake CH12	Lake gyttja/sediment	6.7	Organic material	HiSeq2500	S.S.	102.48	(Schulte et al., 2021)
Lake-9	Spring Lake	Lake gyttja/sediment	9.2	Organic material	HiSeq 100	D.S.	73.02	(Pedersen et al., 2016)
Shelter-39	Abri Pataud	Rock shelter	39.4	Sediment	MiSeq	S.S.	0.097	(Braadbaart et al., 2020)
Cave-22	Chiquihuite cave	Cave sediment	22.2	Carbonate rock	HiSeq4000	D.S.	4.75	(Ardelean et al., 2020)
Cave-100	Eustatas Cave	Cave sediment	100	Carbonate rock	HiSeq2500	S.S.	13.37	(Vernot et al., 2021)
Cave-102	Pesturina Cave	Neanderthal tooth	102	Dental calculus	HiSeq4000	D.S.	10.79	(Fellows Yates et al., 2021)

## 2 | METHODS & MATERIALS

88 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  
 90 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,  
 92 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  
 94 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  
 96 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  
 98 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  
 100 separate modern from ancient (data with PMD) (Skoglund et al., 2014). Expanding from the sin-  
 gle read all reads for a sequencing experiment and genome to tabulate the overall substitution or  
 102 mismatch rates to obtain a statistic of the damage (Borry et al., 2021) or even estimate the four  
 Briggs parameters that is traditionally used to characterize the damage signal (Jónsson et al., 2013).

We build a general ancient DNA damage toolkit, that accepts metagenomic datasets, and which  
 106 implements and expands on existing methods by implementing novel state-of-the-art methodolo-  
 gies.

## 108 2.1 | Damage at single read level

Firstly, we implemented the approach given in (Skoglund et al., 2014) which allows for extraction of  
 110 only damaged DNA reads. Three non-mutually exclusive events can lead to an observation of  $C \rightarrow$   
 $T$  or  $G \rightarrow A$  (Skoglund et al., 2014), namely (i) a true biological polymorphism (occurring at rate  $\pi$ ), (ii)  
 112 a sequencing errors (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  
 114 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). The error probability of the postmortem nucleotide misincorporation is  
 116 under the pmdtools model given by:

$$118 D_x = C + p(1 - p)^{|x|}, \quad (1)$$

here  $C = 0.01$  and  $p = 0.3$  are both suitable constants. Skoglund et al., 2014 defines the likelihood  
 120 ratio of a strand between the PMD model and the NULL model as its postmortem damage score  
 (PMDS),

$$122 \text{PMDS} = \log \left\{ \frac{\prod_x L(\text{PMD} | S_x)}{\prod_x L(\text{NULL} | S_x)} \right\}, \quad (2)$$

124 The reads with the PMDS exceeding an empirical p-value threshold can then be used for filtering  
 intensively damaged fragments. See Appendix A for details and derivation. The test statistic mimics  
 126 and reimplements the method described (Skoglund et al., 2014) and we defer to that article for  
 discussion in which scenarios this method is suitable.

## 128 2.2 | Mismatch matrices/nucleotide misincorporation patterns

We next expanded the singular read level to a summary statistic of all aligned sequencing data  
 130 across multiple reads by generating what is called a mismatch matrix or nucleotide misincorpora-  
 tion matrix. This matrix represent the nucleotide substitution frequencies across reads and pro-  
 132 vides us with the position dependent mismatch matrices,  $\underline{\underline{M}}(x)$ , with  $x$  denoting the position in the  
 read, starting from 1. At a specific position,  $M_{\text{ref} \rightarrow \text{obs}}(x)$  describes the number of nucleotides that

134 was mapped to a reference base  $B_{ref}$  but observed to be  $B_{obs}$ , where  $B \in \{A, C, G, T\}$ . The number  
136 of C to T transitions, e.g., is denoted as  $M_{C \rightarrow T}(x)$ . When tabulating these counts it is possible to  
138 disregard an entire read if it has a low mapping quality or the specific nucleotide if the basequality  
140 score fall below some threshold. Theoretical it is also possible to take into account the mapping  
quality and quality by not using an integer count but use the mapping quality and quality score  
as *weights* (Yi Wang et al., 2013), but given the four bin discretization of quality scores on modern  
day sequencing machines we do not model the counts probabilistically but solely use the mapping  
quality and quality scores as filters. See Appendix Table S2,S3 for an example of a mismatch matrix.

### 142 2.3 | Regression framework

The nucleotide misincorporation frequencies are routinely used as basis for assessing whether or  
144 not a given library is ancient. This is done by standardizing the substitution frequencies of the  
reference being C for the first few cycles and validating that the library exhibits the expected drop  
146 of C to T frequencies as we move through the position of the reads. This signal is caused by the  
higher deamination frequency in the single stranded part of the damaged fragment. Under the  
148 assumption of vast amounts of data we have defined a full multinomial regression model building  
on the method in (Cabanski et al., 2012).

150 This 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  
152 knowledge there are no currently available methods that is geared towards damage analysis in  
a metagenomic setting and existing approaches are essentially based on remapping against the  
154 single target organism and does not take into account any possible issues with regards to reads be-  
ing well assigned or specified. Our solution called `metaDMG` (pronounced metadamage), estimates  
156 the damage patterns in metagenomic samples in a three step approach. First, the lowest com-  
mon ancestor for each read (mapped to a multi-species reference database) is computed and the  
158 the mismatch matrix for each leaf node (e.g. taxonomic ID or contig, depending on the database  
used) is computed based on the mapped reads. Second, `metaDMG` fits a damage model to each leaf  
160 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  
162 interaction with the results for further downstream analysis and visualization.

## 2.4 | Lowest Common Ancestor and Mismatch matrices

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

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

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

When calculating the mismatch matrix, two different approaches can be taken. Either all align-  
190 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).

## 192 2.5 | Damage Estimation

The damage pattern observed in aDNA has several features which are well characterized. By modelling these, one can construct observables sensitive to aDNA signal. We model the damage patterns seen in ancient DNA by looking exclusively at the C→T transitions in the forward direction (5') and the G→A transitions in the reverse direction (3'). For each taxa, we denote the number of transitions,  $k(x)$ , as:

$$198 \quad k(x) = \begin{cases} M_{C \rightarrow T}(x) & \text{for } x > 0 \text{ (forward)} \\ M_{G \rightarrow A}(x) & \text{for } x < 0 \text{ (reverse)} \end{cases} \quad (3)$$

200 and the number of the reference counts  $N(x)$ :

$$202 \quad N(x) = \begin{cases} \sum_{i \in \{A,C,G,T\}} M_{C \rightarrow i}(x) & \text{for } x > 0 \text{ (forward)} \\ \sum_{i \in \{A,C,G,T\}} M_{G \rightarrow i}(x) & \text{for } x < 0 \text{ (reverse).} \end{cases} \quad (4)$$

The damage frequency is thus  $f(x) = k(x)/N(x)$ .

204 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 206 (overdispersion) we found in the data due to bad references and misalignments.

To accommodate overdispersion, we instead apply a beta-binomial distribution,  $P_{\text{BetaBinomial}}$ , which 208 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 210 function:

$$212 \quad P_{\text{BetaBinomial}}(k | N, \mu, \phi) = \binom{N}{k} \frac{B(k + \mu\phi, N - k + \phi(1 - \mu))}{B(\mu\phi, \phi(1 - \mu))}, \quad (5)$$

where  $B$  is defined as the beta function:

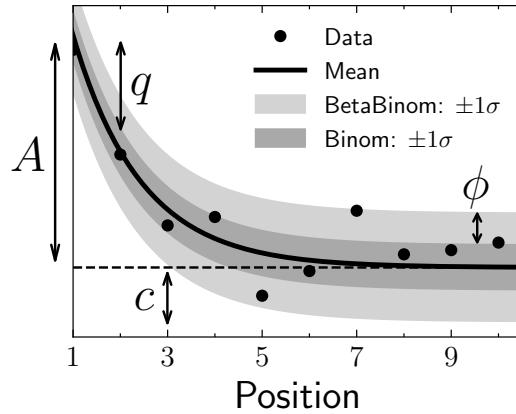
$$214 \quad B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}, \quad (6)$$

<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).

216 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 218 variance of a random variable  $k$  following a beta-binomial distribution,  $k \sim P_{\text{BetaBinomial}}(N, \mu, \phi)$ :

$$\begin{aligned} \mathbb{E}[k] &= N\mu \\ \mathbb{V}[k] &= N\mu(1 - \mu) \frac{\phi + N}{\phi + 1}. \end{aligned} \quad (7)$$



**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.

220 The expected value of  $k$  is similar to that of a binomial distribution and the variance of the beta-  
 221 binomial distribution reduces to a binomial distribution as  $\phi \rightarrow \infty$ . The beta-binomial distribution  
 222 can thus be seen as a generalization of the binomial distribution.

Note that both equation (5) and (7) relates to the damage at a specific base position, i.e. for a single  $k$  and  $N$ . To estimate the overall damage in the entire read using the position dependent counts,  
 224  $k(x)$  and  $N(x)$ , we model  $\mu$  as being position dependent,  $\mu(x)$ , and assume a position-independent  
 225 concentration,  $\phi$ . We model the damage frequency with a modified geometric sequence, i.e. exponentially decreasing for discrete values of  $x$ :

$$228 \quad \tilde{f}(x; A, q, c) = A(1 - q)^{|x| - 1} + c. \quad (8)$$

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

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

<sup>2</sup> Parameterized as  $(\mu, \phi)$

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

$$\begin{aligned} \text{[A prior]} & A \sim \text{Beta}(0.1, 10) \\ \text{240 [q prior]} & q \sim \text{Beta}(0.2, 5) \\ \text{[c prior]} & c \sim \text{Beta}(0.1, 10) \\ \text{242 [\phi prior]} & \phi \sim 2 + \text{Exponential}(1/1000) \\ \text{244 [likelihood]} & k_i \sim \mathcal{P}_{\text{BetaBinomial}}(N_i, \tilde{f}(x_i; A, q, c), \phi), \end{aligned} \quad (9)$$

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

246 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  
248 properties of the beta-binomial distribution, eq. (7), we find the mean and variance of the damage:

$$\begin{aligned} \mathbb{E}[D] & \equiv \bar{D} = A \\ \mathbb{V}[D] & \equiv \sigma_D^2 = \frac{A(1-A)}{N} \frac{\phi + N}{\phi + 1}. \end{aligned} \quad (10)$$

250 Since  $D$  estimates the overexpression of damage due to ancientness, not only is the mean of  
 $D$  relevant but also the certainty of it being non-zero (and positive). We quantify this through the  
252 significance  $Z = \bar{D}/\sigma_D$  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  
254 zero, i.e. containing 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  
256 ( $Z$ ) without having to run multiple models and comparing these.

We perform the Bayesian inference of the parameters models using Hamiltonian Monte Carlo  
258 (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,  
260 2019), a Python package which uses JAX (Bradbury et al., 2018) under the hood for automatic differ-  
entiation and JIT compilation. We treat each taxa as being independent and generate 1000 MCMC  
262 samples after an initial 500 samples as warm up.

Since running the full Bayesian model is computationally expensive, we also allow for a faster,  
264 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,  
266 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.

## 270 2.6 | Visualisation

We provide an interactive dashboard to properly visualise the results from the modelling phase,  
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  
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  
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  
Plotly plots and incorporated into a Dash dashboard (Plotly, 2015).

# 3 | SIMULATION STUDY

280 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  
282 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  
284 to determine the performance of the full metaDMG pipeline, i.e. both LCA and damage model.

## 3.1 | Single-genome Simulations

286 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  
288 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  
290 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,  $\delta_{ss}$ )  
292 in the original Briggs model (Briggs et al., 2007)), the number of reads, and the fragment length  
distribution.

```

./ngsngs -i $genome -r $Nread -ld LogNorm,$lognorm_mean,$lognorm_std -seq SE \
-f fq -q1 $quality_scores -m b,0.024,0.36,$damage,0.0097 -o $fastq
bowtie2 -x $genome -q $fastq.fq --no-unal

```

294 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  
 296 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  
 298 quantify the robustness of the estimates. We simulated eight different sets of damage (approxi-  
 mately 0%, 1%, 2%, 5%, 10%, 15%, 20%, 30%), 13 sets of different number of reads (10, 25, 50, 100, 250,  
 300 500, 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  
 302 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  
 304 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  
 306 genomes by sampling 1.000, 10.000 or 100.000 different basepairs from a uniform categorical dis-  
 tribution of {A, C, G, T}.

308 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.  
 310 The difference in damage frequency at position 1 and 15 is then the value to compare to:

$$D_{\text{known}} = \frac{f(x=1) - f(x=15)}{2} + \frac{f(x=-1) - f(x=-15)}{2}, \quad (11)$$

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

314 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  
 316 HiSeq 2500 (150 bp). The mapping was performed using Bowtie-2 with the –no-unal flag (Langmead  
 and Salzberg, 2012).

### <sup>318</sup> 3.2 | Metagenomic Simulations

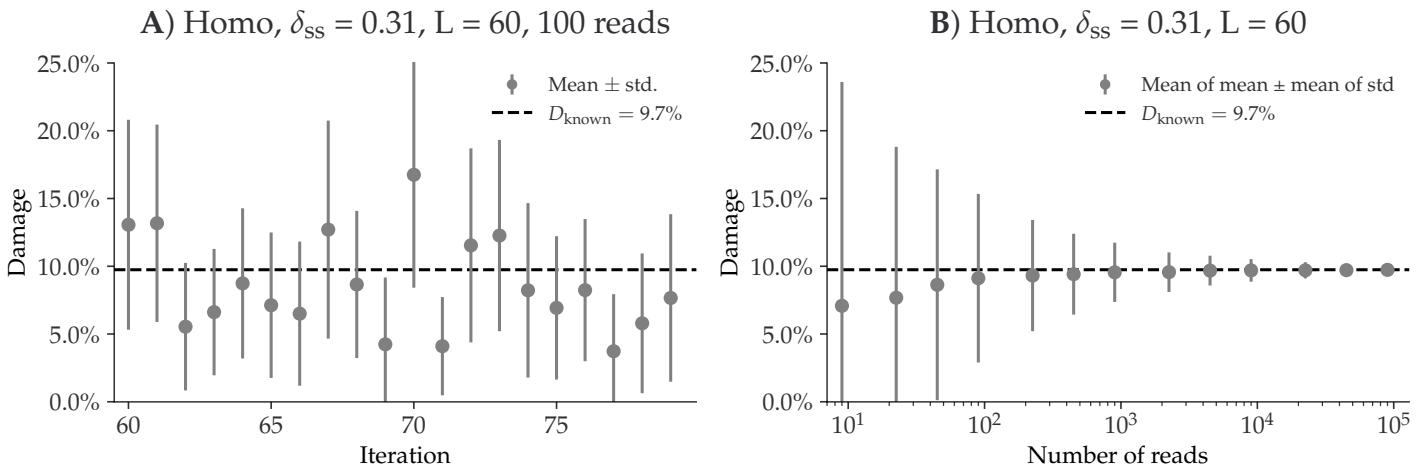
While the previously mentioned simulation study is perfectly aimed at quantifying the performance  
<sup>320</sup> of the damage model in the case of single-reference genomics it does lack the complexity related  
to metagenomic samples. Therefore, we also conduct a more advanced simulation study to deter-  
<sup>322</sup> mine the accuracy of the full `metaDMG` pipeline.

The previously mentioned simulation study quantifies the damage model's performance for  
<sup>324</sup> single-reference genomics, but it does not address the complexity of metagenomic studies. There-  
fore, we also conducted a more advanced simulation study to determine the performance of the  
<sup>326</sup> `metaDMG` pipeline in a standard eDNA setting.

Based on an ancient metagenome scenario, we created a synthetic dataset that mimics the  
<sup>328</sup> composition, fragment length distribution, and damage patterns for each genome. We selected 7  
metagenomes covering several environmental conditions and ages, based on **Table 1**. First, we  
<sup>330</sup> mapped the reads of each metagenome with `bowtie2` against a database that contained the GTDB  
r202 (Parks et al., 2018) species cluster reference sequences, all organelles from NCBI RefSeq (NCBI  
<sup>332</sup> Resource Coordinators, 2018), and the reference sequences 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 distribu-  
<sup>334</sup> tion for each genome and their abundance. The genomes with an observed-to-expected coverage  
ratio greater than 0.75 were kept. The filtered BAM files were processed by `metaDMG` to obtain the  
<sup>336</sup> misincorporation matrices. The abundance tables, fragment length distribution, and misincorpo-  
ration matrices were used in `aMGSIM-smk v0.0.1` (Fernandez-Guerra, 2022), a Snakemake workflow  
<sup>338</sup> (Mölder et al., 2021) that facilitates the generation of many synthetic ancient metagenomes. The  
data used and generated by the workflow can be obtained from Figshare link (XXX). We then per-  
<sup>340</sup> formed taxonomic profiling using the same parameters used for the synthetic reads generated by  
`aMGSIM-smk`.

## <sup>342</sup> 4 | RESULTS

The accuracy of the `metaDMG` pipeline was tested and validated in various simulation scenarios and  
<sup>344</sup> we applied it to a proper real metagenomic dataset. In general and across all scenarios we find that  
`metaDMG` yields accurate, precise damage estimates even in extreme low-coverage data.



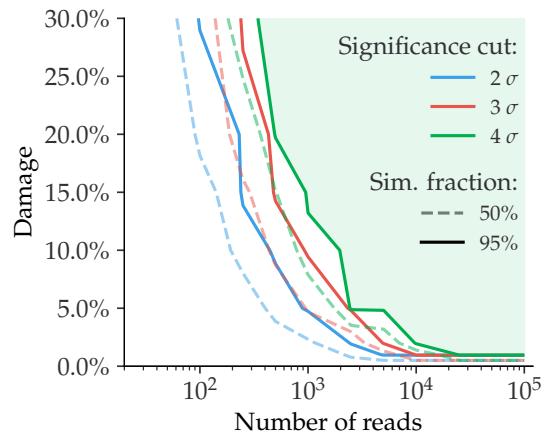
**Figure 2.** Overview of the single-genome simulations based on the homo sapiens genome with the Briggs parameter  $\delta_{SS} = 0.065$  and a fragment length distribution with mean 60. **A)** This plot shows the estimated damage ( $D$ ) of 10 simulations with 100 simulated reads. The grey points shows the mean damage (with its standard deviation as errorbars). The known damage ( $D_{known}$ ) is shown as a dashed line, see eq. (11). **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).

#### 346 4.1 | Single-genome Simulations

The results of the single-genome simulations can be seen in *Figure 2*. The left part of the figure  
 348 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  
 350 100 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  
 352 going into negative damage, which represents un-physical values. The right hand side of the figure  
 visualizes the average amount of damage across a varying number of reads. This shows that the  
 354 damage estimates converge to the known value with more data, and that one needs more than  
 356 100 reads to even get strictly positive damage estimates (when including uncertainties).

358 Across multiple simulations, each with 8 different damage levels, 13 different numbers of reads,  
 and 100 replications, we find no significant difference in the damage estimates across different  
 species (*Figure S2* and *Figure S3*), across different GC-levels (*Figure S4–Figure S6*), different frag-  
 360 ment length distributions (*Figure S7–Figure S9*), or different contig lengths (*Figure S10–Figure S12*),  
 see *Appendix 3*.

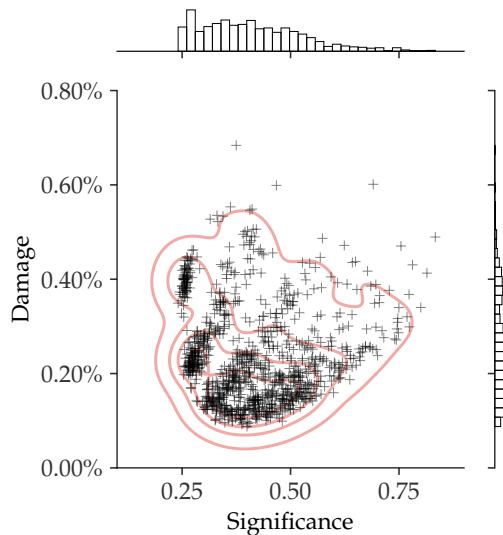
Based on the single-genome simulations, we can compute the relationship between the amount



**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.

362 of damage in a species and the number of reads required to correctly infer that the given species  
 363 is damaged, see **Figure 3**. If we want to find damage with a significance of more than 2 (solid blue  
 364 line) in a sample with around 5% expected damage, it requires about 1000 reads to be 95% certain  
 365 that we will find results this good. Said in other words: given 100 different fits, each with 1000  
 366 reads and around 5% damage, one would expect to find damage (with a  $Z > 2$ ) in 95 of the total  
 367 100 samples, on average. If we loose the requirement such that it is okay to only find it in every  
 368 second fit, it would be enough with only around 250 reads in each fit (dashed blue line).

Finally, to quantify the risk of incorrectly assigning damage to a non-damaged species, we cre-  
 370 ated 1000 independent simulations for a varying number of reads, where none of them had any  
 371 artificial ancient damage applied, only sequencing noise. **Figure 4** shows the damage ( $D$ ) as a  
 372 function of the significance ( $Z$ ) for the case of 1000 simulated reads. Even though the estimated  
 373 damage is larger than zero, the damage is non-significant since the significance is less than one.  
 374 When looking at all the figures across the different number of reads, see **Appendix 4**, we note that  
 375 a loose cut requiring that  $D > 1\%$  and  $Z > 2$  would filter out all of non-damaged points. Overall the  
 376 conclusion being that our devised test statistic is conservative and has low false positive rate.

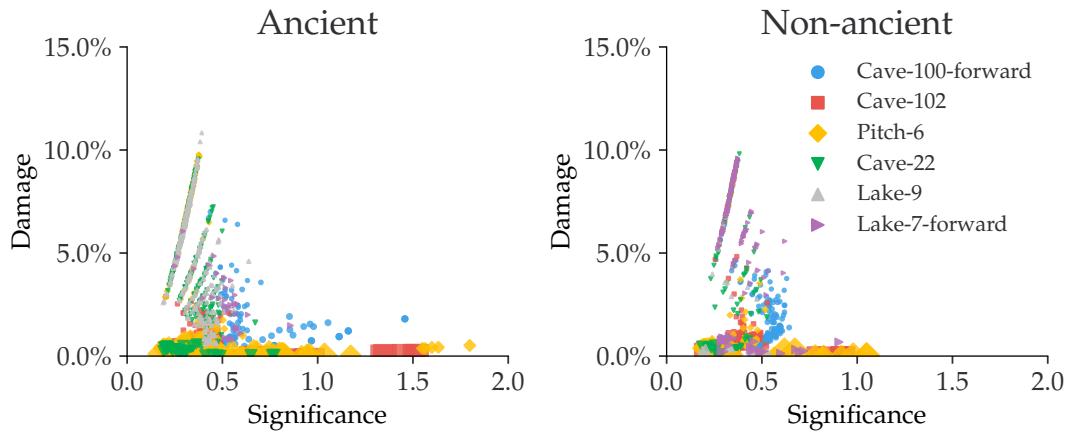


**Figure 4.** This figure shows the inferred damage estimates of 1000 independent simulations, each with 1000 reads and no artificial ancient damage applied, with the inferred damage shown on the y-axis and the significance on the x-axis. Each simulation is shown as a single cross and the red lines show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

## 4.2 | Metagenomic Simulations

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 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 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 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 filters out all of non-damaged points.

In comparison we can look at *Figure 6* which shows the same plot, but after the deamination (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 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



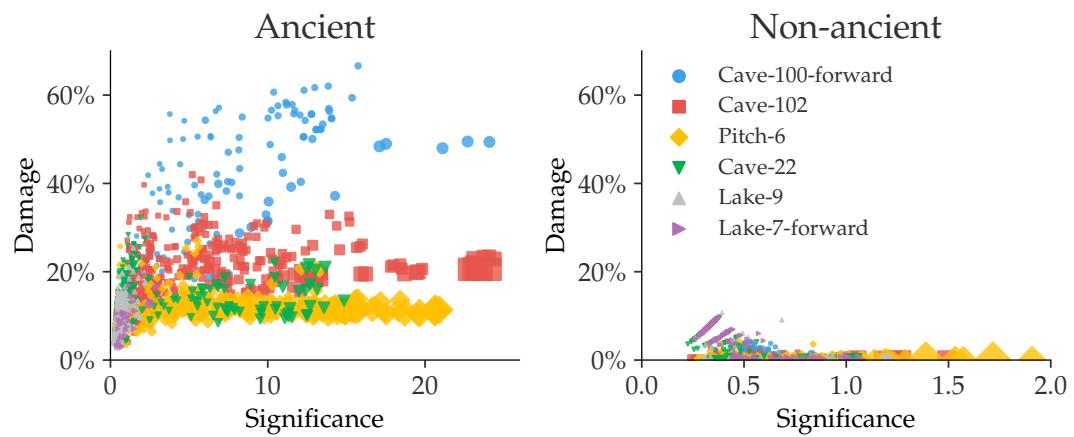
**Figure 5.** Estimated amount of damage as a function of significance using the fragSim data. The left figure shows the damage of the species that we simulated to be ancient (however with no deamination added yet) and the right figure shows the same for the species that are not going to have deamination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

are summarised in Table 2 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%) 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 ancient species correctly identified as ancient increases to more than 90% for most the samples.

To better understand the damage estimates, we can look at them individually. *Figure 7* shows 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 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 pass the loose cut. Furthermore, we notice that the error bars decrease with simulation size, as expected.

### 4.3 | Real Data

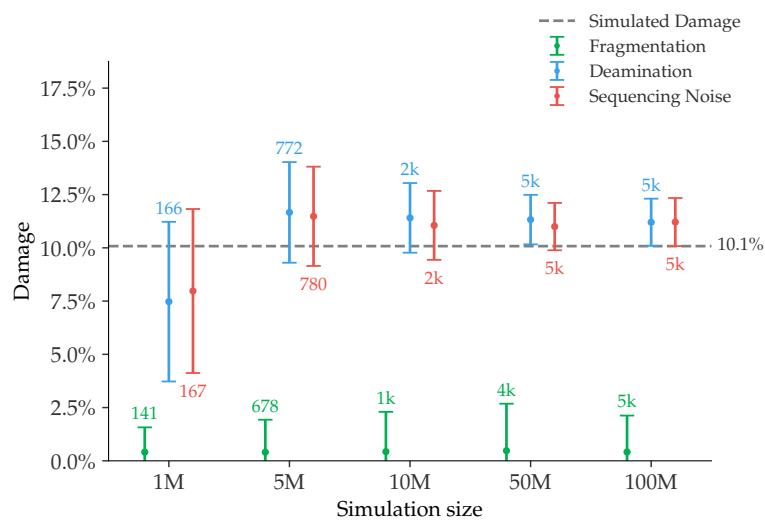
The results from running the full metaDMG pipeline on real data can be seen in *Figure 8*. The figures shows Blablabla, real life data here, XXX, Mikkel. We find that the loose cut ( $D > 1\%$ ,  $Z > 2$ ) accepts only one of the fits from the control test Library-0, which would not have been accepted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



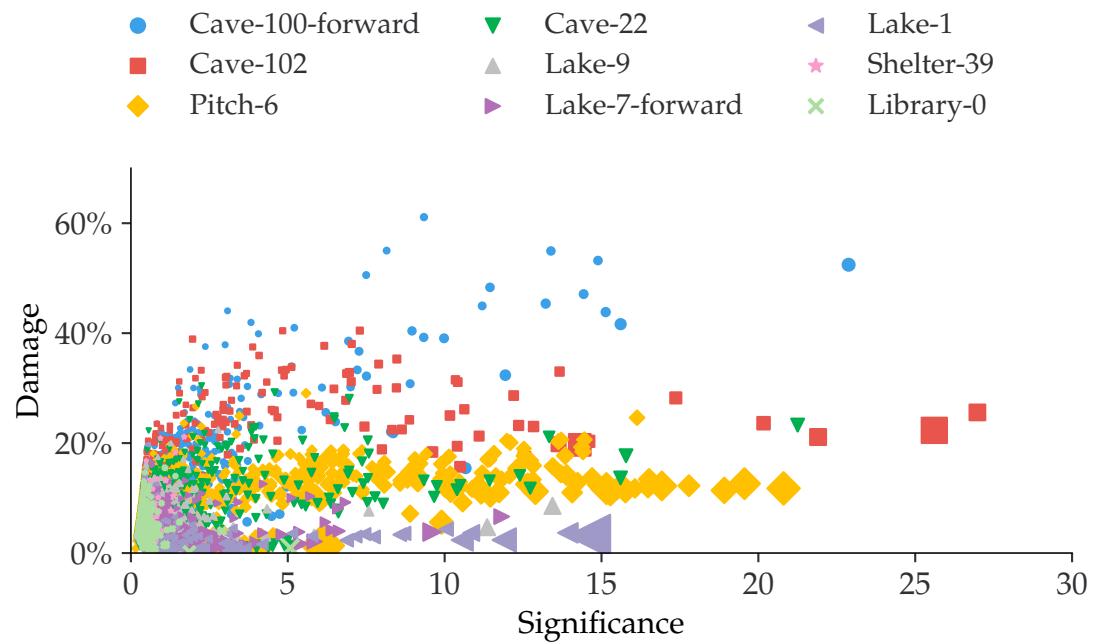
**Figure 6.** Estimated amount of damage as a function of significance using the ART data. The left figure shows the damage of the species that we simulated to be ancient and the right figure shows the same for the species that have not had deamination added.

**Table 2.** 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.

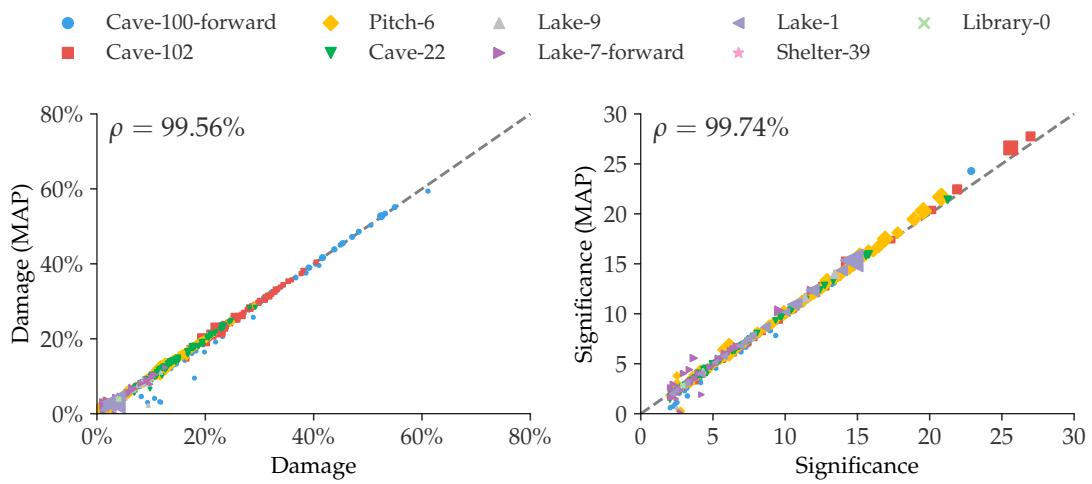
Sample	Total	Pass	+100 Reads	+100 Reads and Pass		
Cave-100-forward	135	107	79.3%	88	87	98.9%
Cave-102	500	326	65.2%	309	285	92.2%
Pitch-6	415	260	62.7%	274	260	94.9%
Cave-22	393	71	18.1%	73	69	94.5%
Lake-9	410	2	0.5%	8	0	0%
Lake-7-forward	32	4	12.5%	6	4	66.7%



**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.



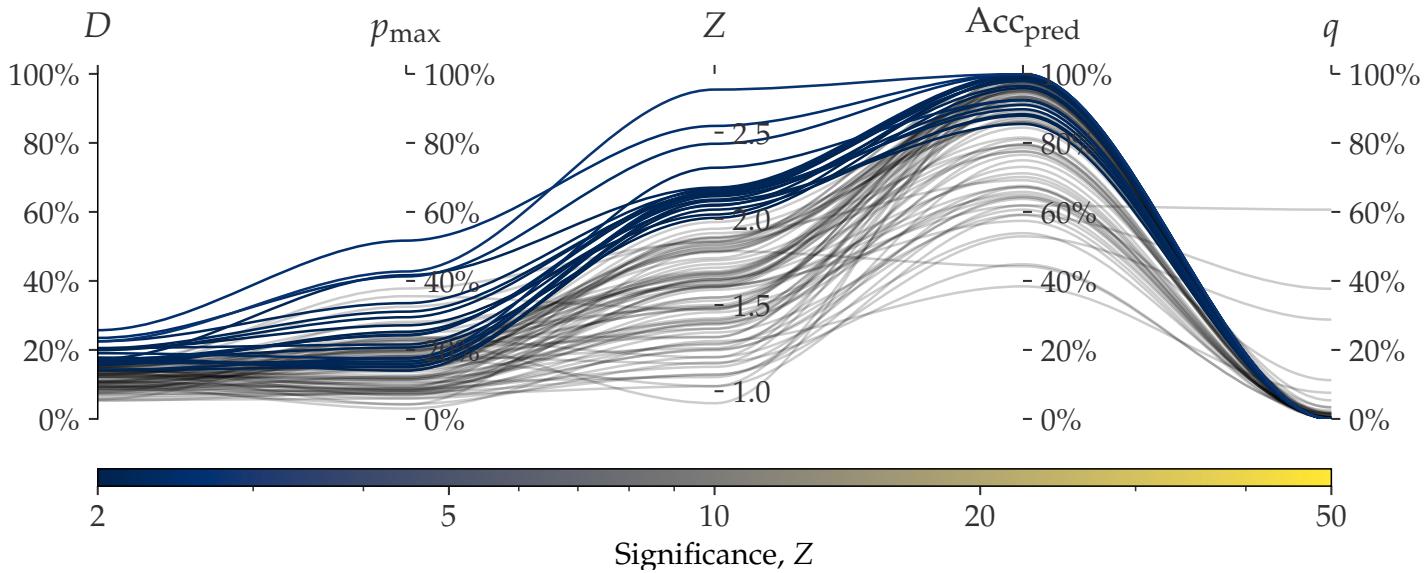
**Figure 9.** Comparison between the full Bayesian model and the fast, approximate, MAP model for the estimated damage and significance. The figure shows data after a loose cut of  $D > 1\%$ ,  $Z > 2$  and more than 100 reads. The dashed, grey line shows the 1:1 ratio and the correlation,  $\rho$ ,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

#### 4.4 | Bayesian vs. MAP

410 Due to increased computational burden of running the full Bayesian model compared to faster,  
 411 approximate MAP model, in samples with several thousand species, the MAP model is often the  
 412 most realistic model to use due to time constraints. In this case, it is of course important to know  
 413 that the damage estimates are indeed trustworthy. **Figure 9** compares the estimated damage  
 414 between the Bayesian model and the MAP model and the estimated significances for species with  
 415  $D > 1\%$ ,  $Z > 2$  and more than 100 reads. The figure shows that the vast majority of species map  
 416 1:1 between the Bayesian and the MAP model. One should note, though, that the few species with  
 417 the highest mismatch, all are based on forward-only fits, i.e. with no information from the reverse  
 418 strand, which thus leads to less data to base the fits on. For the comparison with no cuts, see  
**Figure S1** in appendix.

#### 420 4.5 | Existing Methods

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



**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.

426 One example of this can be found in **Figure 10**, which shows both the `metaDMG` and `PyDamage`  
 427 results of the 100 *Homo Sapiens* single-genome simulations with 100 reads and 15% added artificial  
 428 damage (and a fragment length distribution with mean 60). This figure shows that the `metaDMG`  
 429 estimates are between 5% and 25% damage, while `PyDamage` estimates up to more than 50%  
 430 damage, in a sample with 15% artificially added damage.

431 To compare the computational performance, we use the Pitch-6 sample, see **Table 1**. This align-  
 432 ment file (compressed to BAM-format) takes up 857 MB of space and has 3.7 millions reads with a  
 433 total of 19 million alignments to 11.433 unique taxa. When using only a single core, `PyDamage` took  
 434 1105 s to compute all fits, while `metaDMG` took 88 s, a factor of 12.6x faster. Out of the 88 s, `metaDMG`  
 435 spent 53 s on the actual fits, the rest was for loading and reading the alignment file and computing  
 436 the mismatch matrices. This makes `metaDMG` more than 20x faster than `PyDamage` for the fit compu-  
 437 tation. For the rest of the timings, see Table 3. `PyDamage` requires the alignment file to be sorted  
 438 by chromosome position and be supplied with an index file, allowing it to iterate fast through the  
 alignment file, at the expense of computational load before running the actual damage estimation.

**Table 3.** 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

440 metaDMG on the other hand requires the reads to be sorted by name to minimize the time it takes  
 441 to run the LCA, which however, is not tested in this comparison.

## 442 5 | DISCUSSION

443 As metaDMG is the first framework designed specifically for the use of estimating ancient damage  
 444 in a metagenomic context, multiple areas of future improvements exists. Currently, the damage  
 445 estimation is based on a statistical model which treats each leaf node in the taxonomic tree as  
 446 being fully independent, even for closely related species. This could be improved upon with the  
 447 use of a hierarchical model were information across taxonomic leaf nodes is shared. The current  
 448 implementation, however, allows for easy parallelization of the individual fits which reduces the  
 449 time spent on the inference. In addition to this basic assumption, another improvement would be  
 450 to include the read length distribution as a covariate in the damage model, as, in addition to deam-  
 451 ination, the fragment length distribution is also an indicator of ancient damage (Dabney, Meyer,  
 452 and Pääbo, 2013; Peyrégne and Prüfer, 2020).

453 We show that the damage estimates that metaDMG provides are accurate across different dam-  
 454 age levels and different number of reads. In the single-genome reference case, we further show  
 455 that the estimates are stable across different species and fragment length distributions. In addition  
 456 to this, we find that the results are independent of the contig size, in contrast to PyDamage (Borry  
 457 et al., 2021). Our research indicate that the metaDMG results are conservative with very low false pos-

itive rates. This is particularly important with metagenomic samples as the number of species, and thus the number of damage estimates, tend to be large. We have tested metaDMG using a state of the art metagenomic simulation pipeline based on multiple metagenomic real-life sample from a variety of different environments. In future studies, the simulation setup can further be improved by XXX (Mikkel, Antonio). We hope that metaDMG can improve the knowledge about DNA damage degradation in different environments and be the foundation of a more general, metagenomic ancient damage study.

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).

- no linkage
- weight
- improved fishing (pmdtools)

## 6 | DATA AVAILABILITY

Source code is hosted at GitHub: <https://github.com/metaDMG-dev>. Sequencing data and supporting material used in simulations can be found at: [https://sid.elda.dk/cgi-sid/ls.py?share\\_id=I7NGWfSkXq](https://sid.elda.dk/cgi-sid/ls.py?share_id=I7NGWfSkXq).

## 7 | COMPETING INTERESTS

The authors declare that they have no competing interests.

## 8 | FUNDING

CM and TP is funded by the Lundbeck Foundation. MWP and LZ is funded by the Lundbeck Foundation Centre for Disease Evolution Grant id: R302-2018-2155. TSK is funded by Carlsberg grant CF19-0712. AFG is funded by ?.

## 9 | AUTHOR CONTRIBUTIONS

- <sup>482</sup> CM developed and implemented the damage model and all aspect of the python code. TP helped developing the model and with statistical discussions. TSK implemented the C/C++ code relating to the lowest common ancestor and nucleotide misincorporation matrices. LZ implemented the PMDtools and full multinomial regression subfunctionality. AFG and MWP designed the metagenomic simulation study and the application of metaDMG to real data. CM and MWP ran all analyses.
- <sup>484</sup> CM, MWP and TSK initiated and devised the overall project. All authors contributed to writing the manuscript.
- <sup>486</sup> manuscript.

## REFERENCES

- 490 Ardelean, Ciprian F. et al. (2020). "Evidence of human occupation in Mexico around the Last Glacial Maximum". en. In: *Nature* 584.7819. Number: 7819 Publisher: Nature Publishing Group, pp. 87–92. ISSN: 1476-4687. DOI: [10.1038/s41586-020-2509-0](https://doi.org/10.1038/s41586-020-2509-0). URL: <https://www.nature.com/articles/s41586-020-2509-0> (visited on 2022).
- 492
- 494 Betancourt, Michael (2018). "A Conceptual Introduction to Hamiltonian Monte Carlo". In: *arXiv:1701.02434 [stat]*. arXiv: 1701.02434.
- 496 Bezanson, Jeff et al. (2017). "Julia: A fresh approach to numerical computing". In: *SIAM review* 59.1. Publisher: SIAM, pp. 65–98. URL: <https://julialang.org/>.
- 498 Borry, Maxime et al. (2021). "PyDamage: automated ancient damage identification and estimation 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).
- 500
- 502 Braadbaart, F. et al. (2020). "Heating histories and taphonomy of ancient fireplaces: A multi-proxy case study from the Upper Palaeolithic sequence of Abri Pataud (Les Eyzies-de-Tayac, France)". en. In: *Journal of Archaeological Science: Reports* 33, p. 102468. ISSN: 2352-409X. DOI: [10.1016/j.jasrep.2020.102468](https://doi.org/10.1016/j.jasrep.2020.102468). URL: <https://www.sciencedirect.com/science/article/pii/S2352409X20302595> (visited on 2022).
- 504
- 506 Bradbury, James et al. (2018). *JAX: composable transformations of Python NumPy programs*. URL: [http://github.com/google/jax](https://github.com/google/jax).
- 508 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 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>.
- 510
- 512 Cabanski, Christopher R. et al. (2012). "ReQON: a Bioconductor package for recalibrating quality 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).
- 514
- 516 Cappellini, Enrico et al. (2018). "Ancient Biomolecules and Evolutionary Inference". In: *Annual Review of Biochemistry* 87.1. \_eprint: <https://doi.org/10.1146/annurev-biochem-062917-012002>, pp. 1029–

- 518 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).
- 520 Cepeda-Cuervo, Edilberto and MARÍA VICTORIA Cifuentes-Amado (2017). "Double Generalized Beta-Binomial and Negative Binomial Regression Models". en. In: *Revista Colombiana de Estadística*
- 522 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).
- 524 Champlot, Sophie et al. (2010). "An efficient multistrategy DNA decontamination procedure of PCR reagents for hypersensitive PCR applications". eng. In: *PLoS One* 5.9, e13042. ISSN: 1932-6203.
- 526 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*
- 528 *Harbor Perspectives in Biology* 5.7, a012567. ISSN: 1943-0264. DOI: [10.1101/cshperspect.a012567](https://doi.org/10.1101/cshperspect.a012567).
- URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3685887/> (visited on 2022).
- 530 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).
- 532 Fellows Yates, James A. et al. (2021). "The evolution and changing ecology of the African hominid oral microbiome". en. In: *Proceedings of the National Academy of Sciences* 118.20, e2021655118.
- 534 ISSN: 0027-8424, 1091-6490. DOI: [10.1073/pnas.2021655118](https://doi.org/10.1073/pnas.2021655118). URL: <https://doi.org/10.1073/pnas.2021655118> (visited on 2022).
- 536 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>.
- 538 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://doi.org/10.1016/j.tree.2005.07.005> (visited on 2022).
- 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).
- 542 544 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>.
- 546 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).

- 548 Jensen, Theis Z. T. et al. (2019). "A 5700 year-old human genome and oral microbiome from chewed  
birch pitch". en. In: *Nature Communications* 10.1. Number: 1 Publisher: Nature Publishing Group,  
550 p. 5520. ISSN: 2041-1723. DOI: [10.1038/s41467-019-13549-9](https://doi.org/10.1038/s41467-019-13549-9). URL: <https://www.nature.com/articles/s41467-019-13549-9> (visited on 2022).
- 552 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.  
554 DOI: [10.1093/bioinformatics/btt193](https://doi.org/10.1093/bioinformatics/btt193).
- 556 Lam, Siu Kwan, Antoine Pitrou, and Stanley Seibert (2015). "Numba: a LLVM-based Python JIT com-  
piler". 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.  
558 DOI: [10.1145/2833157.2833162](https://doi.org/10.1145/2833157.2833162). URL: <https://github.com/numba/numba>.
- 560 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  
562 2022).
- 564 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: [10.1080/20548923.2016.1258824](https://doi.org/10.1080/20548923.2016.1258824). URL: <https://www.tandfonline.com/doi/full/10.1080/20548923.2016.1258824> (visited  
on 2022).
- 568 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-  
570 13991-9.
- 572 Mölder, Felix et al. (2021). *Sustainable data analysis with Snakemake*. en. Tech. rep. 10:33. Type: arti-  
cle. 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).
- 574 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  
576 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).

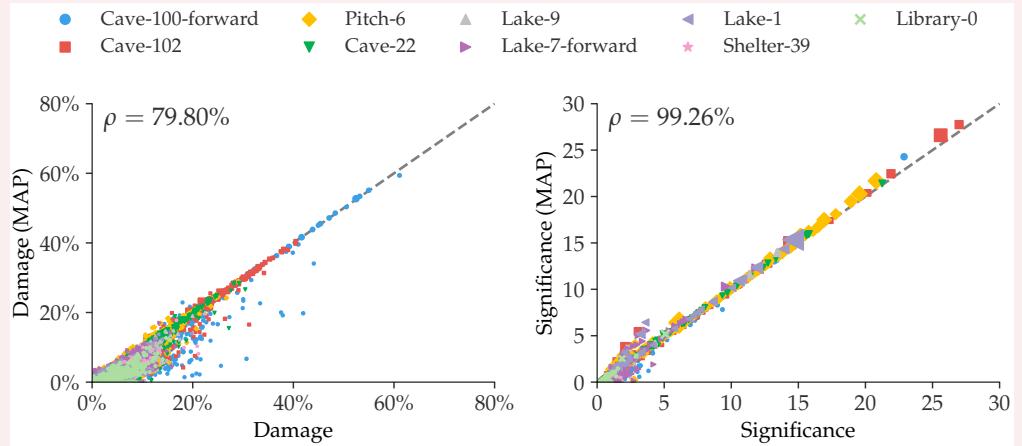
- 578 Nayfach, Stephen et al. (2021). "CheckV assesses the quality and completeness of metagenome-  
assembled viral genomes". en. In: *Nature Biotechnology* 39.5. Number: 5 Publisher: Nature Pub-  
lishing Group, pp. 578–585. ISSN: 1546-1696. DOI: [10.1038/s41587-020-00774-7](https://doi.org/10.1038/s41587-020-00774-7). URL: <https://www.nature.com/articles/s41587-020-00774-7> (visited on 2022).
- 580  
582 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).
- 584  
586 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).
- 588 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:  
590 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).
- 592 Pedersen, Mikkel 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–  
594 49. ISSN: 1476-4687. DOI: [10.1038/nature19085](https://doi.org/10.1038/nature19085). URL: <https://www.nature.com/articles/nature19085>  
(visited on 2022).
- 596 Peyrégne, Stéphane and Kay Prüfer (2020). "Present-Day DNA Contamination in Ancient DNA Datasets".  
en. In: *BioEssays* 42.9. \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1002/bies.202000081>,  
598 p. 2000081. ISSN: 1521-1878. DOI: [10.1002/bies.202000081](https://doi.org/10.1002/bies.202000081). (Visited on 2022).
- 600 Phan, Du, Neeraj Pradhan, and Martin Jankowiak (2019). "Composable Effects for Flexible and Accel-  
erated Probabilistic Programming in NumPyro". In: *arXiv:1912.11554 [cs, stat]*. arXiv: 1912.11554.
- 602 Plotly, Technologies (2015). *Collaborative data science*. Place: Montreal, QC Publisher: Plotly Tech-  
nologies Inc. URL: <https://plot.ly>.
- 604 Schulte, Luise et al. (2021). "Hybridization capture of larch (*Larix Mill.*) chloroplast genomes from  
sedimentary ancient DNA reveals past changes of Siberian forest". en. In: *Molecular Ecology Re-  
sources* 21.3. \_eprint: <https://onlinelibrary.wiley.com/doi/pdf/10.1111/1755-0998.13311>, pp. 801–  
606 815. ISSN: 1755-0998. DOI: [10.1111/1755-0998.13311](https://doi.org/10.1111/1755-0998.13311). URL: <https://onlinelibrary.wiley.com/doi/abs/10.1111/1755-0998.13311> (visited on 2022).

- 608 Skoglund, Pontus et al. (2014). "Separating endogenous ancient DNA from modern day contami-  
nation in a Siberian Neandertal". In: *Proceedings of the National Academy of Sciences* 111.6. Pub-  
lisher: 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).
- 610  
612 Vernot, Benjamin et al. (2021). "Unearthing Neanderthal population history using nuclear and mi-  
tochondrial DNA from cave sediments". In: *Science* 372.6542. Publisher: American Association  
for the Advancement of Science, eabf1667. DOI: [10.1126/science.abf1667](https://doi.org/10.1126/science.abf1667). URL: <https://www.science.org/doi/full/10.1126/science.abf1667> (visited on 2022).
- 614  
616 Wang, Yi et al. (2013). "An integrative variant analysis pipeline for accurate genotype/haplotype  
inference in population NGS data". eng. In: *Genome Research* 23.5, pp. 833–842. ISSN: 1549-5469.  
DOI: [10.1101/gr.146084.112](https://doi.org/10.1101/gr.146084.112).
- 618  
620 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". In: *Methods  
in Ecology and Evolution* n/a.n/a. Publisher: John Wiley & Sons, Ltd. ISSN: 2041-210X. DOI: [10.1111/2041-210X.14006](https://doi.org/10.1111/2041-210X.14006). URL: <https://doi.org/10.1111/2041-210X.14006> (visited on 2022).
- 622  
624 Wang, Yucheng, Mikkel 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).
- 626  
628 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).
- 630

## Appendix 1

632

### BAYES VS. MAP



634

**Appendix 1—figure S1.** 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.

636

## Appendix 2

### EXAMPLE TABLE

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

**Appendix 2—table S1.** An example table.

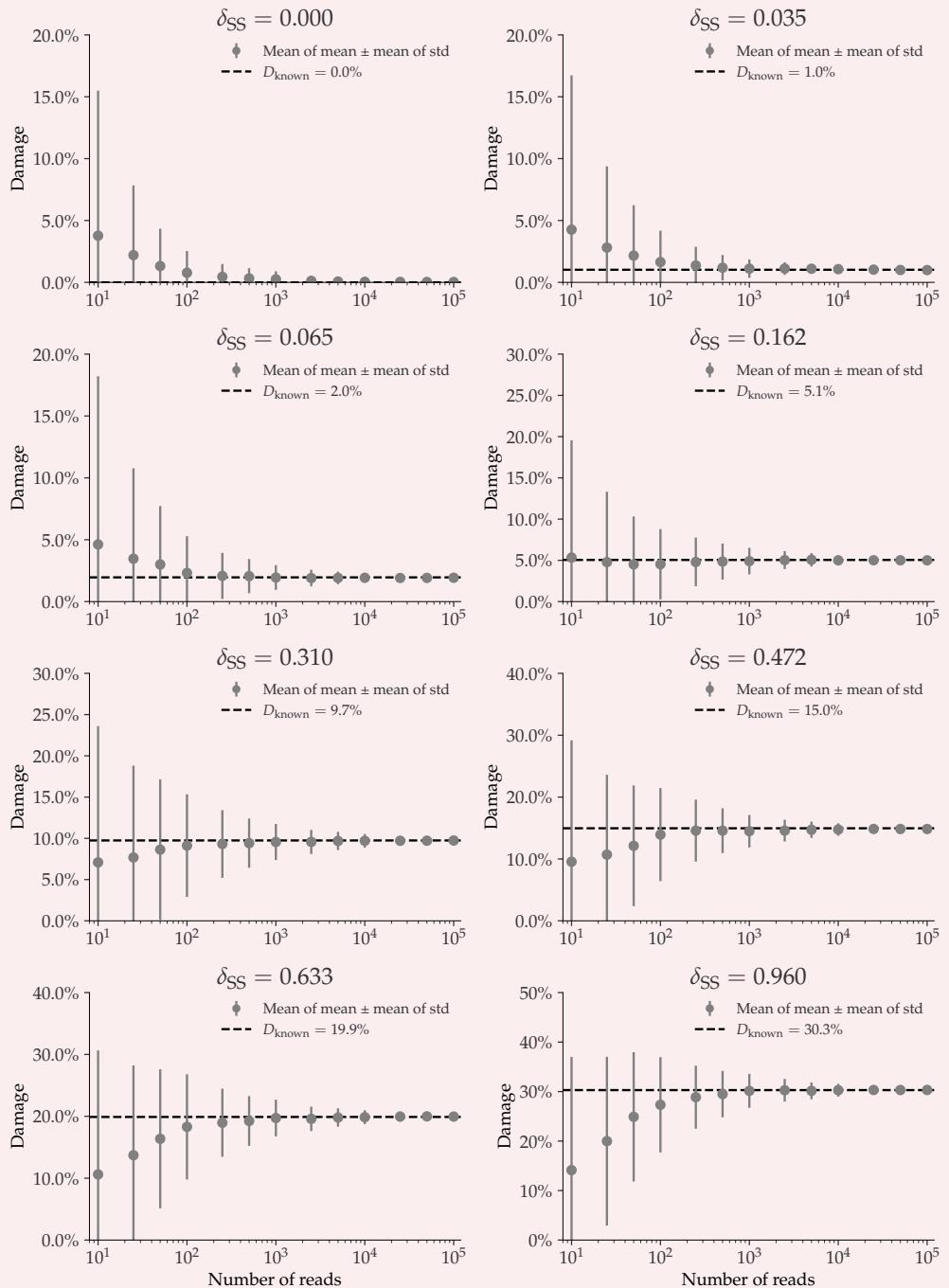
Speed (mph)	Driver	Car	Engine	Date
407.447	Craig Breedlove	Spirit of America	GE J47	8/5/63
413.199	Tom Green	Wingfoot Express	WE J46	10/2/64
434.22	Art Arfons	Green Monster	GE J79	10/5/64
468.719	Craig Breedlove	Spirit of America	GE J79	10/13/64
526.277	Craig Breedlove	Spirit of America	GE J79	10/15/65
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

## Appendix 3

### 644 NGSNGS SIMULATIONS

646 The following figures show the `metaDMG` damage estimates for the different NGSNGS simu-  
648 lations. These simulations include different species (*homo sapiens* and *betula*), different  
GC-levels (low, middle, high), different fragment length distributions (with mean 35, 60, and  
90), and different contig lengths (length 1.000, 10.000, 100.000).

## Homo



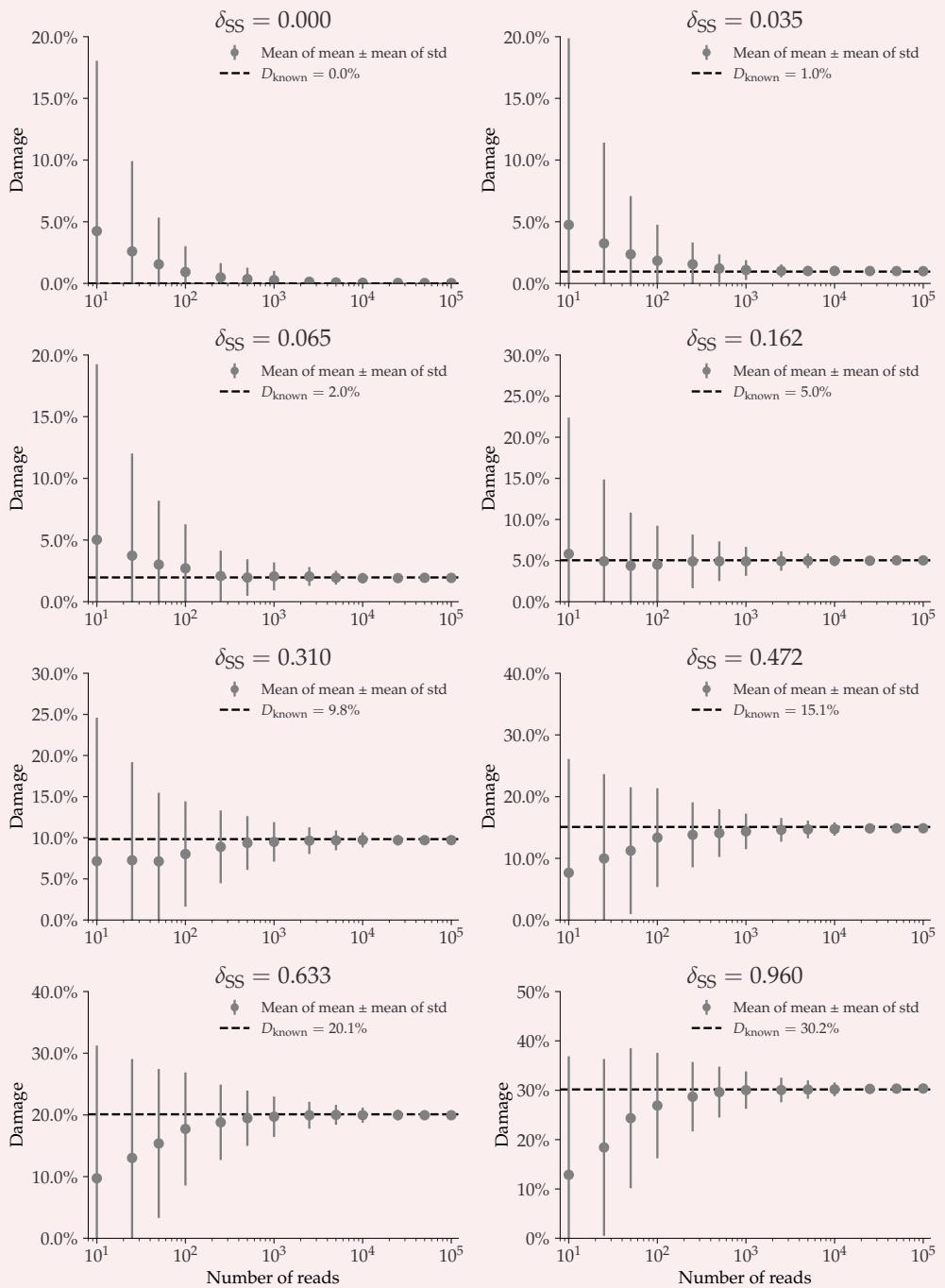
650

**Appendix 3—figure S2.** This plot shows the average damage as a function of the number of reads.

652

The grey points show the average of the individual means (with the average of the standard deviations as errors).

## Betula



654

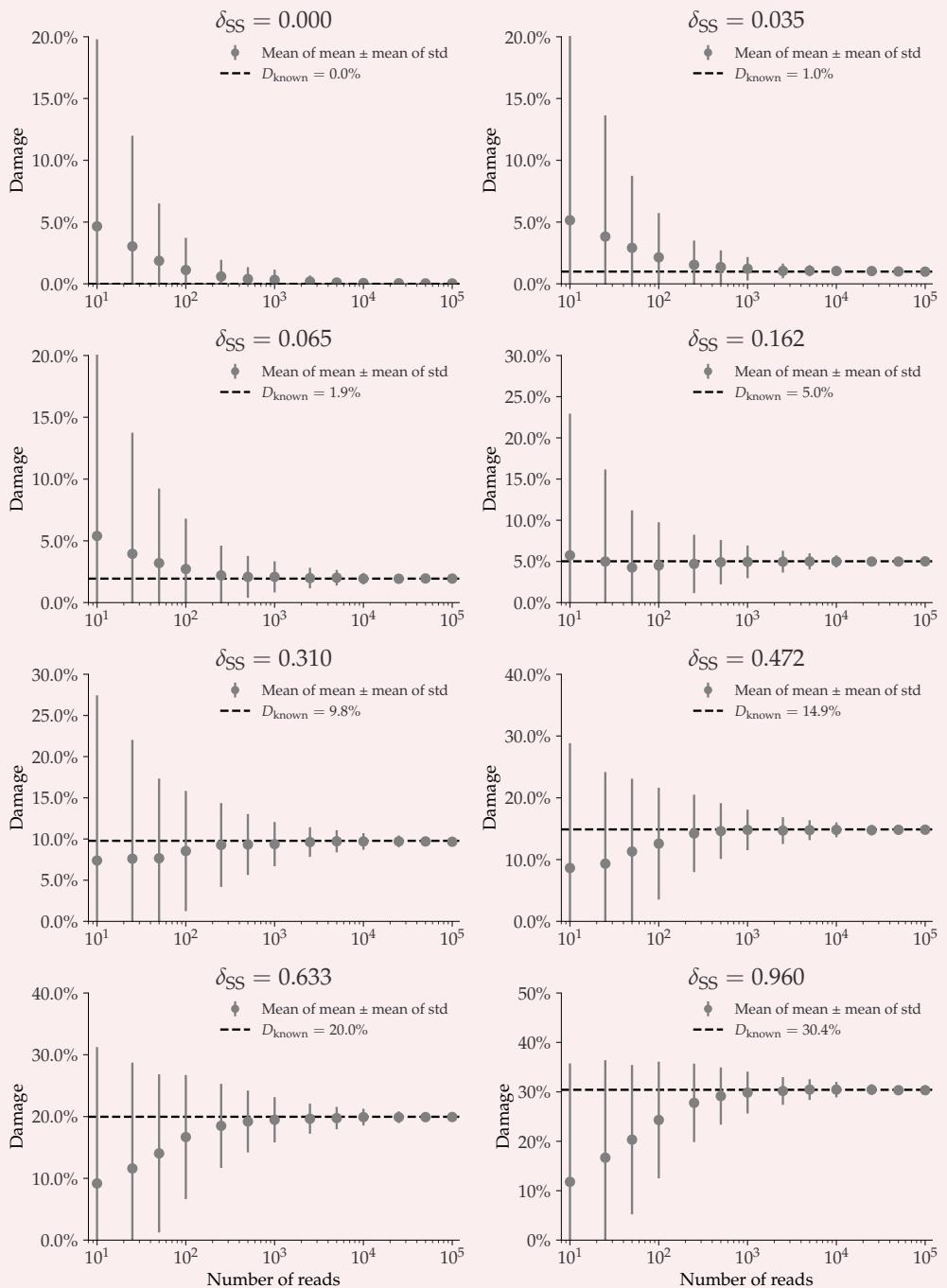
**Appendix 3—figure S3.** This plot shows the average damage as a function of the number of reads.

656

The grey points show the average of the individual means (with the average of the standard deviations as errors).

658

## GC-low



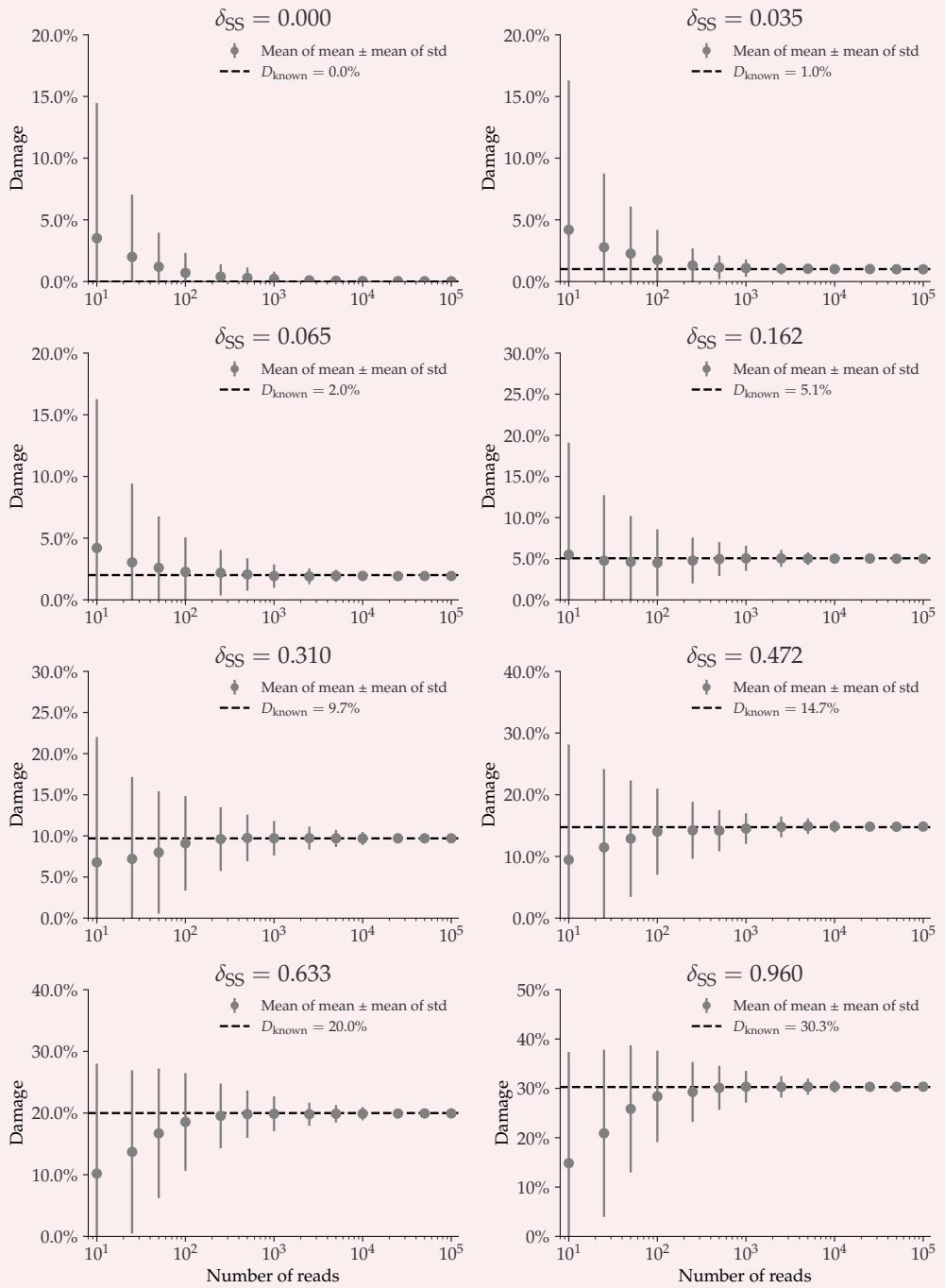
660

**Appendix 3—figure S4.** This plot shows the average damage as a function of the number of reads.

662

The grey points show the average of the individual means (with the average of the standard deviations as errors).

## GC-mid



664

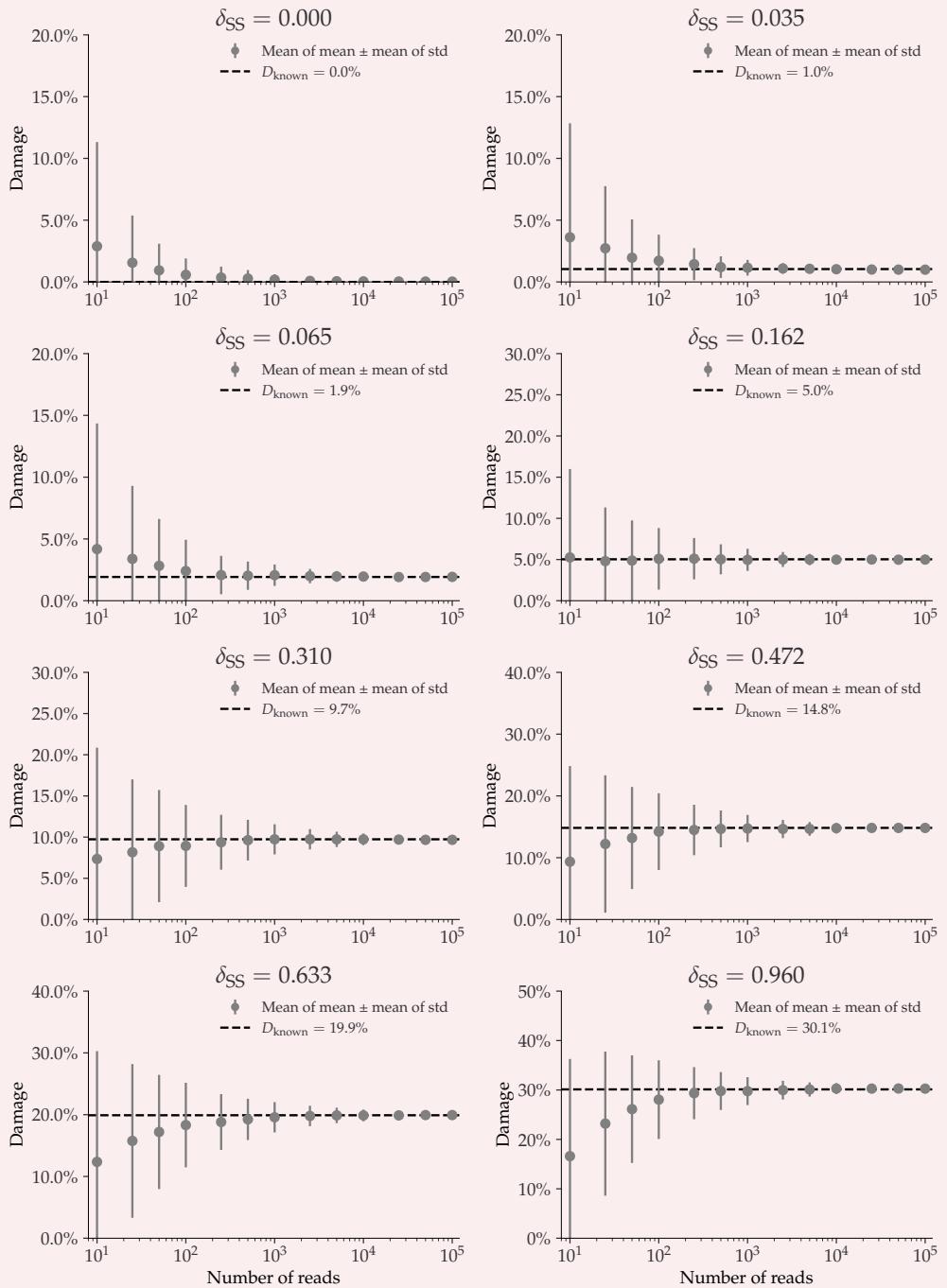
**Appendix 3—figure S5.** This plot shows the average damage as a function of the number of reads.

666

The grey points show the average of the individual means (with the average of the standard deviations as errors).

668

## GC-high



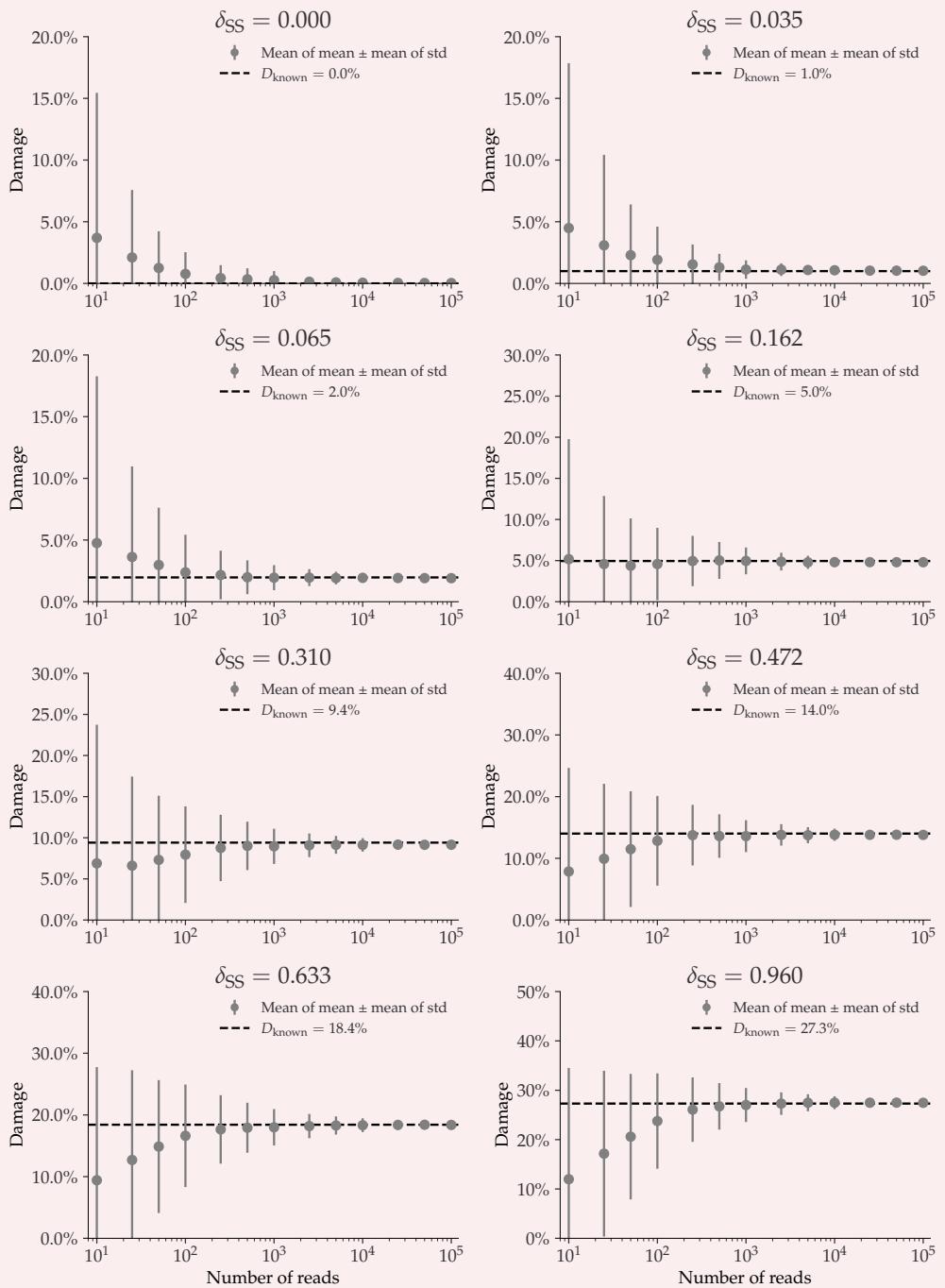
670

**Appendix 3—figure S6.** This plot shows the average damage as a function of the number of reads.

672

The grey points show the average of the individual means (with the average of the standard deviations as errors).

## Fragment Length Average: 35



674

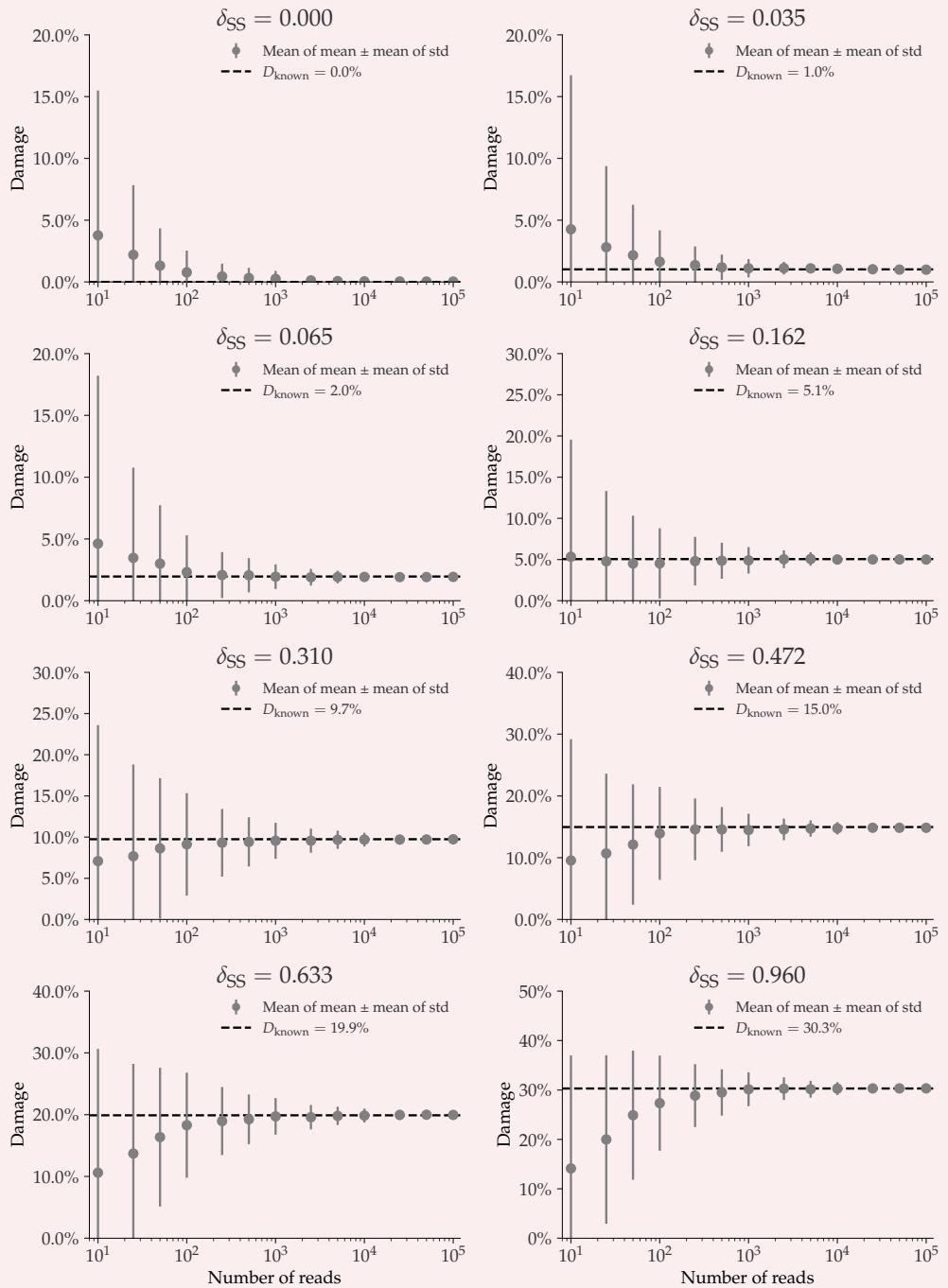
**Appendix 3—figure S7.** This plot shows the average damage as a function of the number of reads.

676

The grey points show the average of the individual means (with the average of the standard deviations as errors).

678

## Fragment Length Average: 60



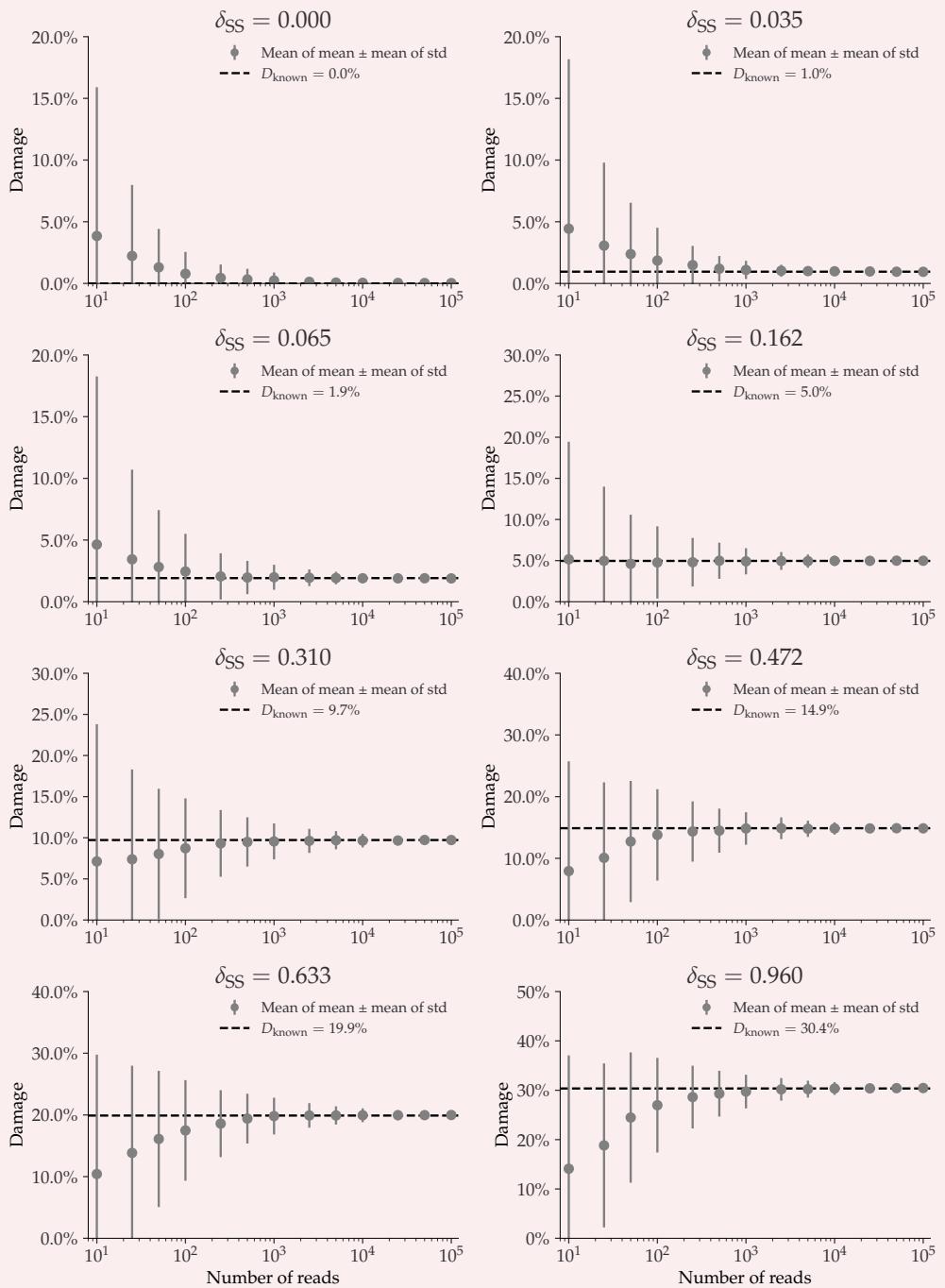
680

**Appendix 3—figure S8.** This plot shows the average damage as a function of the number of reads.

682

The grey points show the average of the individual means (with the average of the standard deviations as errors).

## Fragment Length Average: 90



684

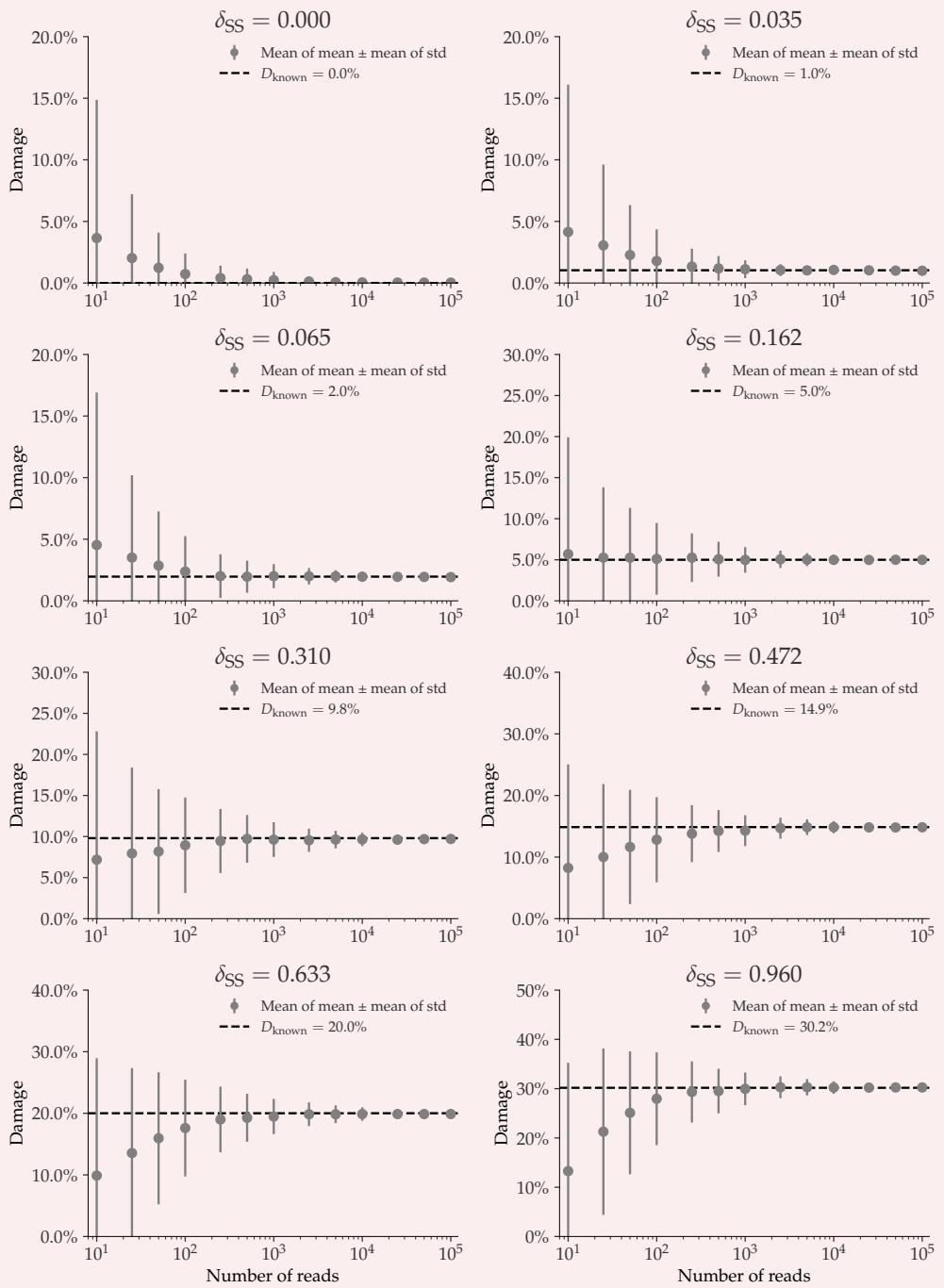
**Appendix 3—figure S9.** This plot shows the average damage as a function of the number of reads.

686

The grey points show the average of the individual means (with the average of the standard deviations as errors).

688

## Contig length: 1 000



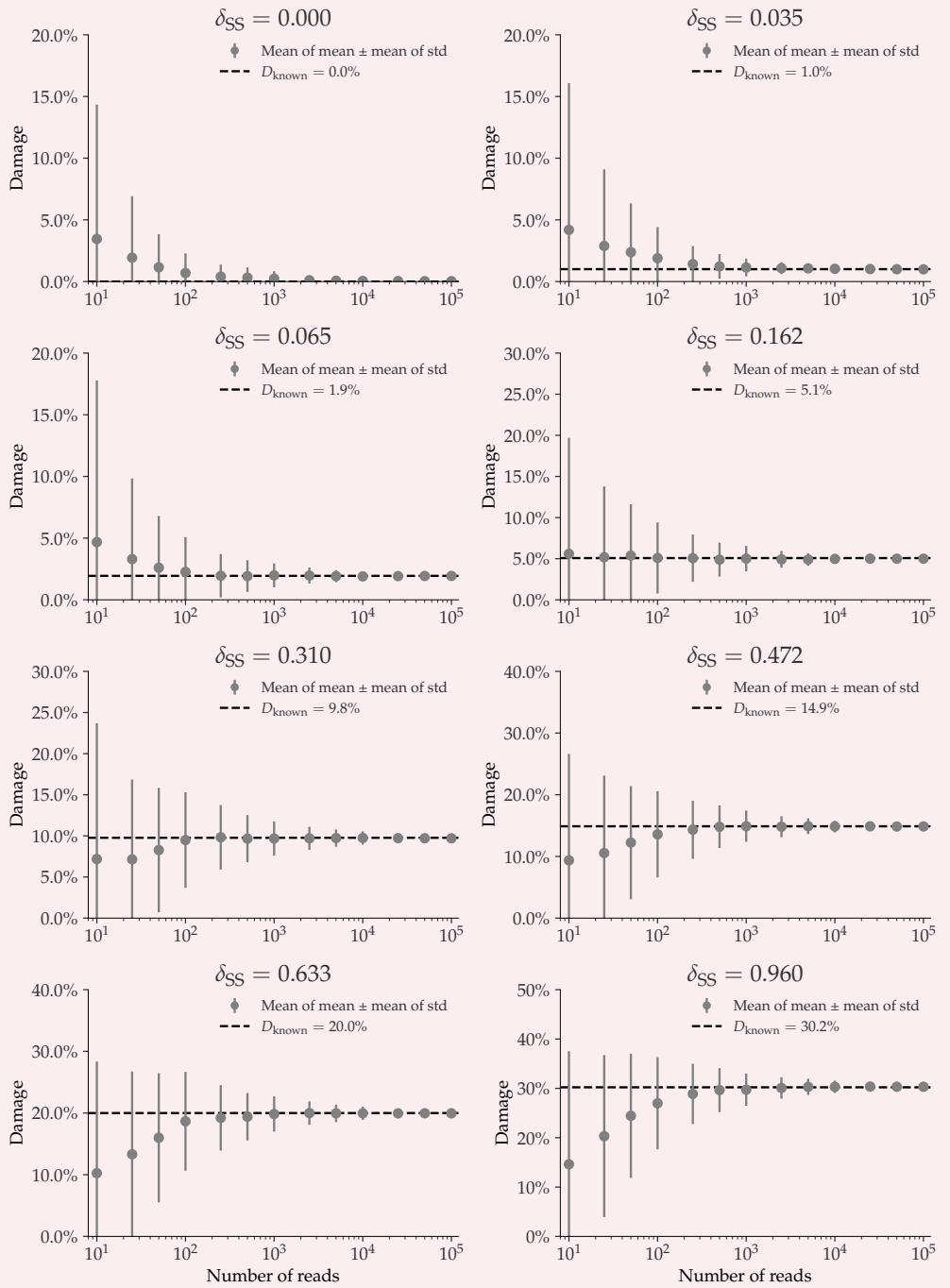
690

**Appendix 3—figure S10.** This plot shows the average damage as a function of the number of reads.

692

The grey points show the average of the individual means (with the average of the standard deviations as errors).

## Contig length: 10 000



694

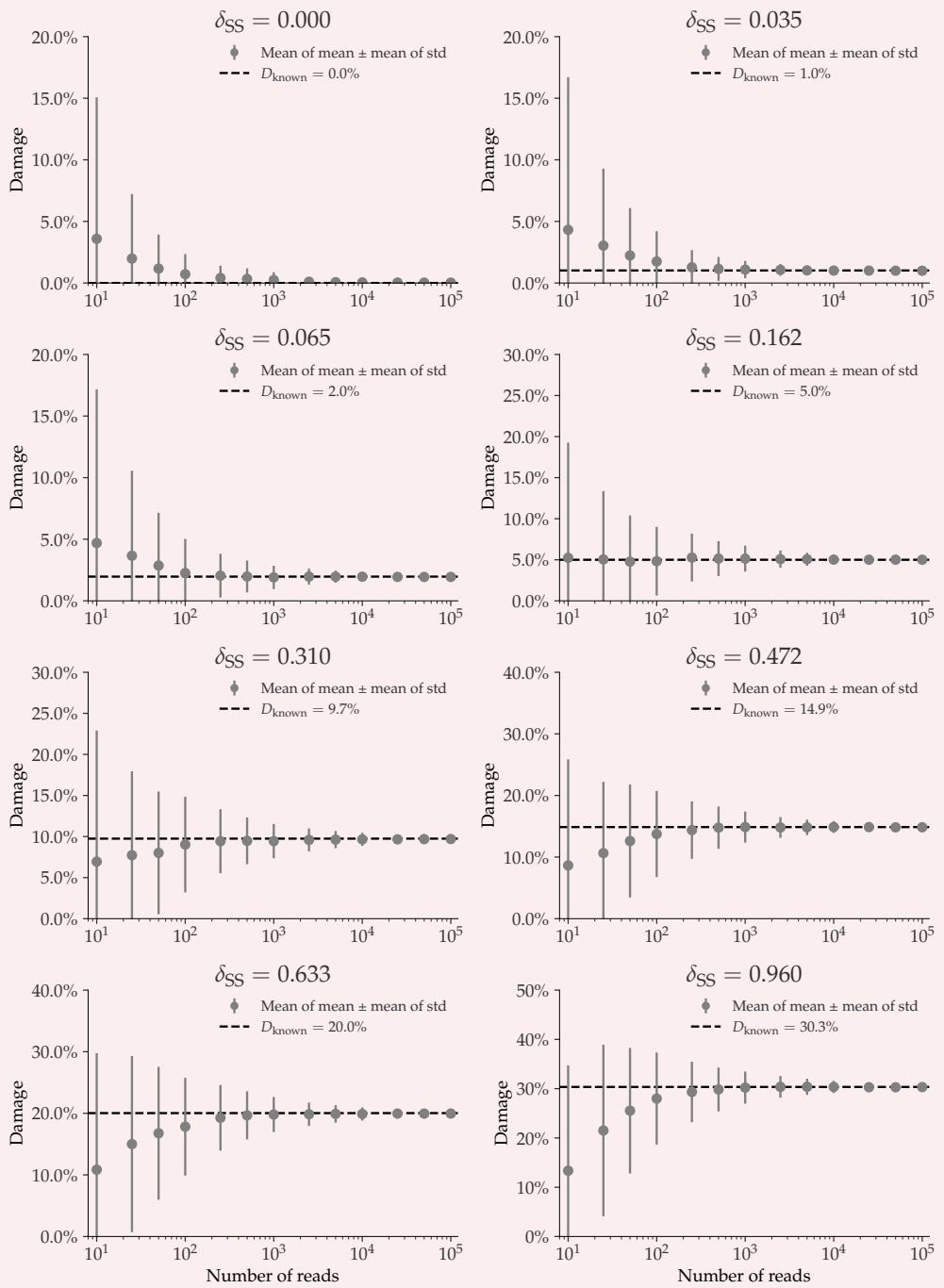
**Appendix 3—figure S11.** This plot shows the average damage as a function of the number of reads.

696

The grey points show the average of the individual means (with the average of the standard deviations as errors).

698

## Contig length: 100 000



700

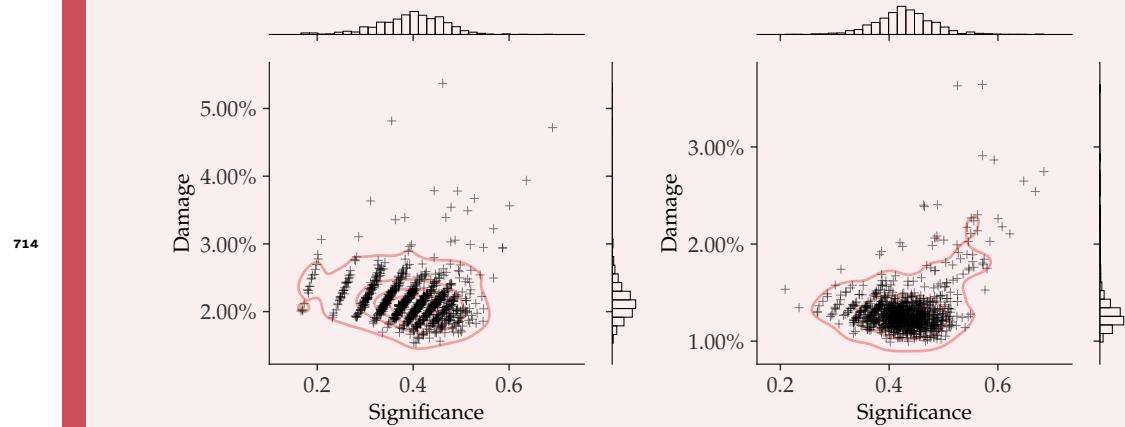
**Appendix 3—figure S12.** This plot shows the average damage as a function of the number of reads.

702

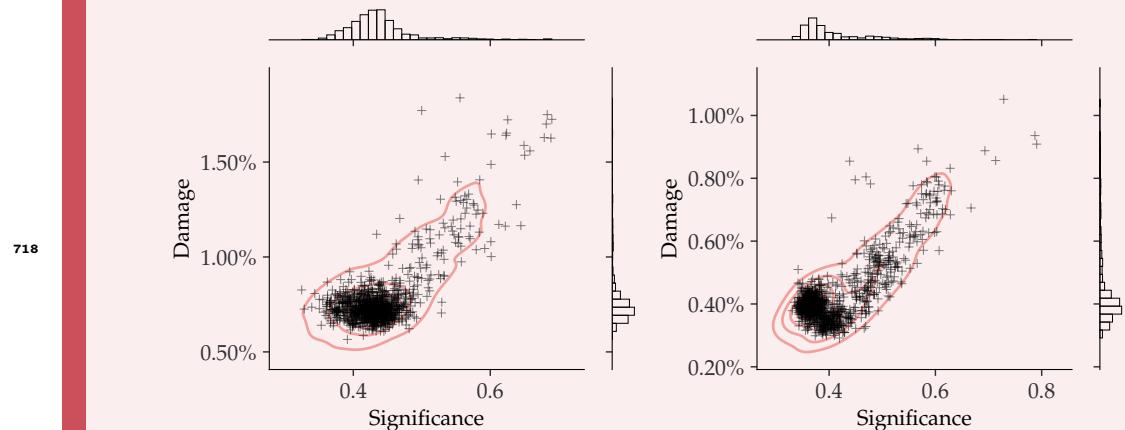
The grey points show the average of the individual means (with the average of the standard deviations as errors).

## NGSNGS SIMULATIONS – ZERO DAMAGE

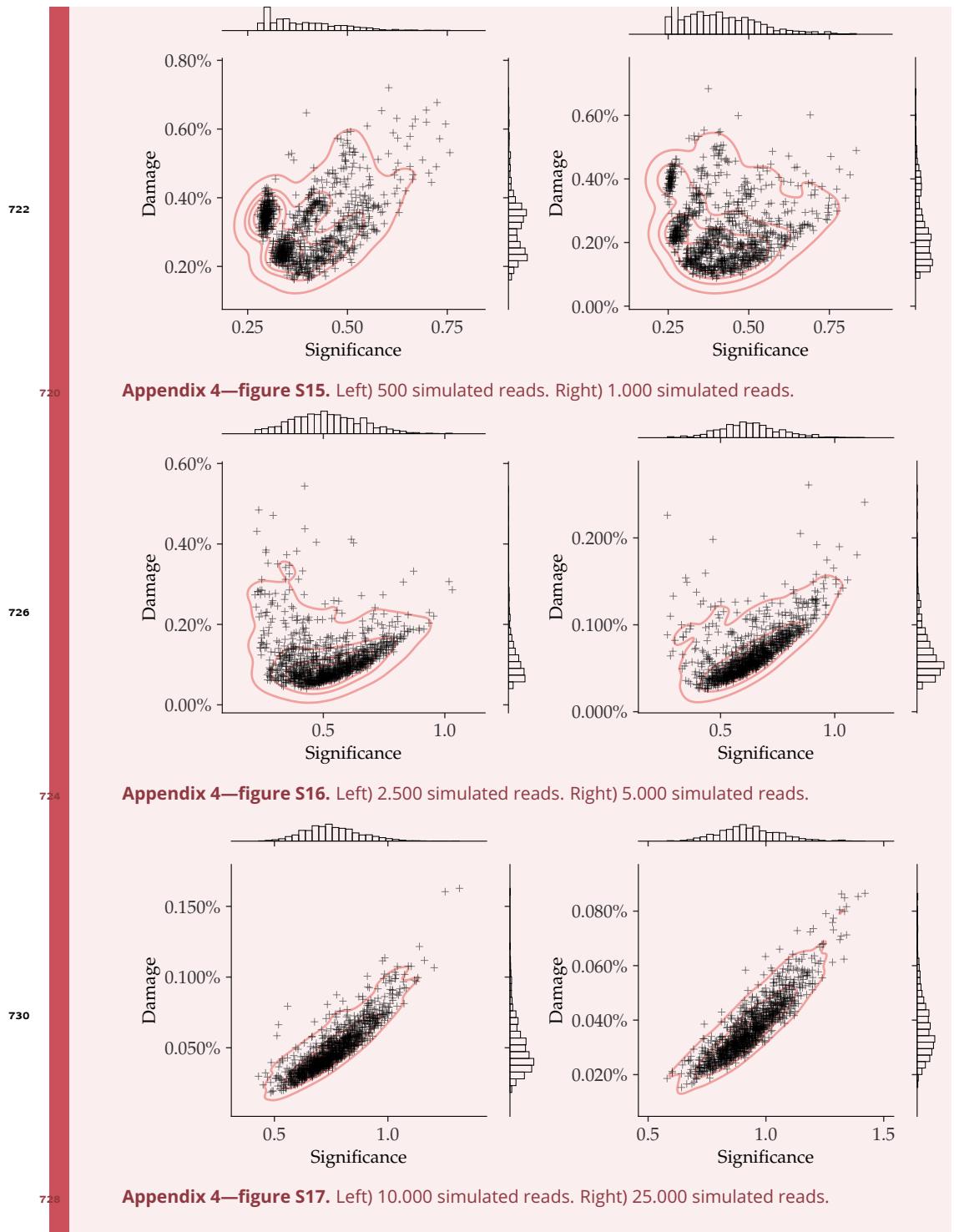
706 Damage estimates for non-damaged simulated data, each with 1000 replications. The in-  
 ferred damage is shown on the y-axis and the significance on the x-axis. Each simulation  
 708 is shown as a single cross and the red lines show the kernel density estimate (KDE) of the  
 damage estimates. The marginal distributions are shown as histograms next to the scatter  
 710 plot.

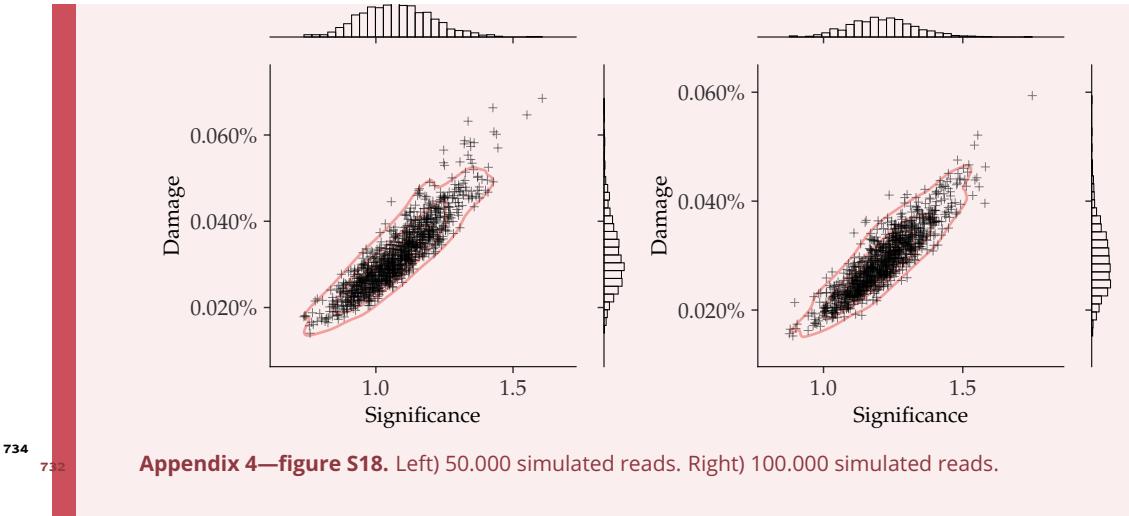


712 Appendix 4—figure S13. Left) 25 simulated reads. Right) 50 simulated reads.



718 Appendix 4—figure S14. Left) 100 simulated reads. Right) 250 simulated reads.





## MULTINOMIAL LOGISTIC REGRESSIONS

### Full Multinomial Logistic Regression models

Postmortem damages will have impacts on the NGS (next generation sequencing) reads. A common phenomenon is the calling error rates increases from nucleotide C to T due to the cytosine deamination process. Unawareness of this will lead to inaccurate inferences. Evidences show that the magnitude of such changes are related to the positions the site is within a read (the fraction of the ancient DNA). Here we present 3 slightly different ways to 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.,  $\underline{\underline{M}}(x)$ , which is a table which contains the counts of reads of different reference/read categories (in total 16) and positions on the forward/reversed strand (15 positions on each direction). Table S2 and Table S3 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	55598	
-15	12660634	4746	11753	6732	14815	9561359	5184	59633	

**Appendix 5—table S2.** 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  $M_{A-i}(x)$  or  $M_{C-i}(x)$ .

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 5—table S3.** The read counts per position given the reference nucleotides are G or T of the same human data as in Table S2. 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  $M_{T \rightarrow i}(x)$ .

760

762

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

772 relationship is as follows:

774

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

776

780

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).

778

782

### Unconditional Regression Likelihood

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

784

786

$$\begin{aligned} l_1 &= \sum_x \sum_{i,j \in \mathcal{B}} M_{i \rightarrow j}(x) \log P_{i \rightarrow j}(x) \\ &= \sum_x \left[ M(x) \log P_{T \rightarrow T}(x) + \sum_{(i,j) \neq (T,T)} M_{i \rightarrow j}(x) \log \frac{P_{i \rightarrow j}(x)}{P_{T \rightarrow T}(x)} \right], \end{aligned} \quad (12)$$

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

788

$$\log \frac{P_{i \rightarrow j}(x)}{P_{T \rightarrow T}(x)} = \sum_{n=0}^{\text{order}} \alpha_{i,j,x,n} x^n \quad (13)$$

790

Applying Equation 13 to Equation 12, we have

792

$$l_1 = \sum_x \left\{ -M(x) \log \left[ 1 + \sum_{(i,j) \neq (T,T)} \exp \left( \sum_{n=0}^{\text{order}} \alpha_{i,j,x,n} x^n \right) \right] + \sum_{(i,j) \neq (T,T)} M_{i \rightarrow j}(x) \sum_{n=0}^{\text{order}} \alpha_{i,j,x,n} x^n \right\} \quad (14)$$

794

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,

$$\frac{\partial l_1}{\partial \alpha_{i,j,x,n}} = -M(x) \frac{x^n \exp \left( \sum_{n=0}^{\text{order}} \alpha_{i,j,x,n} x^n \right)}{1 + \sum_{(i,j) \neq (T,T)} \exp \left( \sum_{n=0}^{\text{order}} \alpha_{i,j,x,n} x^n \right)} + M_{i \rightarrow j}(x) x^n. \quad (15)$$

796

### Conditional Regression Likelihood

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

$$\begin{aligned} l_2 &= \sum_{i \in \mathcal{B}} \sum_x \sum_{j \in \mathcal{B}} M_{i \rightarrow j}(x) \log P_{j|i}(x) \\ &= \sum_{i \in \mathcal{B}} \sum_x \left[ M_i(x) \log P_{T|i}(x) + \sum_{j \neq T} M_{i \rightarrow j}(x) \log \frac{P_{j|i}(x)}{P_{T|i}(x)} \right], \end{aligned} \quad (16)$$

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

$$\log \frac{P_{j|i}(x)}{P_{T|i}(x)} = \sum_{n=0}^{\text{order}} \beta_{i,j,x,n} x^n \quad (17)$$

By applying Equation 17 to Equation 16, we can obtain,

$$l_2 = \sum_{i \in B} \sum_x \left\{ -M_i(x) \log \left[ 1 + \sum_{j \neq T} \exp \left( \sum_{n=0}^{\text{order}} \beta_{i,j,x,n} x^n \right) \right] + \sum_{j \neq T} M_{i \rightarrow j}(x) \sum_{n=0}^{\text{order}} \beta_{i,j,x,n} x^n \right\} \quad (18)$$

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,

$$\frac{\partial l_2}{\partial \beta_{i,j,x,n}} = -M_i(x) \frac{x^n \exp \left( \sum_{n=0}^{\text{order}} \beta_{i,j,x,n} x^n \right)}{1 + \sum_{j \neq T} \exp \left( \sum_{n=0}^{\text{order}} \beta_{i,j,x,n} x^n \right)} + M_{i \rightarrow j}(x) x^n. \quad (19)$$

## Folded Regression

The folded regressions use the same log-likelihood functions as the full regression (i.e., Equation 14 and 18) 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].

$$\alpha_{i,j,x,n} = \alpha_{c(i),c(j),-x,n}, \quad (20)$$

$$\beta_{i,j,x,n} = \beta_{c(i),c(j),-x,n}, \quad (21)$$

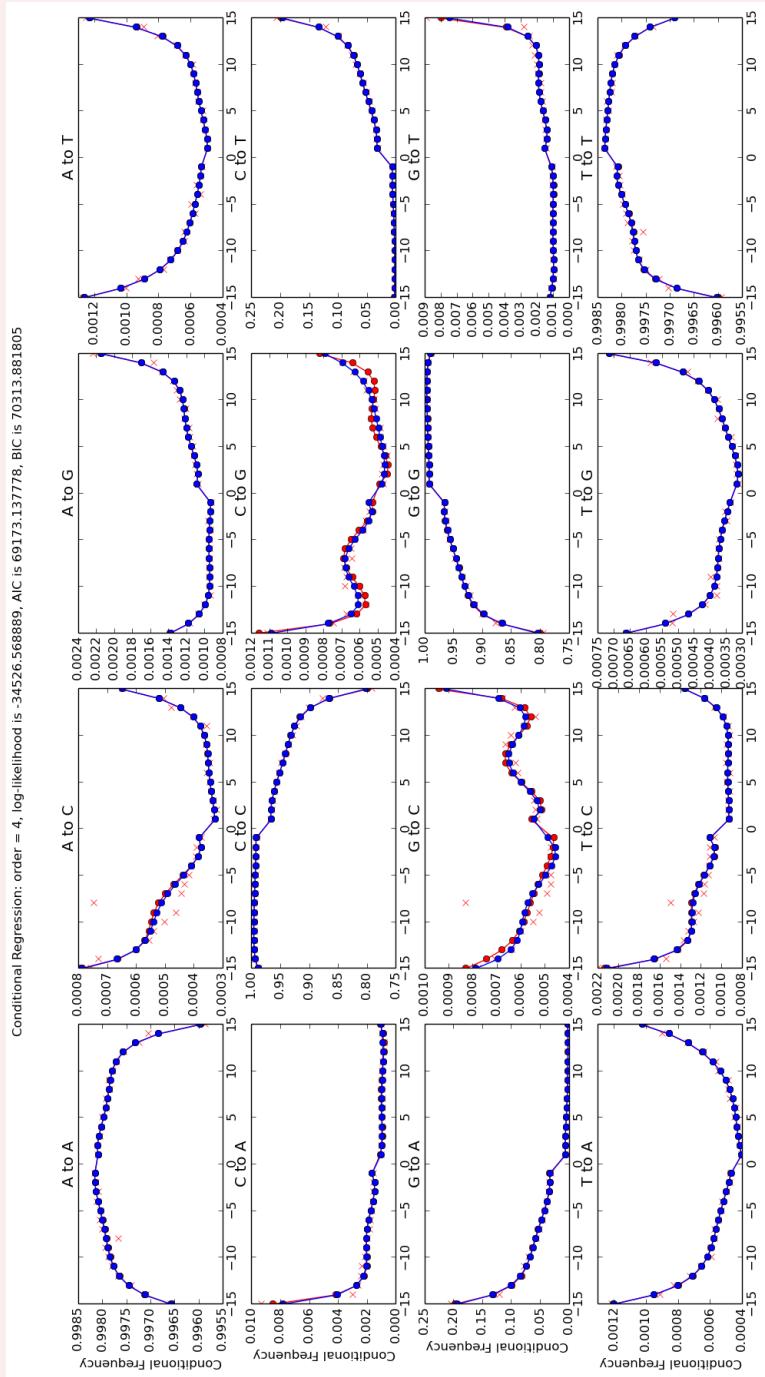
where  $c(i)$  means the complimentary nucleotide of the nucleotide  $i$ , e.g.,  $c(A) = T$  and  $c(G) = C$ .

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. S19 and Fig. S20 are the conditional regression results of the ancient and control human data correspondingly. And Fig. S21 and Fig. S22 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.

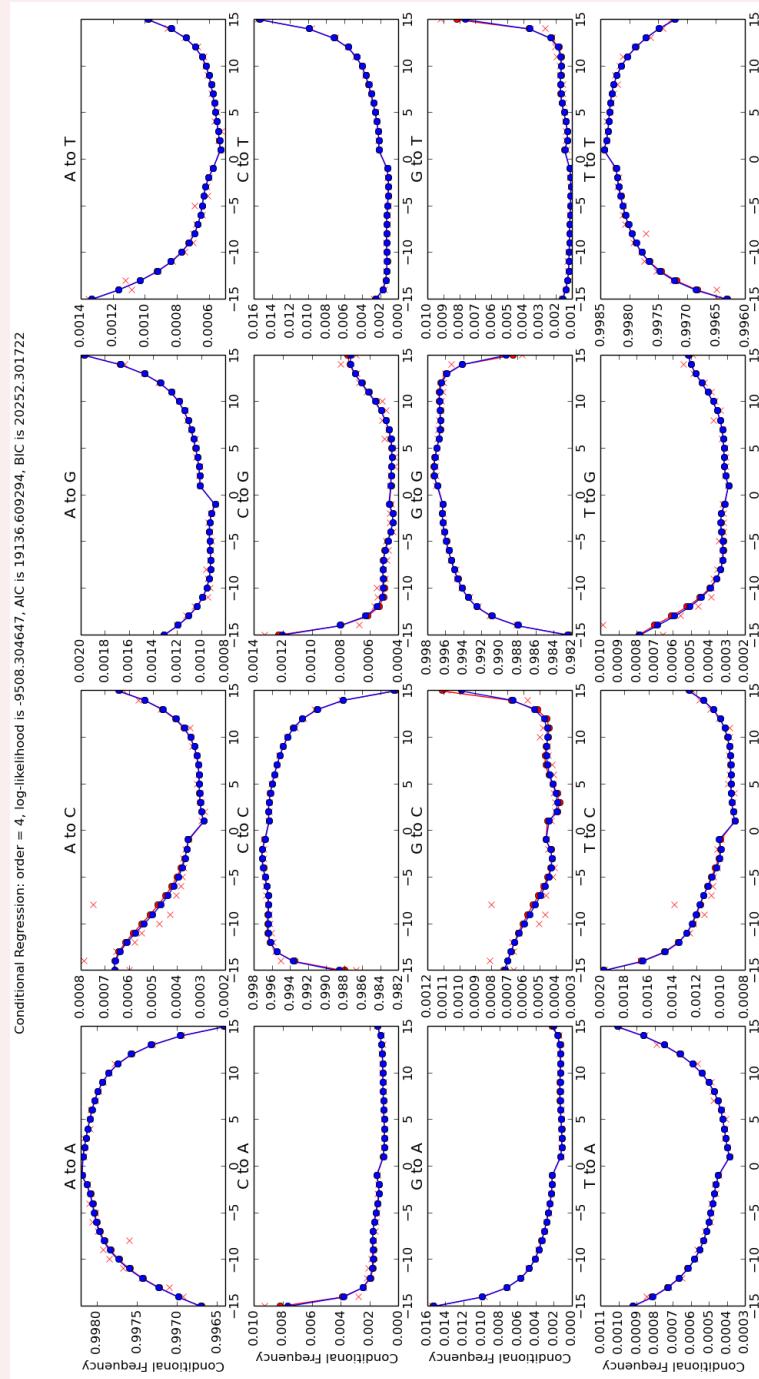


**Appendix 5—figure S19.** Conditional regression results with the order 4 of the ancient human data.

842

844

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  $1$ .



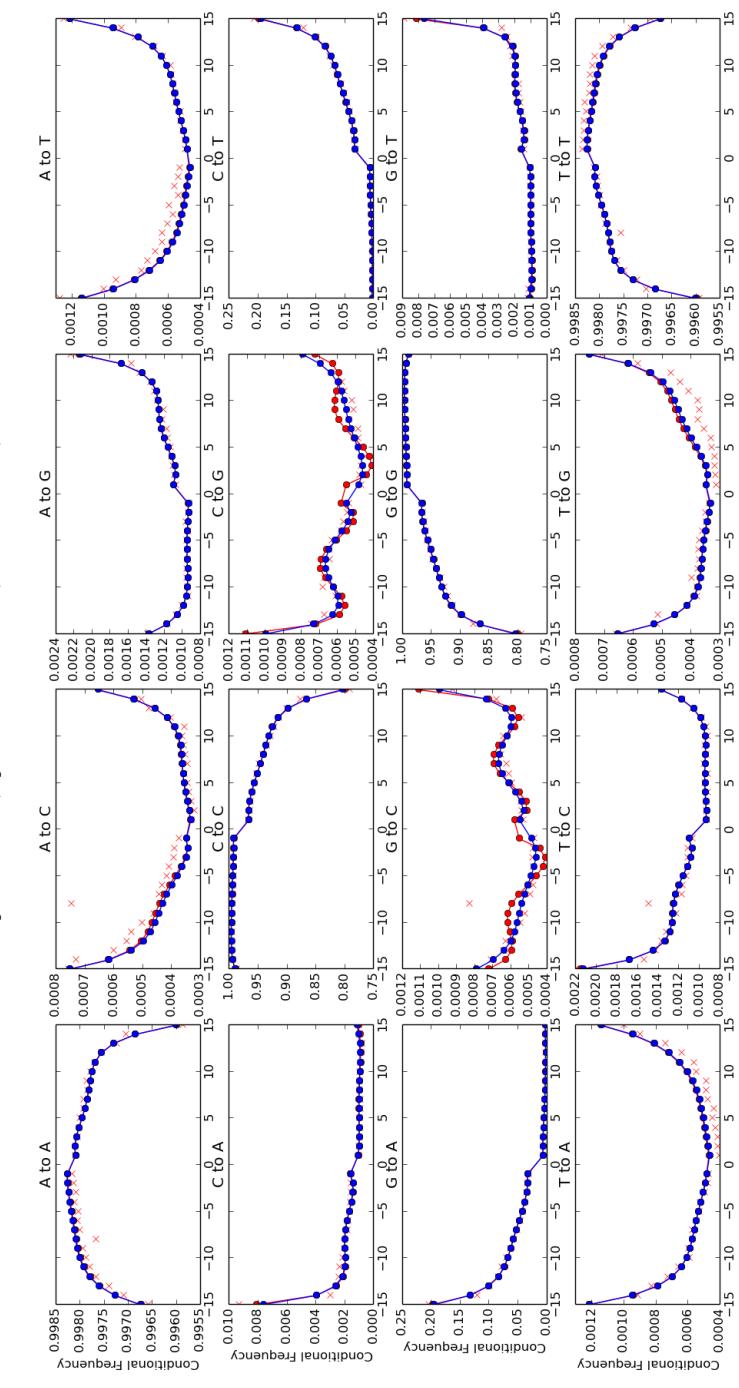
846

848

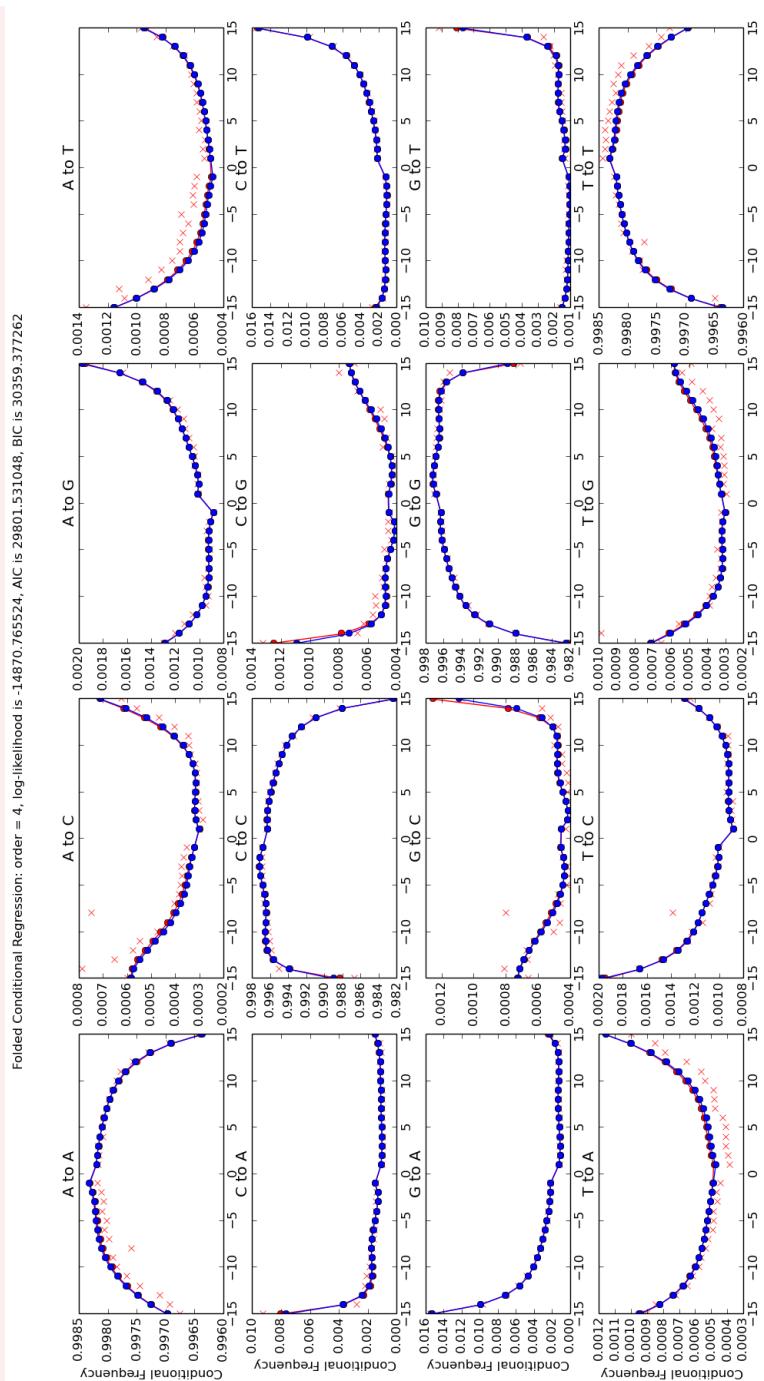
**Appendix 5—figure S20.** Conditional regression results with the order 4 of the control human data.

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  $1$ .

Folded Conditional Regression, order = 4, log-likelihood is -39581.989001, AIC is 79223.978001, BIC is 79794.350015



**Appendix 5—figure S21.** 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  $1$ .



856

**Appendix 5—figure S22.** 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  $1$ .

858

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-

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

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

## A | PMDTOOLS

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

870 According to (Skoglund et al., 2014), three nonmutually exclusive events can lead to an  
 872 observation of  $C \rightarrow T$  or  $G \rightarrow A$ , namely (i) a true biological polymorphism (occurring at  
 874 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),

$$876 D_x = C + p(1-p)^{|x|}, \quad (22)$$

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

878 The observation "Match" is defined as the case when we observe a  $C$  at a position whose  
 880 reference is also a  $C$  or a  $G$  at a position whose reference is also a  $G$ . And the observation  
 882 "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.

### 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

890

divergence contribute to a "Match" observation,

892

$$P(\text{Match} | x, \text{PMD}) = (1 - \pi)(1 - \epsilon)(1 - D_x) + (1 - \pi)\epsilon D_x + \pi\epsilon(1 - D_x), \quad (23)$$

894

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

896

$$L(\text{PMD} | S_x) = \chi_{S_x=\text{Match}} P(\text{Match} | x, \text{PMD}) + \chi_{S_x=\text{Mismatch}} P(\text{Mismatch} | x, \text{PMD}), \quad (25)$$

898

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

## Model without PMD

900

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.

902

And we have the following equations,

904

$$P(\text{Match} | x, \text{NULL}) = (1 - \pi)(1 - \epsilon) + \pi\epsilon \quad (26)$$

906

$$P(\text{Mismatch} | x, \text{NULL}) = 1 - P(\text{Match} | x, \text{NULL}) \quad (27)$$

908

$$L(\text{NULL} | S_x) = \chi_{S_x=\text{Match}} P(\text{Match} | x, \text{NULL}) + \chi_{S_x=\text{Mismatch}} P(\text{Mismatch} | x, \text{NULL}) \quad (28)$$

910

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),

$$\text{PMDS} = \log \left\{ \frac{\prod_x L(\text{PMD} | S_x)}{\prod_x L(\text{NULL} | S_x)} \right\}, \quad (29)$$

912

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