

Supporting Information

Rebolledo-Jaramillo et al. 10.1073/pnas.1409328111

SI Materials and Methods

Numt Reads Simulation. We obtained the hg19 numt genomic coordinates from the numt track of the University of California, Santa Cruz (UCSC) Genome Browser (1), extended them by 1 kb on each side and downloaded the corresponding sequences in FASTA format from the UCSC Genome Browser (<http://genome.ucsc.edu/>). We used these numt sequences to simulate 1 million paired-end reads of different lengths (50, 100, 150, 200, 250, and 300 bp) with wgsim (2). The numt-derived paired-end reads were aligned to the human reference genome hg19 (chrM replaced by the revised Cambridge Reference Sequence, NC_012920) with bwa (3). Reads were simulated and aligned with 1,000 simulations at each read length. The number of reads mapped to the mitochondrial genome was recorded.

Likelihood Function for Evaluating Point Heteroplasmies. We used the tool ANGSD (4) (www.github.com/ANGSD) to obtain a likelihood ratio test statistic (LRT) for each site in each sample. The tool implements the log likelihood function for the allele frequencies $f = (f_A, f_C, f_G, f_T)$, $f_A + f_C + f_G + f_T = 1$, for k reads covering a position, given by

$$l(\mathbf{f}) = \sum_{i=1}^k \text{Log} \left(\sum_{b \in S} \text{Pr}(x_i | B_i = b) f_b \right),$$

where $\text{Pr}(x_i | B_i = b)$ is the position-specific probability of the data read, x_i , from the i^{th} read covering the position, given that the true nucleotide in read i , B_i , equals b , $b \in S = \{A, C, G, T\}$. The tool calculates the LRT for the hypothesis $H_0: f_2 = f_3 = f_4 = 0$, where f_j is the frequency of the allele with the j^{th} highest allele frequency. Strong statistical evidence against H_0 indicates that the site is heteroplasmic. The tool was applied independently to each sample's filtered reads. The LRT statistic was transformed into a P value using the χ^2 approximation with 3 degrees of freedom.

Germ-Line Bottleneck Size and Mutation Rate Estimation Accounting for Mitotic Segregation. One can argue that the estimates of the germ-line bottleneck size do not account for the variance owing to developmental bottlenecks and mitotic mtDNA segregation, and this can lead to a lower germ-line bottleneck size. To take this into account, we subtracted the variance between two tissues for the same individual from our estimate of the genetic variance for each quartet, that is, we computed $\sigma_{\text{gen}}^2 = ((\text{MAF}_{m1} - \text{MAF}_{c1})^2 + (\text{MAF}_{m1} - \text{MAF}_{c2})^2 + (\text{MAF}_{m2} - \text{MAF}_{c1})^2 + (\text{MAF}_{m2} - \text{MAF}_{c2})^2 - 2(\text{MAF}_{m1} - \text{MAF}_{m2})^2 - 2(\text{MAF}_{c1} - \text{MAF}_{c2})^2)/4$, where m is mother, c is child, 1 is buccal tissue, and 2 is blood tissue. Applying this approach to the same 51 quartets as above (while excluding two quartets with negative N) led to the median estimate of $N = 35.0$ (with interquartile range 10.0–138.0; Fig. S15B), a value slightly higher but still very similar to our other estimates of N (discussed in the main text). Although this will have to be evaluated in further studies in more detail, our results argue for a smaller effect of developmental bottlenecks and mitotic mtDNA segregation than of the germ-line bottleneck on determining heteroplasmy levels in tissues.

Indel Analysis. The reads used to assemble were first mapped to hg19, rCRS, pUC18, and PhiX174 as described in Fig. S4. Only nonduplicate read pairs that mapped to the rCRS were used as

input to the assembler, SPAdes (5). Assemblies were curated by aligning the contigs to the rCRS, then discarding contigs that fell entirely within the alignment of a larger contig. One assembly was chosen to use as the reference for mapping all four samples in each family. The assembly was chosen by first discarding assemblies with no LASTZ hits to the rCRS, ones with over 500 contigs, and ones with an erroneous full-genome duplication. Then the assemblies were narrowed to those without contigs with non-rCRS flanks, and finally the remaining assembly with the lowest number of contigs was chosen.

The same reads used to build the assemblies were mapped back to them using the same bwa version and options as in Dataset S1, Table S19. The only difference is that we did not limit the number of mismatches to the quartet-specific reference, which otherwise could have biased against indels. The alignments for each quartet were merged into one file, with samples marked with read groups. Then, several of the filtering steps described above were applied, specifically PicardTools MarkDuplicates, selecting reads properly mapped in a pair and above a minimum read length, and removing chimeric alignments. The Naive Variant Caller was used to find indels and their read counts, using the same settings as above, except the region restriction. The resulting indels were filtered by quartet to eliminate ones that were below 0.75% frequency or 1,000× coverage in all members. Then indels were eliminated that were above 1.0 strand- or mate-bias in all samples in the quartet. The strand-bias metric used was the “SB” formula in Guo et al. (6) The mate-bias metric used the same formula, but with the first and second mate in the pair replacing forward and reverse strand reads. Indels in the low-complexity regions of 302–316 and 16,183–16,193 were excluded. Finally, indels above 1.0% frequency in any sample, with less than 1.0 strand- and mate-bias were considered real.

To estimate the background noise in indel frequency estimation, the same pipeline was applied to the PCR-amplified pGEMTeasy-derivative Z1-1 clone described below. The frequency cutoff was lowered to 0.2% to observe spurious indel calls. Minor allele frequencies at microsatellites similar to those containing our putative indel heteroplasmies were used as a baseline error rate.

Preparation of Artificial Heteroplasmy Standards from PCR Amplicons and Clonal DNA. To determine the heteroplasmy detection thresholds of Sanger sequencing and ddPCR, we created artificial mixtures with known allelic ratios. To do so we mixed mtDNA amplicons from a sample M9 (7), who is heteroplasmic at site 8,992 (C = 65.9%, T = 34.1%), with mtDNA amplicons of a sample MSu homoplasmic at that same site (C = 100%, T = 0%). Several mixtures were prepared with the frequency of T allele ranging from 0 to 34%. The resulting mixtures were analyzed with ddPCR and Sanger sequencing (Dataset S1, Tables S5 and S7 and Fig. S8 A and B).

To determine the measurement error ($\sigma_{\text{measurement}}^2$) for computing the bottleneck size we cloned the D-loop in pGEM-T-Easy (clone Z1-1). Next, we ran two independent PCR reactions amplifying the whole clone (with primers located next to each other but facing opposite directions). Each of such PCR fragments was then sequenced twice in two independent MiSeq runs. We computed the averaged (among sites) squared difference in MAFs for all sites between the runs corresponding to the two PCR reactions. Results are shown in Dataset S1, Table S20.

1. Calabrese FM, Simone D, Attimonelli M (2012) Primates and mouse Numt5 in the UCSC Genome Browser. *BMC Bioinformatics* 13(Suppl 4):S15.

2. Li H, et al.; 1000 Genome Project Data Processing Subgroup (2009) The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25(16):2078–2079.

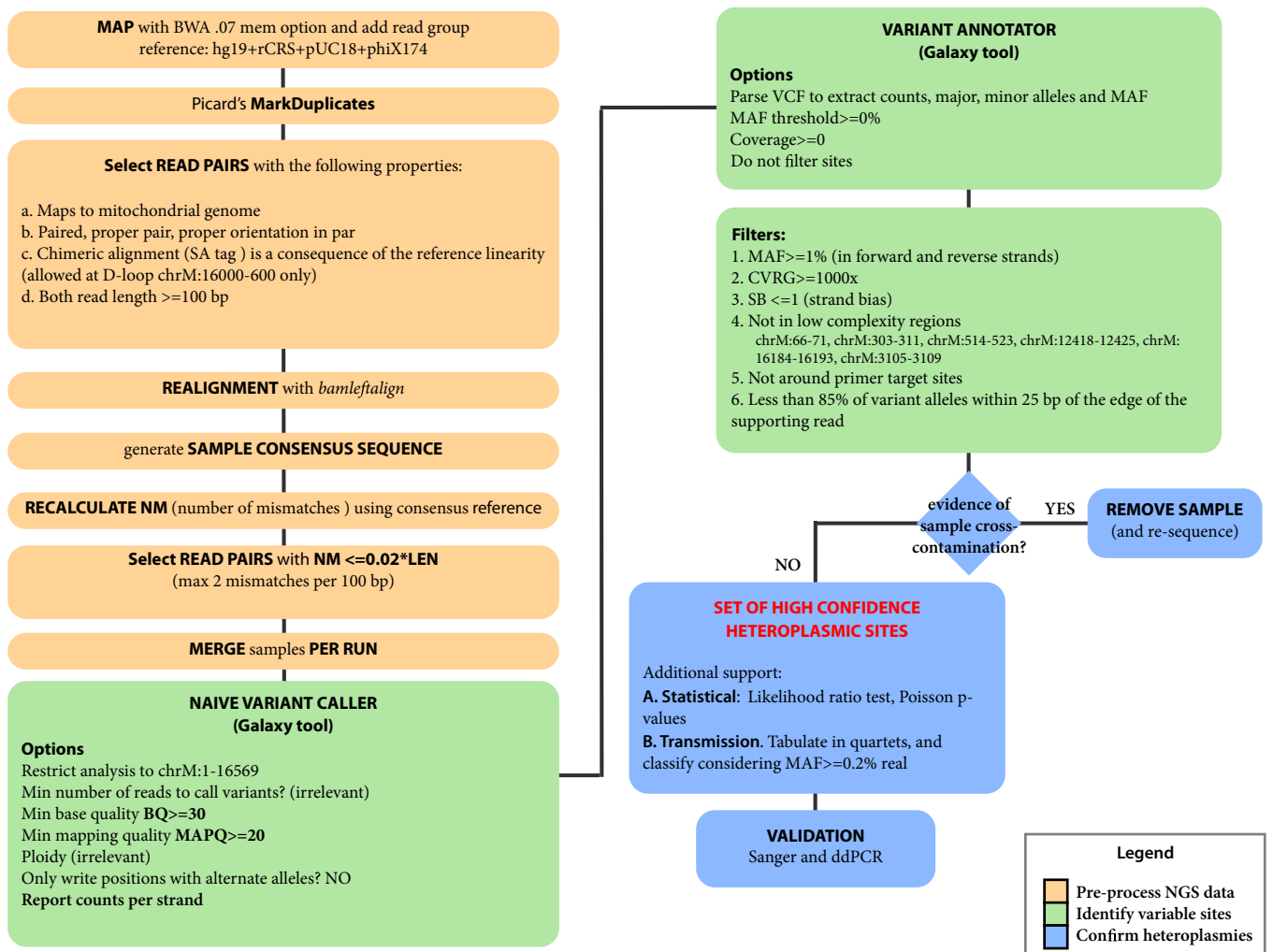


Fig. S4. Sequencing reads processing pipeline for the detection of point heteroplasmies. Reads are processed as pairs taking into consideration several filters to minimize the interference by numts, as well as alignment and sequencing artifacts. Strand bias was computed as in ref. 1. The parameters for each software is given in [Dataset S1](#), [Table S19](#).

1. Guo Y, et al. (2012) The effect of strand bias in Illumina short-read sequencing data. *BMC Genomics* 13:666.

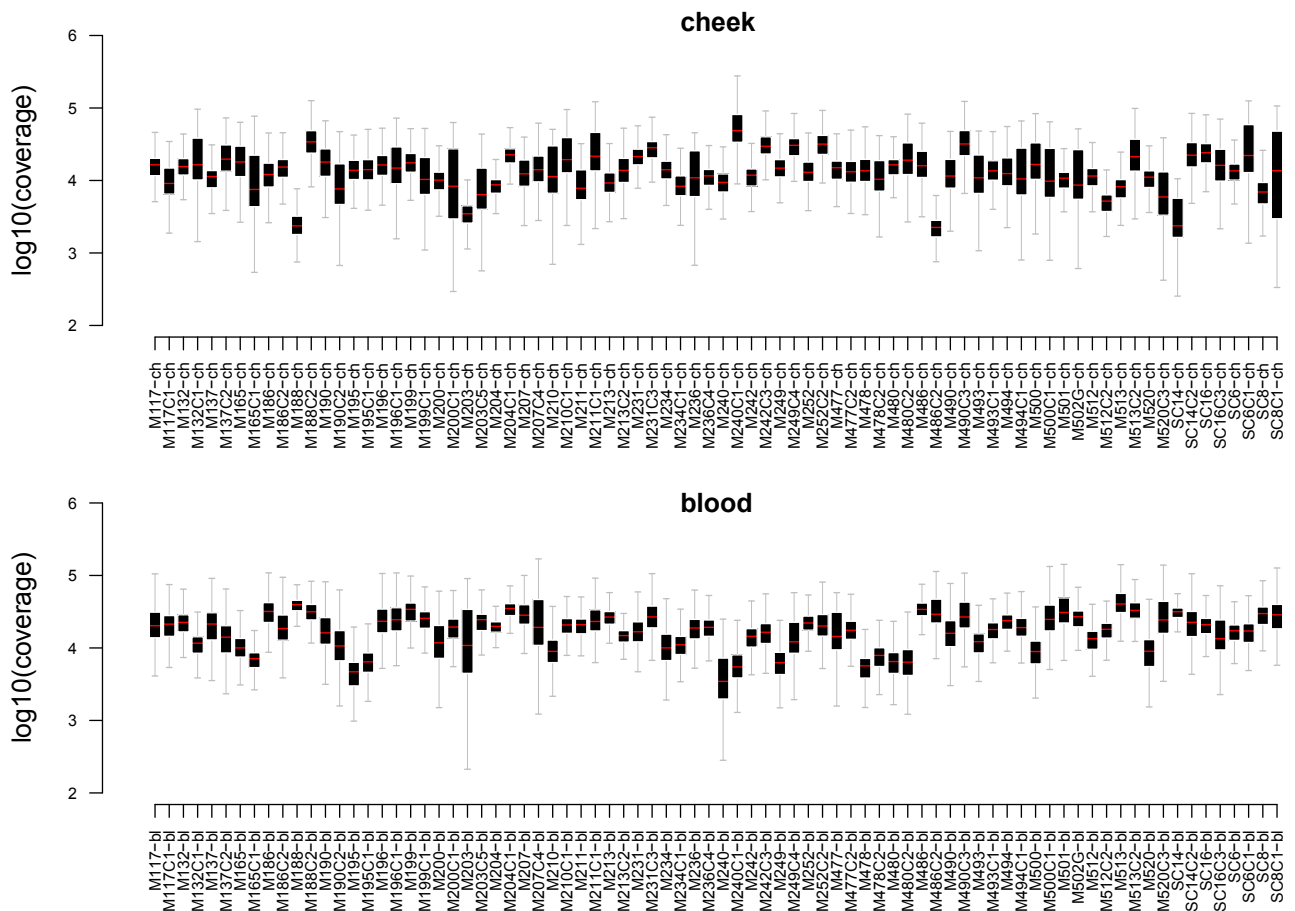


Fig. S7. Cheek and blood sample sequencing depth distribution. Distribution of the number of bases with sequencing quality ≥ 30 per site after applying several quality filters to the reads (Fig. S5). Red color indicated the median of the distribution. The y axis is on the \log_{10} scale.

expectation. (C) Determination of heteroplasmy detection limit with ddPCR. Observed MAFs were measured with ddPCR in artificial mixtures with known MAFs of (a) 0.0%, (b) 0.11%, and (c) 0.21%. The box plots show the distribution of observed MAFs in each artificial heteroplasmy standard across at least eight technical replicates. The 0% and 0.21% are well separated, and thus 0.21% is considered to be the detection limit.

1. Goto H, et al. (2011) Dynamics of mitochondrial heteroplasmy in three families investigated via a repeatable re-sequencing study. *Genome Biol* 12(6):R59.

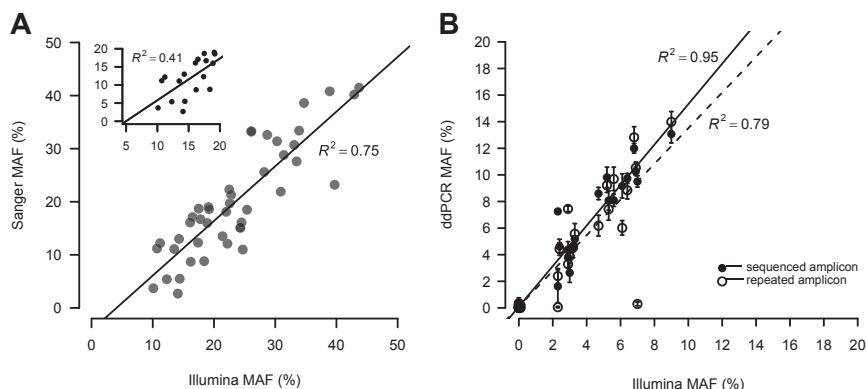


Fig. 59. Comparison of MAFs detected by three distinct experimental approaches: Sanger sequencing, Illumina sequencing (MiSeq instrument), and ddPCR. (A) The correlation between MAFs detected with Sanger (y axis) and Illumina sequencing (x axis). Black circles represent MAFs for 84 sites (21 sites in four samples; for each value the intensity was averaged across at least two sequencing runs). (Inset) A magnification of the area with MAFs between 10 and 20%. (B) The correlation between MAFs detected with ddPCR (y axis) and Illumina sequencing (x axis). For sequenced amplicons, the same exact long-range PCR was used to perform MiSeq sequencing and ddPCR (aliquots from the same tube). For repeated amplicons, ddPCR was performed on long-range PCR amplicons produced by an independent reaction (not used for MiSeq sequencing). Data are shown as mean \pm SD of four technical replicates generated from PCR amplicons per sample.

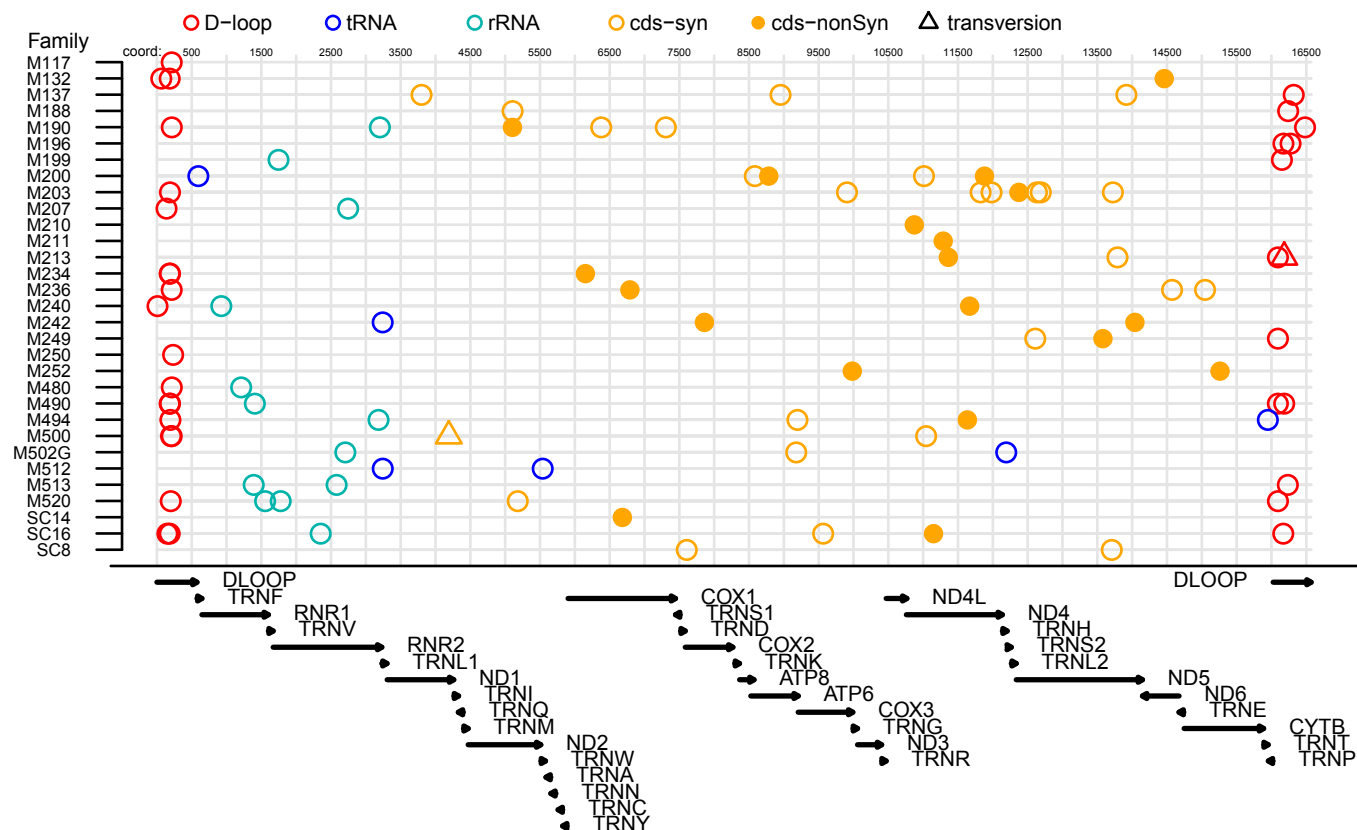


Fig. 510. The distribution of high-quality heteroplasms from 98 quartets (two tissues from a mother and two tissues from her child) along the mitochondrial genome (x axis) stratified by families (y axis). Arrows at the bottom of the graph represent strandedness of mitochondrial protein-coding and RNA genes. Circles represent transitions and triangles transversions. cds, protein-coding regions.

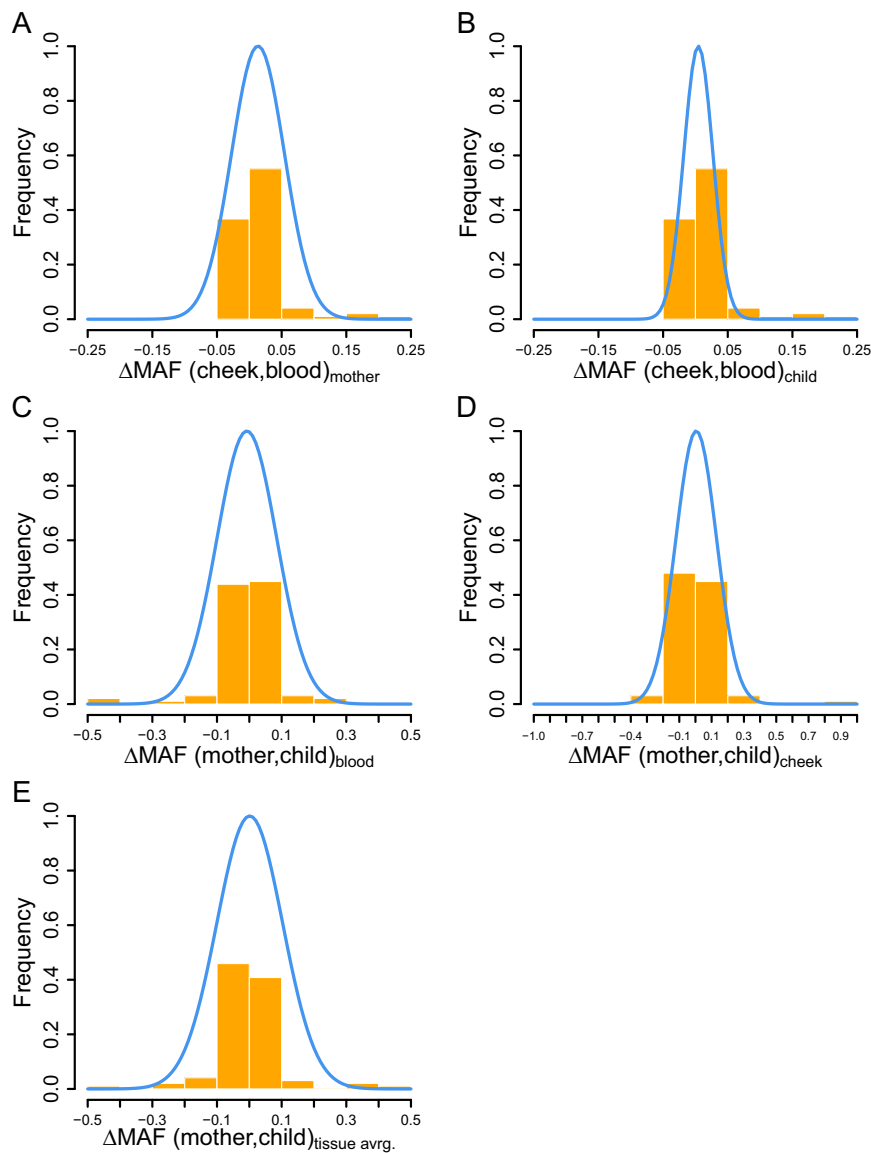


Fig. S12. The distribution of the differences in heteroplasmy allele frequencies between (A) maternal buccal tissue and blood, (B) child buccal tissue and blood, (C) mother and child blood, (D) mother and child buccal tissues, and (E) the average (between blood and buccal) maternal heteroplasmy allele frequency and the average (between blood and buccal) child heteroplasmy allele frequency. For C–E we assume maternal major allele to be ancestral for the family.

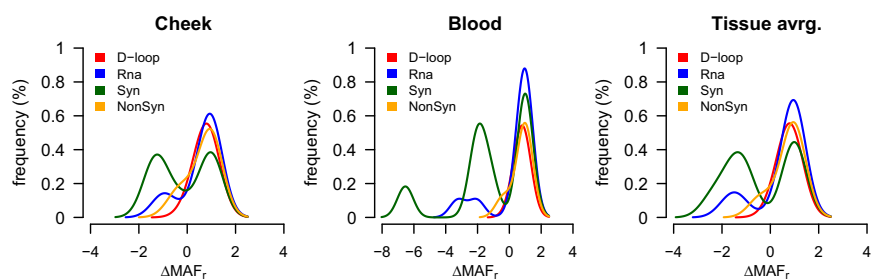


Fig. S13. Relative change (ΔMAF_r) in heteroplasmy allele frequency grouped by categories corresponding to genomic features within mitochondrial DNA: D-loop, RNA (rRNA and tRNA genes), Syn (synonymous sites), and NonSyn (nonsynonymous sites). Here maternal major allele is assumed to be ancestral for the family. $\Delta\text{MAF}_r = (\text{MAF}_{\text{mother}} - \text{MAF}_{\text{child}}) / \text{MAF}_{\text{mother}}$.

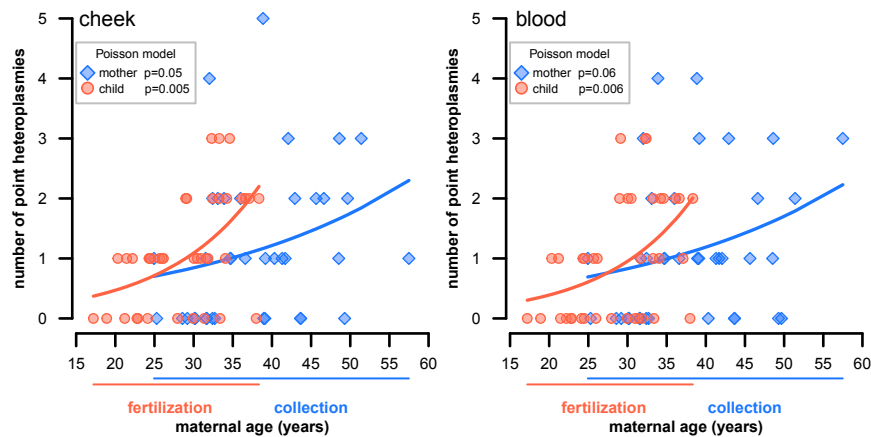


Fig. S14. Maternal age effect per tissue. Age of the mother at the time of tissue collection vs. the number of point heteroplasmies found in the corresponding tissue of the mother (blue) and age of the mother at the time of conception of the child (fertilization age) vs. the number of point heteroplasmies found in her child in the matching tissue (red). Poisson generalized linear models predicting the number of sites in the mother or child and the corresponding *P* values for the predictor (age at collection or fertilization, respectively) are indicated for each comparison.

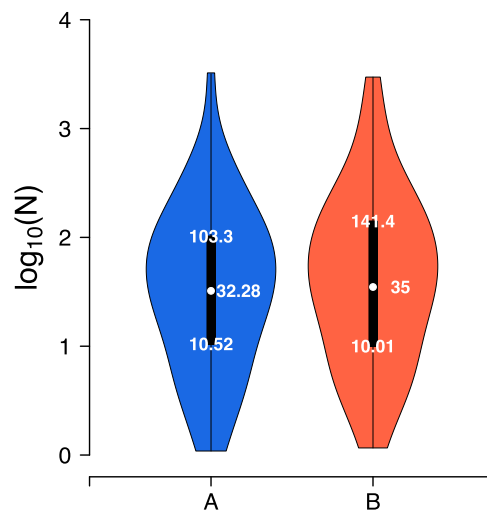


Fig. S15. Bottleneck size estimations. (A) The size of the germ-line bottleneck (*N*) was estimated according to ref. 1 using the minor alleles frequencies of the quartets in the "all," "mother," and "somatic-loss" categories only considering instances when maternal minor allele frequency was above 1% in one tissue and above 0.2% in another tissue (a total of 51 quartets, leftmost violin plot). (B) A correction for variance owing to developmental bottleneck and mitotic mtDNA segregation was applied at the rightmost plot. See text for details. The median and interquartile range are indicated inside the violin plots.

1. Millar CD, et al. (2008) Mutation and evolutionary rates in adélie penguins from the antarctic. *PLoS Genet* 4(10):e1000209.

Other Supporting Information Files

[Dataset S1 \(XLS\)](#)