

Title: Germline selection shapes the landscape of human mitochondrial DNA

Authors: Wei Wei^{1,2}, Salih Tuna^{3,4}, Michael J Keogh¹, Katherine R Smith^{5†}, Timothy J Aitman^{6,7}, Phil L Beales^{8,9}, David L Bennett¹⁰, Daniel P Gale¹¹, Maria A K Bitner-Glindzicz^{8,9,12}, Graeme C Black^{13,14}, Paul Brennan^{15,16,17}, Perry Elliott^{18,19}, Frances A Flinter^{20,21}, R Andres Floto^{22,23,24}, Henry Houlden²⁵, Melita Irving²¹, Ania Koziell^{26,27}, Eamonn R Maher^{28,29}, Hugh S Markus³⁰, Nicholas Morrell^{4,22}, William G Newman^{13,14}, Irene Roberts^{31,32,33}, John A Sayer^{16,34}, Kenneth G C Smith²², Jenny C Taylor^{33,35}, Hugh Watkins^{35,36,37}, Andrew R Webster^{38,39}, Andrew O Wilkie⁴⁰, Catherine Williamson^{41,42}, on behalf of the NIHR BioResource - Rare Diseases⁺ and the 100,000 Genomes Project - Rare Diseases Pilot⁺, Sofie Ashford^{4,43}, Christopher J Penkett^{3,4}, Kathleen E Stirrups^{3,4}, Augusto Rendon^{3,5†}, Willem H Ouwehand^{3,4,44,45,46¶}, John R Bradley^{4,22,24,29,47¶}, F Lucy Raymond^{4,28¶}, Mark Caulfield^{5,48†}, Ernest Turro^{3,4,49*}, Patrick F Chinnery^{1,2,4*}

Affiliations:

- ¹Department of Clinical Neurosciences, School of Clinical Medicine, University of Cambridge, Cambridge Biomedical Campus, Cambridge, UK.
- ²Medical Research Council Mitochondrial Biology Unit, Cambridge Biomedical Campus, Cambridge, UK.
- ³Department of Haematology, University of Cambridge, Cambridge Biomedical Campus, Cambridge, UK.
- ⁴NIHR BioResource, Cambridge University Hospitals NHS Foundation, Cambridge Biomedical Campus, Cambridge, UK.
- ⁵Genomics England, Charterhouse Square, London, UK.
- ⁶MRC Clinical Sciences Centre, Faculty of Medicine, Imperial College London, London, UK.
- ⁷Institute of Genetics and Molecular Medicine, University of Edinburgh, Edinburgh, UK.
- ⁸Genetics and Genomic Medicine Programme, UCL Great Ormond Street Institute of Child Health, London, UK.
- ⁹Great Ormond Street Hospital for Children NHS Foundation Trust, London, UK.
- ¹⁰The Nuffield Department of Clinical Neurosciences, University of Oxford, John Radcliffe Hospital, Oxford, UK.
- ¹¹UCL Centre for Nephrology, University College London, London, UK.
- ¹²University College London, London, UK.
- ¹³Evolution and Genomic Sciences, Faculty of Biology, Medicine and Health, University of Manchester, Manchester, UK.
- ¹⁴Manchester Centre for Genomic Medicine, St Mary's Hospital, Manchester Universities Foundation NHS Trust, Manchester, UK.
- ¹⁵Newcastle University, Newcastle upon Tyne, UK.

- ¹⁶Newcastle upon Tyne Hospitals NHS Foundation Trust, Newcastle upon Tyne, UK.
- ¹⁷Northern Genetics Service, Newcastle upon Tyne Hospitals NHS Foundation Trust, Newcastle upon Tyne, UK.
- ¹⁸UCL Institute of Cardiovascular Science, University College London, London, UK.
- 5 ¹⁹Barts Heart Centre, St Bartholomew's Hospital, Barts Health NHS Trust, London, UK.
- ²⁰Guy's and St Thomas' Hospital, Guy's and St Thomas' NHS Foundation Trust, London, UK.
- ²¹Clinical Genetics Department, Guy's and St Thomas NHS Foundation Trust, London, UK.
- ²²Department of Medicine, School of Clinical Medicine, University of Cambridge, Cambridge Biomedical Campus, Cambridge, UK.
- 10 ²³Royal Papworth Hospital NHS Foundation Trust, Cambridge, UK.
- ²⁴Addenbrookes Hospital, Cambridge University Hospitals NHS Foundation Trust, Cambridge, UK.
- ²⁵Department of Molecular Neuroscience, UCL Institute of Neurology, London, UK.
- ²⁶King's College London, London, UK.
- 15 ²⁷Department of Paediatric Nephrology, Evelina London Children's Hospital, Guy's & St Thomas' NHS Foundation Trust, London, UK.
- ²⁸Department of Medical Genetics, Cambridge Institute for Medical Research, University of Cambridge, Cambridge Biomedical Campus, Cambridge, UK.
- ²⁹NIHR Cambridge Biomedical Research Centre, Cambridge Biomedical Campus, Cambridge, UK.
- 20 ³⁰Stroke Research Group, Department of Clinical Neurosciences, University of Cambridge, Cambridge Biomedical Campus, Cambridge, UK.
- ³¹MRC Molecular Haematology Unit, MRC Weatherall Institute of Molecular Medicine, University of Oxford, Oxford, UK.
- 25 ³²Department of Paediatrics, Weatherall Institute of Molecular Medicine, University of Oxford, Oxford, UK.
- ³³NIHR Oxford Biomedical Research Centre, Oxford University Hospitals Trust, Oxford, UK.
- ³⁴Institute of Genetic Medicine, Newcastle University, Newcastle upon Tyne, UK.
- ³⁵Wellcome Centre for Human Genetics, University of Oxford, Oxford, UK.
- 30 ³⁶Department of Cardiovascular Medicine, Radcliffe Department of Medicine, University of Oxford, Oxford, UK.
- ³⁷Oxford University Hospitals NHS Foundation Trust, Oxford, UK.
- ³⁸Moorfields Eye Hospital NHS Foundation Trust, London, UK.
- ³⁹UCL Institute of Ophthalmology, University College London, London, UK.
- 35 ⁴⁰MRC Weatherall Institute of Molecular Medicine, University of Oxford, John Radcliffe Hospital, Oxford, UK.

⁴¹Division of Women's Health, King's College London, London, UK.

⁴²Institute of Reproductive and Developmental Biology, Surgery and Cancer, Hammersmith Hospital, Imperial College Healthcare NHS Trust, London, UK.

⁴³Department of Public Health and Primary Care, University of Cambridge, Cambridge, UK.

5 ⁴⁴Wellcome Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridge, UK.

⁴⁵NHS Blood and Transplant, Cambridge Biomedical Campus, Cambridge, UK.

⁴⁶British Heart Foundation Cambridge Centre of Excellence, University of Cambridge, Cambridge, UK.

10 ⁴⁷Department of Renal Medicine, Addenbrookes Hospital, Cambridge University Hospitals NHS Foundation Trust, Cambridge, UK.

⁴⁸William Harvey Research Institute, NIHR Biomedical Research Centre at Barts, Queen Mary University of London, London, UK.

⁴⁹MRC Biostatistics Unit, Cambridge Institute of Public Health, University of Cambridge, Cambridge, UK.

15

*Corresponding author. Email: pfc25@cam.ac.uk & et341@cam.ac.uk

⁺ a full list of the authors is included in the Supplementary Material under 'Extended Authors'

[¶] = equal contribution

20 [†] = equal contribution

= deceased

* = corresponding authors

Abstract: ~2.4% of the human mitochondrial genome (mtDNA) shows common homoplasmic genetic variation. Analyzing 12,975 whole genome sequences we show that 45.1% of individuals from 1,526 mother-offspring pairs harbor a mixed population of mtDNA (heteroplasmy), but the propensity for maternal transmission differs across the mitochondrial genome. Over one
5 generation, we observe selection both for and against variants in specific genomic regions, and previously seen variants were more likely to be transmitted. New heteroplasmy were more likely to match the nuclear genetic ancestry than the mitochondrial genome on which the mutations occurred, validating our findings in 40,325 individuals. Thus, human mtDNA at the population level is shaped by selective forces within the female germline under nuclear genetic
10 control to ensure consistency between the two independent genetic lineages.

One Sentence Summary: Human mitochondrial DNA (mtDNA) undergoes selection in the female germ line which is shaped by the nuclear genome.

Primarily inherited from the maternal line, the 16.5Kb human mitochondrial DNA (mtDNA) genome acquired mutations sequentially following the emergence of modern humans out of Africa (1-3). Pedigree and phylogenetic analyses have estimated a *de novo* mtDNA nucleotide substitution rate of $\sim 10^{-8}$ /base pair/year (4). However, from 30,506 mitochondrial genome sequences from across the globe (5), only 2.4% of nucleotides show genetic variation with frequencies greater than 1% within a population (**Fig. 1**). Although contentious (6, 7), selection could explain the non-random distribution of common variants across the mitochondrial genome in the human population.

Heteroplasmic mtDNA variants are common and maternally inherited

We analyzed high-depth mtDNA sequences from 1,526 mother-offspring pairs (mean depth in the mothers = 1,880x, range 249x-7,454x; mean depth in the offspring = 1,901x, range 259x-7,475x; mothers vs. offspring, $P=0.49$, two-sample *t*-test) (**fig. S1**). We called homoplasmic and heteroplasmic mtDNA variants from whole-blood DNA sequence data (8, 9) and filtered out heteroplasmic calls likely to be due to errors (9, 10, 11). We identified a mixed population of mtDNA (heteroplasmic variants) with a heteroplasmic variant allele frequency (VAF) >1% with high confidence in 47.8% of mothers (1,043 heteroplasmic variants at 812 sites) and 42.5% of offspring (893 heteroplasmic variants at 693 sites) (**Fig. 1, table S1 and Data S1**). In 22 individuals, where the whole genome was independently sequenced twice, the heteroplasmic mtDNA calls were 96.4% concordant (**fig. S2**) (9). As expected (12, 13), there was a small but significant positive correlation between the number of heteroplasmic variants and age of mother ($P=6.42 \times 10^{-11}$, $R^2=0.17$, CI= 0.12 - 0.23, Pearson's correlation) (**fig. S3**), with mothers having more heteroplasmic variants than offspring (mean number in the mothers = 0.68, range 0-6;

mean number in the offspring = 0.58, range 0-4; $P=0.002$, effect size = 0.68, Wilcoxon rank sum test) (**Fig. 2A**).

We defined three categories of heteroplasmic variants: (1) transmitted/inherited, if the variant was present in the mother and the offspring and was heteroplasmic in at least one of the two; (2) lost, if the heteroplasmic variant was present in the mother but not detectable in the offspring; and (3) *de novo*, if the heteroplasmic variant was present in the offspring but not detectable in the mother (**table S1**) (9). Note that very low level heteroplasmies (<1% VAF) may be missed by our sequencing and bioinformatics pipeline. Hence, “lost” and “*de novo*” variants could potentially be present at very low levels in, respectively, the offspring’s and mother’s germline. The heteroplasmic fraction (HF) of transmitted heteroplasmic variants (mean HF = 19.5%, sd = 13.9%) was significantly higher than the HF of lost variants (mean HF = 5.6%, sd = 6.3%) in the mothers ($P<2.2 \times 10^{-16}$, effect size = 4.24, Wilcoxon rank sum test); and the HF of inherited heteroplasmic variants (mean HF = 19.8%, sd = 14.1%) was significantly higher than the HF of *de novo* variants (mean HF = 6.2%, sd = 7.4%) in the offspring ($P<2.2 \times 10^{-16}$, effect size = 4.06, Wilcoxon rank sum test) (**Fig. 2B and table S1**). The HF of transmitted variants in the offspring strongly correlated with the corresponding maternal level ($P=1.52 \times 10^{-93}$, $R^2=0.79$, CI=0.75 - 0.82, Pearson's correlation) (**Fig. 2C**). In total, 477 *de novo* heteroplasmic variants were observed at >1% HF in the offspring that were not seen in the mother, in keeping with previous estimates (13). To ensure these data were not due to technical errors, we determined whether any heteroplasmic variants in the offspring were also present in their fathers. Amongst 313 father-offspring pairs, the offspring harbored 196 heteroplasmic variants with HF >1%, and only one of these was also observed in the corresponding father. This was a common population variant (population minor allele frequency (MAF) = 25.8% (5)) in the D-loop region

(m.152T>C) which was homoplasmic in the father and had an HF of 12.4% in his child. The alternate allele was not detected in the mother, suggesting this is a recurrent site of mutation or conceivably due to the paternal transmission of mtDNA.

The difference between HF in mothers and their offspring can be measured in percentage points (**Fig. 2D**) (14). This metric is limited by the difference between the HF of the mother and the boundaries 0 and 100% and the magnitude of the percentage change does not correspond with the magnitude of the fold-change in VAF. For example, a change from 50% to 55% would be given the same value as a change from 1% to 6%, even though the latter implies 6-fold increase in the proportion of mtDNA carrying the alternate allele. We therefore studied the log₂ ratio of HF between offspring and mothers after imputation of HF values below 1% to our detection threshold of 1% (subsequently termed the heteroplasmy shift (HS), **Fig. 2, E and F**) (9), which shrunk HSs towards zero only when the true HF in either the mother or the offspring was below 1%.

Overall, there was no significant difference between the number of heteroplasmic variants with a positive (n=731) and a negative (n=798) HS ($P=0.091$, binomial test). The HS distribution around zero was moderately symmetric and gave a marginal P value for asymmetry ($P=0.05$, one sample t -test) (**Fig. 2, D and E**), consistent with random segregation of mitochondria during meiosis (14, 15). All of the HSs were <6 in magnitude, corresponding to a <64-fold increase or decrease in HF across one generation, with 3 exceptions. *De novo* variants at m.57T>C (HF=99.3%), m.8993T>G (HF=82.1%) and m.14459G>A (HF=93.6%) were detected in three unrelated offspring and not present in the corresponding mothers (**figs. S4 to S6**). m.14459G>A is a non-synonymous (NS) variant in *ND6* which, on the basis of evidence from previously published pedigrees (16, 17), causes Leber hereditary optic neuropathy (LHON) and Leigh

syndrome/dystonia. m.8993T>G is a NS variant in *ATP6* (L156R), which has been observed on many independent occasions in Leigh syndrome or neurogenic ataxia with retinitis pigmentosa (18-20). Although these extreme HSs could reflect differences in the mechanism of transmission for pathogenic mtDNA mutations (21), ascertainment is a more likely explanation because

5 childhood-onset neurodegenerative diseases were recruited as part of this study (22).

Ascertainment bias is unlikely to explain the *de novo* occurrence of m.57T>C, but these findings indicate that extreme HSs at moderate HF are not typical of human populations.

As expected, the non-coding displacement (D)-loop had the highest substitution frequency (7.64×10^{-5} /base/genome/transmission) of all the regions in the mitochondrial genome (**Fig. 3A**

10 **and table S2**)(13). In total, we observed 16 out of 57 previously defined (5) pathogenic mutations in the 1,526 mother-offspring pairs (**Fig. 3B**). After excluding m.14459G>A and m.8993T>G, where the extreme HS likely reflects ascertainment bias, the mean HS for the remaining 14 pathogenic mutations was not significantly different from zero ($P=0.22$, one sample *t*-test), nor from the mean HS for the remaining 1,076 non-pathogenic variants ($P=0.11$,

15 two sample *t*-test). Thus, overall we did not see a strong signature of selection for or against pathogenic alleles, although our statistical analysis does not preclude that a subset of the observed pathogenic alleles may be under selection. Intriguingly, only three mothers carried the most common heteroplasmic pathogenic mutation m.3243A>G (23), each with a low HF (5.2%, 3.6% and 1.7%), which decreased in the corresponding offspring, to levels falling below our

20 detection threshold in two of the three offspring (3.9%, <1% and <1%). Six of the 16 pathogenic mutations were not detectable in the mothers, giving a *de novo* mutation rate for known pathogenic mutations of 393/100,000 live births (95% CI 144 – 854), which is ~3.7-fold higher than previously reported (24).

To gain insight into possible mutational mechanisms, we determined the trinucleotide mutational signature. As shown previously, C>T and T>C substitutions were the most common type of substitution in homoplasmic variants (5) and cancer somatic mtDNA mutations (25). For heteroplasmic variants, C>T and T>C substitutions were also predominant, although we also observed a small but significant excess of C>A, C>G, T>A and T>G substitutions ($P < 2.2 \times 10^{-16}$, odd ratio = 0.36, CI = 0.29 - 0.44, Fisher's exact test) (**fig. S7**). Given that the heteroplasmic variant signature was not identical to the homoplasmic variant signature, this suggests that the germline transmission shapes the mutational signatures seen in homoplasmic variants at the population level. Also of note, *de novo* mutations were more likely to involve a CpG-containing trinucleotide ($P = 3.01 \times 10^{-6}$, odd ratio = 0.50, CI = 0.38 - 0.66, Fisher's exact test) (**Fig. 3C**). Although controversial (26), this could be because methylation of NpCpG sites on the mtDNA genome predisposes to *de novo* mtDNA mutations, as seen in the nuclear genome.

Known mtDNA variants are more likely to be transmitted than novel

We then compared heteroplasmic variants which have been seen before in the general population (known) and those not previously observed (novel). Variants were considered novel if they were absent from the 1000 Genomes datasets and dbSNP and were seen in at most one individual amongst 30,506 NCBI mtDNA sequences (5). Novel heteroplasmic variants were 4.7-fold less commonly transmitted from mother to offspring than known variants ($P = 3.55 \times 10^{-13}$, odd ratio=2.60, CI=1.97 - 3.45, Fisher's exact test), and the HS for transmitted known variants was more likely to be positive ($P = 0.0002$, probability=0.40, CI=0.35 - 0.45, binomial test) (**Fig. 3, D and E**). Also, the transmitted heteroplasmic variants were more likely to affect known haplogroup-specific sites (27) compared to the lost and *de novo* heteroplasmic variants ($P = 7.86 \times 10^{-11}$, odds ratio=0.40, CI=0.30-0.53, and $P = 0.0016$, odds ratio=0.62, CI=0.46-0.84, respectively,

Fisher's exact test) (**Fig. 3F**). This suggests that factors may modulate the transmission of mtDNA heteroplasmy within the female germline over a single generation and influence the likelihood that they become established within human mtDNA populations. As heteroplasmic variants are acquired throughout life, they must be removed at transmission to offspring at a higher rate than they appear *de novo* as, otherwise, each generation would be accompanied by an expected increase in the number of heteroplasmic variants which may be deleterious (28). In keeping with this, the number of novel variants present in the mother but not transmitted (lost variants), exceeded the number of *de novo* novel variants detected in the offspring ($P=7.93 \times 10^{-7}$, probability=0.62, CI= 0.57 - 0.67, binomial test) (**Fig. 3D**), in part reflecting the accumulation of heteroplasmic variants with increasing age in the mothers (**fig. S3**).

Selection for and against heteroplasmy in different genomic regions

We analyzed different functional regions of the genome and found evidence indicating region-specific selection for or against heteroplasmic variants. The distributions of HF in the 1,526 mother-offspring pairs were significantly different between the D-loop, rRNA, tRNA, and coding regions (**Fig. 4A and table S3**). Within the coding region, the NS and synonymous (SS) variants also had different distributions ($P=2.74 \times 10^{-5}$, Kolmogorov-Smirnoff test). The NS/SS ratio was greater for the heteroplasmic variants than for the homoplasmic variants ($P=3.98 \times 10^{-24}$, odds ratio=1.91, CI=1.68 - 2.18, Fisher's exact test), and the *de novo* and lost heteroplasmic variants had a higher NS/SS than the transmitted variants (transmitted *vs de novo*: $P=0.0056$, odds ratio=1.69, CI=1.15 - 2.48; transmitted *vs lost*: $P=0.01$, odds ratio=1.57, CI=1.10 - 2.24, Fisher's exact test) (**Fig. 4B**). The heteroplasmic variants were more often in conserved sites than the homoplasmic variants ($P=3.71 \times 10^{-77}$, odds ratio=3.21, CI=2.86 - 3.60, Fisher's exact test), and the transmitted heteroplasmic variants were less conserved than the *de novo* ($P=0.0018$, odds

ratio=1.62, CI=1.19 - 2.22, Fisher's exact test) and lost ($P=9.60 \times 10^{-9}$, odds ratio=2.25, CI=1.69 - 3.03, Fisher's exact test) heteroplasmic variants (**Fig. 4C**). Also, heteroplasmic variants with a positive HS were less conserved than those with a negative HS ($P=0.03$, odds ratio=1.28, CI=1.01 - 1.61, Fisher's exact test). Variants in the rRNA genes were more likely to show a decrease in the heteroplasmy level on transmission than an increase ($P=1.00 \times 10^{-4}$, probability=0.65, CI=0.57 - 0.72, binomial test) (**Fig. 4D**), and the mean HS was significantly less than zero ($P=8.21 \times 10^{-5}$, $d=0.30$, one sample *t*-test) (**Fig. 4E**).

In order to understand the determinants of transmission of heteroplasmic variants with a reduced risk of confounding, we used multi-variable logistic regression to model the probability of transmission across all 1,526 mother-offspring pairs (9). We modeled the transmission probability of a variant as a function of its HF in the mother, the identity of the mitochondrial genome region containing it, and its known vs novel status (**Fig. 4, F to H and Fig. 3D**) (9). The probability that a heteroplasmic variant in the mother was transmitted to her offspring was associated with its HF in the mother ($P<2.2 \times 10^{-16}$, coefficient estimate=1.17, sd=0.08, logistic regression) (**Fig. 4F**). Variants in the D-loop were more likely to be transmitted ($P=0.04$, coefficient estimate=0.39, sd=0.19, logistic regression) than average and those in the rRNA were less likely to be transmitted ($P=0.0026$, coefficient estimate=-0.94, sd=0.31, logistic regression) than average (**Fig. 4G**). The novel variants were less likely to be transmitted than the known variants ($P=0.028$, coefficient estimate=0.43, sd=0.19, logistic regression) (**Fig. 3D**), even after accounting for all other covariates, including HF in the mothers.

Heteroplasmic variants in the non-coding Displacement (D-) loop

To cast light on the possible effects of selection on the non-coding D-loop, we derived a high-resolution map of heteroplasmic variants in 12,975 individuals, which included the 1,526

mother-offspring pairs (mean mtDNA genome depth = 1832x, sd=945x; mean depth of D-loop = 1569x, sd=819x) (**Fig. 5, A to C and fig. S8**) (9). We found an association between the homoplasmic allele frequency amongst 30,506 NCBI mtDNA sequences and the proportion of individuals heteroplasmic for the same allele ($P < 2.2 \times 10^{-16}$, logistic regression) (**Fig. 5, A to C**) similar to that previously observed (5). Of the 17 regions in the D-loop (**Fig. 5C** bottom - purple and orange bars), two had a significantly greater number of heteroplasmic variants than expected by chance. These regions correspond to the proposed replication fork barrier associated with the D-loop termination sequence (MT-TAS2) (29) and MT-CSB1 (MT-TAS2: $P = 4.5 \times 10^{-11}$, odds ratio=0.40, CI=0.30 - 0.54; MT-CSB1: $P = 7.0 \times 10^{-6}$, odds ratio=0.39, CI=0.24 - 0.61, Fisher's exact test vs remainder of the D-loop).

To help understand the evolution of the D-loop, we identified all the heteroplasmic variants not identified on mtDNA phylogenies across a subset of 10,210 unrelated individuals from the original dataset (9). Five of these heteroplasmic variants were shared by more than one individual and were present exclusively in people with a particular haplogroup (**Fig. 5, D and E**). One variant (m.16237A>T) was present in multiple individuals from two different branches of the phylogeny (L0a1&2 and M35b2) (**Fig. 5, D and E**). Compared to homoplasmic sequences from across the world (5), only m.299C>A was observed previously as a homoplasmic variant (in 3/30,506 individuals), each time on the R30b1 haplogroup background. This suggests that individuals we saw who were heteroplasmic for m.299C>A (**Fig. 5F**), also descended from the same maternal ancestor as the three homoplasmic individuals seen previously (5), but belonged to a closely related maternal lineage that had not yet reached fixation. These recurrent heteroplasmies contributed to the distinct trinucleotide mutational signature of the D-loop ($P = 2.3 \times 10^{-137}$, Stouffer's method for combining Fisher P values), which involves prominent non-

canonical substitutions, and is consistent with the conclusion that the homoplasmic trinucleotide mutational signature of mtDNA is shaped by germline transmission of heteroplasmic variants (**Fig. 5D and fig. S9**).

We observed an absence of low-level heteroplasmic variants in critical sites required for the initiation of mtDNA transcription and replication. These zones include several conserved sequence boxes and the light strand promoter (MT-LSP: $P=7.7 \times 10^{-18}$, odds ratio=10.12, CI=5.43 – 20.31, Fisher's exact test), which are required for mtDNA transcription and mtDNA replication (30). Certain regions with no known function (31) (eg. 16,400-16,500; **Fig. 5C**) also had a complete lack of low-level heteroplasmic variants, which suggests that an intact sequence at these regions is essential for mitochondrial function, perhaps genome propagation. The coordinates of the conserved and non-conserved regions provide a guide for functional studies of the mtDNA D-loop which has been incompletely characterized to date.

The nuclear genetic background influences the heteroplasmy landscape

Most of the ~1,500 known mitochondrial proteins are synthesized from the nuclear genome, including the majority of polypeptide subunits of the oxidative phosphorylation system, and the machinery required to replicate and transcribe the mitochondrial genome *in situ* (1). Selection for or against specific mtDNA variants must therefore occur in the context of a specific nuclear genetic background. To explore this, we identified 12,933 individuals for whom a confident mtDNA haplogroup could be predicted (**fig. S10**). We compared the haplogroup of each individual with the corresponding nuclear genetic ancestry, and identified three distinct groups of individuals: (1) a haplogroup matched group (n=11,867, 91.7%) where the mtDNA haplogroup was concordant with the nuclear ancestry; (2) a mismatched group (n=295, 2.3%) where the nuclear ancestry and mtDNA were from different human populations; and, (3) a group where the

nuclear ancestry could not be reliably determined (n=771, 6.0%) (**Fig. 6, A and B and fig. S10**). Subsequent analyses focused on the haplogroup matched and mismatched groups (9).

8,159 heteroplasmic variants at 3,854 of the 16,569 distinct sites on the mitochondrial genome were present in the matched group, and 195 heteroplasmic variants at 163 distinct sites were present in the mismatched group. The mean number of heteroplasmic variants and mean HF were not statistically different between the matched and mismatched groups (**fig. S11**). Next, we studied distinct heteroplasmic sites in the 10,179 of 12,933 individuals who were not related on the basis of their nuclear genome (9,414 in the matched group, 217 in the mismatched group and 548 in the other group). Distinct heteroplasmic sites were more likely to affect known haplogroup specific sites (27) than the rest of the mitochondrial genome ($P < 2.2 \times 10^{-16}$, Fisher's exact test), particularly within the mismatched group ($P = 0.001$, odds ratio = 1.70, CI = 1.22 - 2.36, Fisher's exact test) (**Fig. 6C**).

We extracted 2,641 haplogroup-specific variants present in only one super-population (European, Asian or African) on the world mtDNA phylogeny (27). We built a predictive model of transmission of these variants using logistic regression in 9,385 unrelated European and Asian nuclear ancestries using 2,215 European (n=940) and Asian (n=1,275) specific variants on the mtDNA phylogeny, omitting the Africans because of the diversity and small number (**figs. S10 and S12**)(9). We included the super-population and the logit population allele frequency as covariates. We also included a dummy variable indicating whether or not the variant matched the mitochondrial ancestry of the individual carrying the variant. Finally, for the matched and mismatched groups, we included a separate variable indicating whether or not the variant super-population matched the nuclear ancestry of the individual who carried the variant.

We fitted the model to 768 heteroplasmic variants in 9,179 unrelated matched individuals and 30 heteroplasmic variants in 206 unrelated mismatched individuals (9). The heteroplasmic variants in the mismatched group were significantly more likely to match the ancestry of the nuclear genetic background than the mtDNA background on which the heteroplasmy occurred ($P=2.9 \times 10^{-4}$, coefficient estimate=0.85, sd=0.24, logistic regression, **Fig. 6D and table S4**).

These findings suggest that the new mtDNA variants underwent selection to match the nuclear genome. Given the high mutation rate of the mitochondrial genome and the patterns we observed over one generation, the selective process is likely to occur within the female germline.

To independently validate this finding, we repeated this analysis with an additional 40,325 WGS recruited through the Genomics England 100,000 Genomes Rare Disease Main Programme (9). There were 36,038 individuals in a haplogroup matched group, 1,098 in a haplogroup mismatched group, and 3,124 in a group where the nuclear ancestry could not be reliably determined (**figs. S12, S13**). As before, we focused on the European and Asian specific variants observed in 23,931 unrelated European and Asian individuals. We fitted the same logistic regression model to 1,942 heteroplasmic variants in 23,277 unrelated matched individuals, and 67 heteroplasmic variants in 654 unrelated individuals where the nuclear and mtDNA had a different ancestral origin. Again, the heteroplasmic variants in the mismatched group were more likely to match the ancestry of the nuclear genetic background than the ancestral background of the mtDNA on which the heteroplasmy occurred ($P=1.33 \times 10^{-3}$, coefficient estimate=0.47, se=0.15, logistic regression, **Fig. 6D and table S4**). An inverse-weighted meta-analysis of the discovery and validation cohorts yielded a significant association across the two datasets ($P=3.3 \times 10^{-6}$, coefficient estimate=0.59, se=0.13). To gain a better understanding of underlying mechanisms we studied the gene location and HF of 97

heteroplasmic variants identified in the mismatched groups across both the discovery and validation studies. Potentially functional variants were found in the non-coding region and RNA genes, and also included 14 non-synonymous protein coding variants in the *MT-ATP*, *MT-COX*, *MT-CYB* and *MT-ND* regions (**fig. S14**). This raises the possibility that differences in oxidative phosphorylation and ATP synthesis are responsible for the association we observed.

Discussion

Several explanations have been proposed for the high substitution rate of the non-coding mtDNA D-loop, including a high intrinsic mutation rate, and/or a permissive sequence relative to the coding regions (31). Here we show that the segregation of mtDNA heteroplasmy likely plays a role in shaping D-loop population polymorphisms by a mechanism operating within the female germline. Similar findings have been seen in *Drosophila* where D-loop variants ‘selfishly’ drive segregation favoring a specific mtDNA genotype (32). These observations have implications for the development of mitochondrial transfer techniques for preventing the inheritance of severe pathogenic mtDNA mutations in humans (33, 34). After mitochondrial transfer, ~15% of human embryonic stem cell lines show reversion to the original mtDNA genotype (34-36). The reasons for this are not fully understood, but the selective propagation of D-loop heteroplasmy is a plausible explanation. Our findings implicate the nuclear genome in this process. This places greater emphasis on matching both nuclear and mtDNA backgrounds when selecting potential mitochondrial donors, in order to minimize the possibility of nuclear-mitochondrial incompatibility following mitochondrial transfer.

In cases of heteroplasmic mtDNA, one allele can be preferentially copied, or segregate to high levels in a population of daughter cells. This can lead to changes in mtDNA allele frequency during the lifetime of an individual cell, tissue or organism through genetic drift (38,

39). A high mtDNA content buffers fluctuations in allele frequency. However, if the number of copies falls to a low level, this creates a ‘genetic bottleneck’, increasing the possibility of large changes in allele frequency.

There is a ~1000-fold reduction in cellular mtDNA content during human germ cell development (40) is followed by a period of intense proliferation and migration when the germ cells migrate to form the developing gonad (41). This process is dependent on oxidative phosphorylation, and is accompanied by a massive increase in mtDNA levels (40). Under these conditions, variants that compromise mitochondrial ATP synthesis will be selected against. On the other hand, variants that promote mtDNA replication will have an advantage, potentially explaining the preferential transmission of specific D-loop variants. Subtle selective pressures will have maximal impact at this time, so the nuclear genetic influence we observed will most likely come into to play during this critical period of development. Thus, human mtDNA at the population level are influenced by selective forces acting within the female germline and modulated by the nuclear genetic background. These are apparent within one generation, and ensure consistency between these two independent genetic systems, shaping the current world mtDNA phylogeny.

Materials and Methods:

Participants, approvals and sequence acquisition

The primary data was whole genome sequencing (WGS) from 13,037 individuals in the NIHR BioResource - Rare Diseases and 100,000 Genomes Project Pilot studies (**table S5**) (22) After
5 quality control (QC, see below and (9)), 12,975 samples including 1,526 mother-offspring pairs were included in this study. For demographics see (9). Ethical approval was provided by the East of England Cambridge South national research ethics committee (REC) under reference number: 13/EE/0325. WGS was performed using the Illumina TruSeq DNA PCR-Free sample preparation kit (Illumina, Inc.) and an Illumina HiSeq 2500 sequencer, generating a mean depth of 45x
10 (range from 34x to 72x) and greater than 15x for at least 95% of the reference human genome (**fig. S8A**).

Extracting mitochondrial sequences, quality control and variant detection

WGS reads were aligned to the Genome Reference Consortium human genome build 37 (GRCh37) using Isaac Genome Alignment Software (version 01.14; Illumina, Inc.). Reads
15 aligning to the mitochondrial genome were extracted from each BAM file and analyzed using MToolBox (v1.0) (8, 9). Variant Call Files and the merged VCF were normalized with bcftools and vt (44, 45, 46), and duplicated variants were dropped with vt. The final VCF was annotated using the Variant Effect Predictor (VEP) (47). Further QC was carried out as described (9). Potential DNA cross-contamination was investigated using verifyBamID (48) in the nuclear
20 genome, and mtDNA variant calls (9).

Determining matched and mismatched groups

The pairwise relatedness and nuclear ancestry were estimated using nuclear genetic markers as described (9). MtDNA haplogroup assignment was performed using HaploGrep2 (27, 57). We then compared the mtDNA phylogenetic haplogroup with the nuclear genetic ancestry in the same individual, and identified three distinct groups of individuals as described in the text.

5

Defining novel variants

Variants were considered to be novel if absent from 1000 Genomes datasets and dbSNP and were seen in at most one individual amongst 30,506 NCBI mtDNA sequences (5).

10 MtDNA mutational spectra and signature

Mutational spectra were derived from the reference and alternative alleles as described (25, 58).

Probability of maternal mtDNA transmission

We modelled the probability of transmission of heteroplasmic variants observed in the mothers
15 using the following logistic regression model:

$$\text{logit } P(y_{ijl} = 1) = \alpha + \beta_1 \mathbf{1}_{j=1} + \beta_2 \mathbf{1}_{j=2} + \beta_3 \mathbf{1}_{j=3} + \beta_4 \mathbf{1}_{j=4} + \gamma w_{ijl} + \eta z_{ijl}$$

where $y_{ijl} = 1$ if the l th variant within mitochondrial genomic region j in mother i was transmitted and zero otherwise; $j = 0, 1, 2, 3$ or 4 denote the coding, Dloop, rRNA, tRNA and the remainder sequences, respectively; w_{ijl} is the logit of the HF of the l th variant within mitochondrial genomic region j in mother i ; and $z_{ijl} = 1$ if the l th variant within mitochondrial
20 genomic region j in mother i was observed in no individuals from the 1000 Genomes datasets, dbSNP and at most one individual amongst 30,506 NCBI mtDNAs, otherwise it was equal to zero.

Homoplasmic allele frequency in the population and heteroplasmic variants

We fitted a logistic regression model to explore the relationship between the homoplasmic allele frequency in the general population and the rate at which individuals who are not homoplasmic for the alternate allele are heteroplasmic (9).

Defining haplogroup specific variants on mtDNA phylogenetic tree

We extracted 4,476 SNVs present on mtDNA phylogenetic tree (27), then focused on SNVs either present in only one super-population (European, Asian or African), or present on two super-populations but commonly seen in one population (>1%) and not seen or extremely rare in the other population in 17,520 mtDNAs (5). This defined 2,641 haplogroup-specific variants, including 426 African variants, 1,275 Asian variants and 940 European variants.

Nuclear genome ancestry and mtDNA heteroplasmic variants

We modelled the presence or absence of a heteroplasmic variant in a particular individual using logistic regression. We considered only the 2,215 mtDNA variants from of over 4,000 haplogroup-specific variants that are present exclusively in European or Asian branches of the world mtDNA phylogeny (27), as this allows unambiguous assignation of mitochondrial ancestry to each variant. To avoid the potential for bias induced by recent shared ancestry between individuals, we considered only the 9,631 unrelated individuals in matched and mismatched groups. We fitted the following logistic regression model:

$$\text{logit } P(y_{ij} = 1) = \alpha + \beta_1 \mathbf{1}_{x_j=1} + \beta_2 \mathbf{1}_{x_j=2} + \gamma w_j + \eta \mathbf{1}_{z_i=w_i} + \omega \mathbf{1}_{x_j=w_i \cap z_i=w_i} + \psi \mathbf{1}_{x_j=w_i \cap z_i \neq w_i}$$

where $y_{ij} = 1$ if variant j is heteroplasmic in individual i , and zero otherwise; $x_j = 0, 1$ or 2 depending on whether the variant ancestry is Asian, African or European, respectively; w_j is the logit of the homoplasmic allele frequency of variant j in 30,506 NCBI samples; $z_i = 0, 1$ or 2 depending on whether the mitochondrial ancestry of individual i is Asian, African or European, respectively; and $w_i = 0, 1$ or 2 depending on whether the nuclear ancestry of individual i is Asian, African or European, respectively. The indicator variable **1** evaluates to 1 if the conditions in its subscript are met and zero otherwise.

Validation dataset

- 10 We repeated the nuclear-mtDNA ancestry analysis in 42,799 WGS from the Genomics England 100,000 Genomes Rare Disease Main Programme aligned to GRCh37 or/and hg38 using the same bioinformatics pipeline. See (9) for details.

Supplementary Materials:

Materials and Methods

Extended authors and all author addresses and affiliations

Figures S1-S15

5 Tables S1-S5

External Database S1

References (*44-59*)

References and Notes:

1. S. B. Vafai, V. K. Mootha, Mitochondrial disorders as windows into an ancient organelle. *Nature* **491**, 374-383 (2012).
2. D. C. Wallace, Mitochondrial DNA variation in human radiation and disease. *Cell* **163**,
5 33-38 (2015).
3. J. B. Stewart, P. F. Chinnery, The dynamics of mitochondrial DNA heteroplasmy: implications for human health and disease. *Nature Reviews Genetics* **16**, 530-542 (2015).
4. P. Soares *et al.*, Correcting for purifying selection: an improved human mitochondrial molecular clock. *Am J Hum Genet* **84**, 740-759 (2009).
- 10 5. W. Wei, A. Gomez-Duran, G. Hudson, P. F. Chinnery, Background sequence characteristics influence the occurrence and severity of disease-causing mtDNA mutations. *PLoS Genet* **13**, e1007126 (2017).
6. B. Cavadas *et al.*, Fine Time Scaling of Purifying Selection on Human Nonsynonymous mtDNA Mutations Based on the Worldwide Population Tree and Mother-Child Pairs.
15 *Hum Mutat* **36**, 1100-1111 (2015).
7. E. Ruiz-Pesini, D. Mishmar, M. Brandon, V. Procaccio, D. C. Wallace, Effects of purifying and adaptive selection on regional variation in human mtDNA. *Science* **303**, 223-226 (2004).
8. C. Calabrese *et al.*, MToolBox: a highly automated pipeline for heteroplasmy annotation and prioritization analysis of human mitochondrial variants in high-throughput
20 sequencing. *Bioinformatics* **30**, 3115-3117 (2014).
9. See Supplementary Methods on line
10. M. Gerstung, E. Papaemmanuil, P. J. Campbell, Subclonal variant calling with multiple samples and prior knowledge. *Bioinformatics* **30**, 1198-1204 (2014).
- 25 11. M. Gerstung *et al.*, Reliable detection of subclonal single-nucleotide variants in tumour cell populations. *Nat Commun* **3**, 811 (2012).
12. M. Wachsmuth, A. Hubner, M. Li, B. Madea, M. Stoneking, Age-Related and Heteroplasmy-Related Variation in Human mtDNA Copy Number. *PLoS Genet* **12**, e1005939 (2016).

13. B. Rebolledo-Jaramillo *et al.*, Maternal age effect and severe germ-line bottleneck in the inheritance of human mitochondrial DNA. *Proc Natl Acad Sci U S A* **111**, 15474-15479 (2014).
14. M. Li *et al.*, Transmission of human mtDNA heteroplasmy in the Genome of the
5 Netherlands families: support for a variable-size bottleneck. *Genome Res* **26**, 417-426 (2016).
15. P. F. Chinnery *et al.*, The inheritance of mitochondrial DNA heteroplasmy: random drift, selection or both? *Trends Genet* **16**, 500-505 (2000).
16. D. M. Kirby, S. G. Kahler, M. L. Freckmann, D. Reddihough, D. R. Thorburn, Leigh
10 disease caused by the mitochondrial DNA G14459A mutation in unrelated families. *Ann Neurol* **48**, 102-104 (2000).
17. J. M. Shoffner *et al.*, Leber's hereditary optic neuropathy plus dystonia is caused by a mitochondrial DNA point mutation. *Ann Neurol* **38**, 163-169 (1995).
18. I. J. Holt, A. E. Harding, R. K. Petty, J. A. Morgan-Hughes, A new mitochondrial disease
15 associated with mitochondrial DNA heteroplasmy. *Am J Hum Genet* **46**, 428-433 (1990).
19. Y. Tatuch *et al.*, Heteroplasmic mtDNA mutation (T---G) at 8993 can cause Leigh disease when the percentage of abnormal mtDNA is high. *Am J Hum Genet* **50**, 852-858 (1992).
20. J. M. Shoffner *et al.*, Lebers Hereditary Optic Neuropathy Plus Dystonia Is Caused by a
20 Mitochondrial-DNA Point Mutation. *Annals of Neurology* **38**, 163-169 (1995).
21. I. J. Wilson *et al.*, Mitochondrial DNA sequence characteristics modulate the size of the genetic bottleneck. *Hum Mol Genet* **25**, 1031-1041 (2016).
22. <http://biorxiv.org/cgi/content/short/507244v1>
23. G. S. Gorman *et al.*, Prevalence of nuclear and mtDNA mutations related to adult
25 mitochondrial disease. *Ann Neurol* **77**, 753-759 (2015).
24. H. R. Elliott, D. C. Samuels, J. A. Eden, C. L. Relton, P. F. Chinnery, Pathogenic mitochondrial DNA mutations are common in the general population. *Am J Hum Genet* **83**, 254-260 (2008).
25. Y. S. Ju *et al.*, Origins and functional consequences of somatic mitochondrial DNA
30 mutations in human cancer. *eLife* **3**, (2014).

26. Y. Matsuda, I. Hanasaki, R. Iwao, H. Yamaguchi, T. Niimi, Estimation of diffusive states from single-particle trajectory in heterogeneous medium using machine-learning methods. *Phys Chem Chem Phys* **20**, 24099-24108 (2018).
27. M. van Oven, M. Kayser, Updated comprehensive phylogenetic tree of global human mitochondrial DNA variation. *Hum Mutat* **30**, E386-394 (2009).
28. H. J. Muller, The relation of recombination to mutational advance. *Mutat Res* **1**, 2-9 (1964).
29. Y. Shi *et al.*, Mitochondrial transcription termination factor 1 directs polar replication fork pausing. *Nucleic Acids Res* **44**, 5732-5742 (2016).
30. M. Falkenberg *et al.*, Mitochondrial transcription factors B1 and B2 activate transcription of human mtDNA. *Nat Genet* **31**, 289-294 (2002).
31. T. J. Nicholls, M. Minczuk, In D-loop: 40 years of mitochondrial 7S DNA. *Exp Gerontol* **56**, 175-181 (2014).
32. H. Ma, P. H. O'Farrell, Selfish drive can trump function when animal mitochondrial genomes compete. *Nat Genet* **48**, 798-802 (2016).
33. H. Ma *et al.*, Metabolic rescue in pluripotent cells from patients with mtDNA disease. *Nature*, (2015).
34. L. A. Hyslop *et al.*, Towards clinical application of pronuclear transfer to prevent mitochondrial DNA disease. *Nature* **534(7607)**, 383-386 (2016).
35. M. Yamada *et al.*, Genetic Drift Can Compromise Mitochondrial Replacement by Nuclear Transfer in Human Oocytes. *Cell stem cell* **18**, 749-754 (2016).
36. E. Kang *et al.*, Mitochondrial replacement in human oocytes carrying pathogenic mitochondrial DNA mutations. *Nature* **540**, 270-275 (2016).
37. C. W. Birky, Relaxed and stringent genomes: why cytoplasmic genes don't obey Mendel's laws. *J Heredity* **85**, 355-365 (1994).
38. P. F. Chinnery, D. C. Samuels, Relaxed replication of mtDNA: a model with implications for the expression of disease. *Am J Hum Genet* **64**, 1158-1165 (1999).
39. P. F. Chinnery, D. C. Samuels, J. Elson, D. M. Turnbull, Accumulation of mitochondrial DNA mutations in ageing, cancer, and mitochondrial disease: is there a common mechanism? *Lancet* **360**, 1323-1325 (2002).

40. V. I. Floros *et al.*, Segregation of mitochondrial DNA heteroplasmy through a developmental genetic bottleneck in human embryos. *Nat Cell Biol*, (2018).
41. M. Ginsburg, M. H. L. Snow, A. McLaren, Primordial germ cells in the mouse embryo during gastrulation. *Development* **110**, 521-528 (1990).
- 5 42. W. J. Kent *et al.*, The human genome browser at UCSC. *Genome Research* **12**, 996-1006 (2002).
43. M. D. Hendy, M. D. Woodhams, A. Dodd, Modelling mitochondrial site polymorphisms to infer the number of segregating units and mutation rate. *Biol Lett* **5**, 397-400 (2009).
44. H. Li, A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987-2993 (2011).
- 10 45. H. Li *et al.*, The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-2079 (2009).
46. A. Tan, G. R. Abecasis, H. M. Kang, Unified representation of genetic variants. *Bioinformatics* **31**, 2202-2204 (2015).
- 15 47. W. McLaren *et al.*, The Ensembl Variant Effect Predictor. *Genome Biol* **17**, 122 (2016).
48. G. Jun *et al.*, Detecting and estimating contamination of human DNA samples in sequencing and array-based genotype data. *Am J Hum Genet* **91**, 839-848 (2012).
49. R. M. Andrews *et al.*, Reanalysis and revision of the Cambridge reference sequence for human mitochondrial DNA. *Nat Genet* **23**, 147 (1999).
- 20 50. C. Genomes Project *et al.*, A global reference for human genetic variation. *Nature* **526**, 68-74 (2015).
51. S. Purcell *et al.*, PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet* **81**, 559-575 (2007).
- 25 52. A. Manichaikul *et al.*, Robust relationship inference in genome-wide association studies. *Bioinformatics* **26**, 2867-2873 (2010).
53. X. Zheng *et al.*, A high-performance computing toolset for relatedness and principal component analysis of SNP data. *Bioinformatics* **28**, 3326-3328 (2012).
54. M. P. Conomos, A. P. Reiner, B. S. Weir, T. A. Thornton, Model-free Estimation of Recent Genetic Relatedness. *Am J Hum Genet* **98**, 127-148 (2016).
- 30

55. M. P. T. Conomos, T.; Gogarten, S. M. , in *GENESIS: GENetic ESTimation and Inference in Structured samples (GENESIS): Statistical methods for analyzing genetic data from samples with population structure and/or relatedness. R package version 2.6.1.* (2017).
56. J. Staples, D. A. Nickerson, J. E. Below, Utilizing graph theory to select the largest set of
5 unrelated individuals for genetic analysis. *Genet Epidemiol* **37**, 136-141 (2013).
57. H. Weissensteiner *et al.*, HaploGrep 2: mitochondrial haplogroup classification in the era of high-throughput sequencing. *Nucleic Acids Res* **44**, W58-63 (2016).
58. L. B. Alexandrov *et al.*, Signatures of mutational processes in human cancer. *Nature* **500**, 415-421 (2013).
- 10 59. M. Krzywinski *et al.*, Circos: an information aesthetic for comparative genomics. *Genome Res* **19**, 1639-1645 (2009).

Acknowledgments: We gratefully acknowledge the patients, families, and health care professionals involved in the NIHR BioResource – Rare Diseases and the 100,000 Genomes projects. We are grateful to N. S. Jones for his critical comments on an early draft of the manuscript. **Funding:** This study makes use of data generated by the NIHR BioResource and the Genomics England Rare Diseases pilot projects. Genotype and phenotype data of both projects are part of the 100,000 Genomes Project. The main source of funding for the BioResource and Genomics England is provided by the National Institute for Health Research of England (NIHR, <http://www.nihr.ac.uk>). This work was also made possible by funding from the UK Medical Research Council (MRC) to create the UK Clinical Genomics Datacentre. PFC is a Wellcome Trust Principal Research Fellow (101876/Z/13/Z & 212219/Z/18/Z), and a NIHR Senior Investigator, who receives support from the Medical Research Council Mitochondrial Biology Unit (MC_UP_1501/2), the Evelyn Trust, and the NIHR Biomedical Research Centre based at Cambridge University Hospitals NHS Foundation Trust and the University of Cambridge. WHO is a NIHR Senior Investigator and his laboratory receives support from the British Heart Foundation, Bristol-Myers Squibb, European Commission, MRC, NHS Blood and Transplant, Rosetrees Trust, and the NIHR Biomedical Research Centre based at Cambridge University Hospitals NHS Foundation Trust and the University of Cambridge. MC is an NIHR Senior Investigator and is funded by the NIHR Biomedical Research Centre at St Bartholomew's Hospital. JeT, JoT and SP are funded by the NIHR Biomedical Research Centre, Oxford. This work was supported in part by Wellcome Trust grant 090532/Z/09/Z. RH is funded by Wellcome Trust grants 201064/Z/16/Z, 109915/Z/15/Z, 203105/Z/16/Z, MRC UK grant MR/N025431/1, ERC grant 309548 and Newton Fund MR/N027302/1. JS is funded by MRC UK grant MR/M012212/1. AM, GA and AW are funded by the Moorfields Eye Charity. GA and AW are

funded by the RP Fighting Blindness. All Moorfields Eye Hospital and Institute of Ophthalmology authors are funded by the UCL Institute of Ophthalmology and Moorfields NIHR Biomedical Resource Centre. The Bristol NIHR Biomedical Research Centre provided infrastructure for BioResource activities in Bristol. Additional NIHR Biomedical Research Centres that contributed include Imperial College Healthcare NHS Trust BRC, Guy's and St Thomas' NHS Foundation Trust and King's College London BRC. The authors listed also represent NephroS, the UK study of Nephrotic Syndrome. AL is a British Heart Foundation Senior Basic Science Research Fellow - FS/13/48/30453. DLB, ACT, NVZ and MIM are members of the DOLORisk consortium funded by the European Commission Horizon 2020 (ID633491). ACT is a member of the International Diabetic Neuropathy Consortium, the Novo Nordisk Foundation (Ref. NNF14SA0006). DLB is a Wellcome clinical scientist (202747/Z/16/Z). ARW is supported by the NIHR-BRC of UCL Institute of Ophthalmology and Moorfields Eye Hospital. IR and EL are supported by the NIHR Translational Research Collaboration- Rare Diseases. HJB works for the Netherlands CardioVascular Research Initiative (CVON). TKB is sponsored by the NHSBT and British Society of Haematology. KGCS holds a Wellcome Investigator Award, MRC Programme Grant (number MR/L019027/1). MIM is a Wellcome Senior Investigator. Support from the Wellcome grant numbers 090532, 0938381. PHD receives funding from ICP Support. HSM receives support from BHF Programme Grant no. RG/16/4/32218. NC is partially funded by Imperial College NIHR BRC. MRW holds a NIHR award to the NIHR Imperial Clinical Research Facility at Imperial College Healthcare NHS Trust. PYWM is supported by grants from MRC UK (G1002570), Fight for Sight (1570/1571), Fight for Sight (24TP171), NIHR (IS-BRC-1215-20002). RH is a Wellcome Trust Investigator (109915/Z/15/Z), who receives support from the Wellcome Centre for

Mitochondrial Research (203105/Z/16/Z), Medical Research Council (UK) (MR/N025431/1), the European Research Council (309548), the Wellcome Trust Pathfinder Scheme (201064/Z/16/Z), the Newton Fund (UK/Turkey, MR/N027302/1) and the European Union H2020 – Research and Innovation Actions (SC1-PM-03-2017, Solve-RD). KF and CVG were supported by the Research Council of the University of Leuven (BOF KU Leuven, Belgium; OT/14/098). JSW is funded by Wellcome Trust [107469/Z/15/Z]; (ii) National Institute for Health Research (NIHR) Cardiovascular Biomedical Research Unit at Royal Brompton & Harefield NHS Foundation Trust and Imperial College London. GA is funded by NIHR-Biomedical Research Centre at Moorfields Eye Hospital and UCL Institute of Ophthalmology, Fight for Sight (UK) Early Career Investigator Award, Moorfields Eye Hospital Special Trustees, Moorfields Eye Charity, Foundation Fighting Blindness (USA) and Retinitis Pigmentosa Fighting Blindness. MCS holds a MRC Clinical Research Training Fellowship, grant ref MR/R002363/1. MAKu holds a NIHR Research Professorship NIHR-RP-2016-07-019 and Wellcome Intermediate Fellowship 098524/Z/12/A. JWhi is a recipient of a Cancer Research UK Cambridge Cancer Centre Clinical Research Training Fellowship. AJM has received funding from a Medical Research Council Senior Clinical Fellowship (MR/L006340/1). DPG is funded by the MRC, Kidney Research UK and St Peters Trust for Kidney, Bladder and Prostate Research. SAJ is funded by Kids Kidney Research. CL received funding from a MRC Clinical Research Training Fellowship (MR/J011711/1). KD is a HSST trainee supported by Health Education England. CHad was funded through a PhD Fellowship by the NIHR Translational Research Collaboration Rare Diseases. MJD receives funding from Wellcome Trust (WT098519MA). KJM is supported by the Northern Counties Kidney Research Fund. ELM a proportion of my work was undertaken at University College London Hospitals/University

College London, which received a proportion of funding from the Department of Health's National Institute for Health Research Biomedical Research Centres funding scheme. KCG is a holder of NIHR – BRC funding. This research was partly funded by the NIHR Great Ormond Street Hospital Biomedical Research Centre. The views expressed are those of the author(s) and

not necessarily those of the NHS, the NIHR or the Department of Health. **Author contributions:**

Study design: PFC, FLR, MC, WHO, E., WW. Data analysis: ET, WW, ST, MJK. Writing: PFC, ET, WW. Experimental and analytical supervision: PFC and ET, Project Supervision: PFC, and ET. The remaining authors contributed to the recruitment of participants, sample logistics and initial data preparation. **Competing interests:** MIM: serves on advisory panels for Pfizer,

NovoNordisk, Zoe Global; has received honoraria from Pfizer, NovoNordisk and Eli Lilly; has stock options in Zoe Global; has received research funding from Abbvie, Astra Zeneca,

Boehringer Ingelheim, Eli Lilly, Janssen, Merck, NovoNordisk, Pfizer, Roche, Sanofi Aventis, Servier, Takeda. TJA has received consultancy payments from AstraZeneca within the last 5

years and has received speaker honoraria from Illumina. KJM previously received funding for

research and currently on the scientific advisory board of Gemini Therapeutics, Boston, USA.

MCS received travel and accommodation fees from NovoNordisk. DML serves on advisory

boards for Agios, Novartis and Cerus. AMK had no competing interests at the time of the study,

since the study has received an educational grant from CSL Behring to attend the ISTH meeting in Berlin in 2017. CVG is holder of the Bayer and Norbert Heimbürger (CSL Behring) Chairs.

Data and materials availability: Heteroplasmy data for the mother-child pairs is provided in the external database S1 (Data S1). Whole genome sequence data from the NIHR BioResource – Rare Diseases project can be found in the European Genome-phenome Archive (EGA) at the EMBL European Bioinformatics Institute (BPD: EGAD00001004519, CSVD:

EGAD00001004513, HCM: EGAD00001004514, ICP: EGAD00001004515, IRD:
EGAD00001004520, MPMT: EGAD00001004521, NDD: EGAD00001004522, NPD:
EGAD00001004516, PAH: EGAD00001004525, PID: EGAD00001004523, PMG:
EGAD00001004517, SMD: EGAD00001004524, SRNS: EGAD00001004518, See **Table S5** for

5 the disease abbreviations). Whole genome sequence data from the UK Biobank samples are
available through a data release process overseen by UK Biobank

(<https://www.ukbiobank.ac.uk/>). Whole genome sequence data from the participants enrolled in
100,000 Genomes Project can be accessed via Genomics England Limited following the
procedure outlined at: <https://www.genomicsengland.co.uk/about-gecip/joining-research->

10 [community/](https://www.genomicsengland.co.uk/about-gecip/joining-research-community/)

Figures

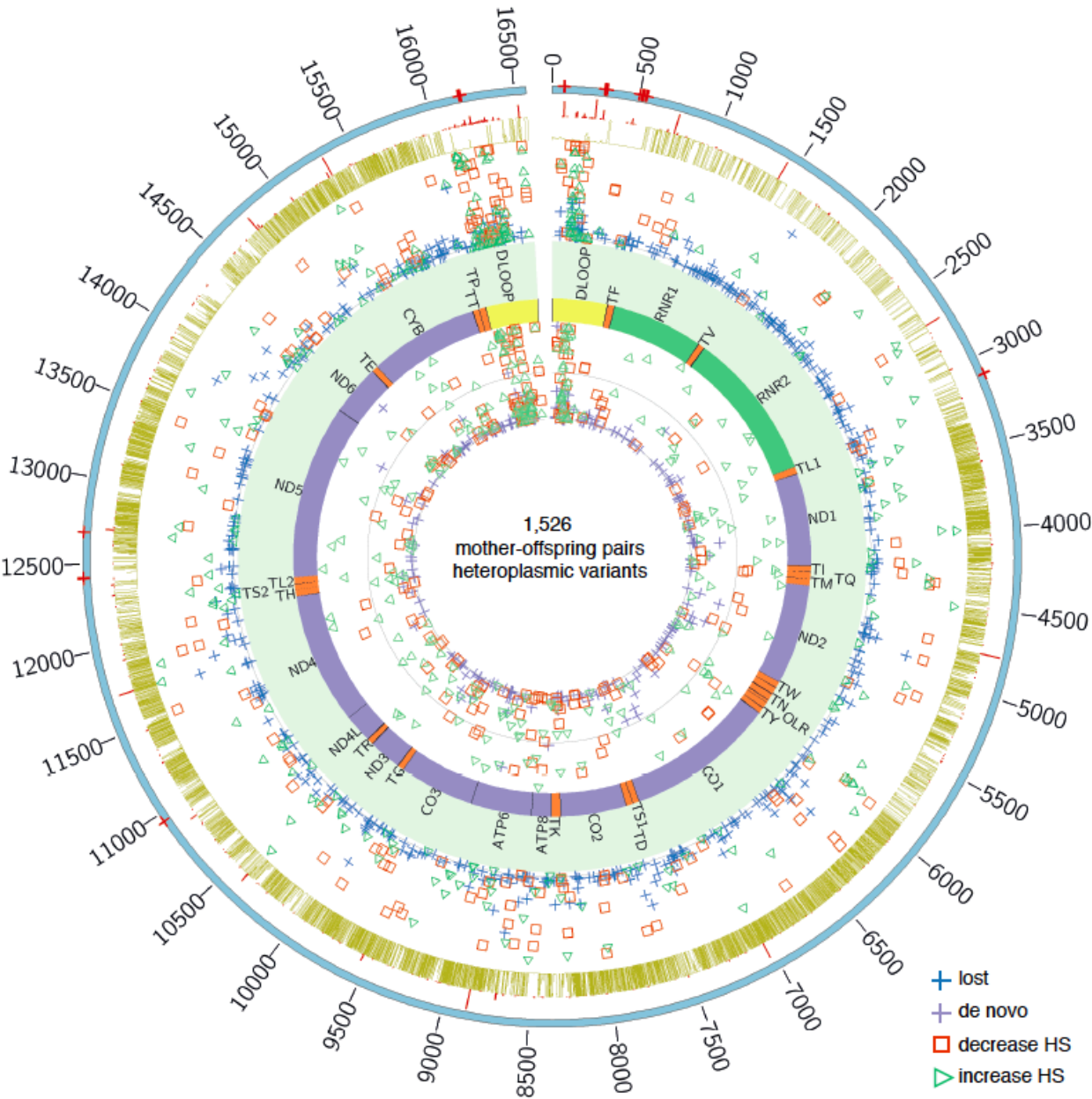


Fig 1. Circos plot of mitochondrial heteroplasmic variants identified in 1,526 mother-offspring pairs.

Circles from the outside to the inside indicate the following: (1) position of a variant on the mtDNA, the removed regions are shown in red crosses; (2) Minor allele frequency for common variants (MAF>1%) derived from 30,506 NCBI mtDNA sequences (5), where the radial axis corresponds to the MAF; (3) phastCons100way scores from UCSC (42) where the radial axis corresponds to the degree of conservation ; (4) heteroplasmic variants identified in the mothers where the radial axis corresponds to the heteroplasmy fraction, HF; (5) regions corresponding to the different mtDNA genes (yellow - D-loop, purple – coding region, green – rRNAs and orange - tRNAs); (6) heteroplasmic variants identified in the offspring where the radial axis corresponds to the heteroplasmy fraction, HF.

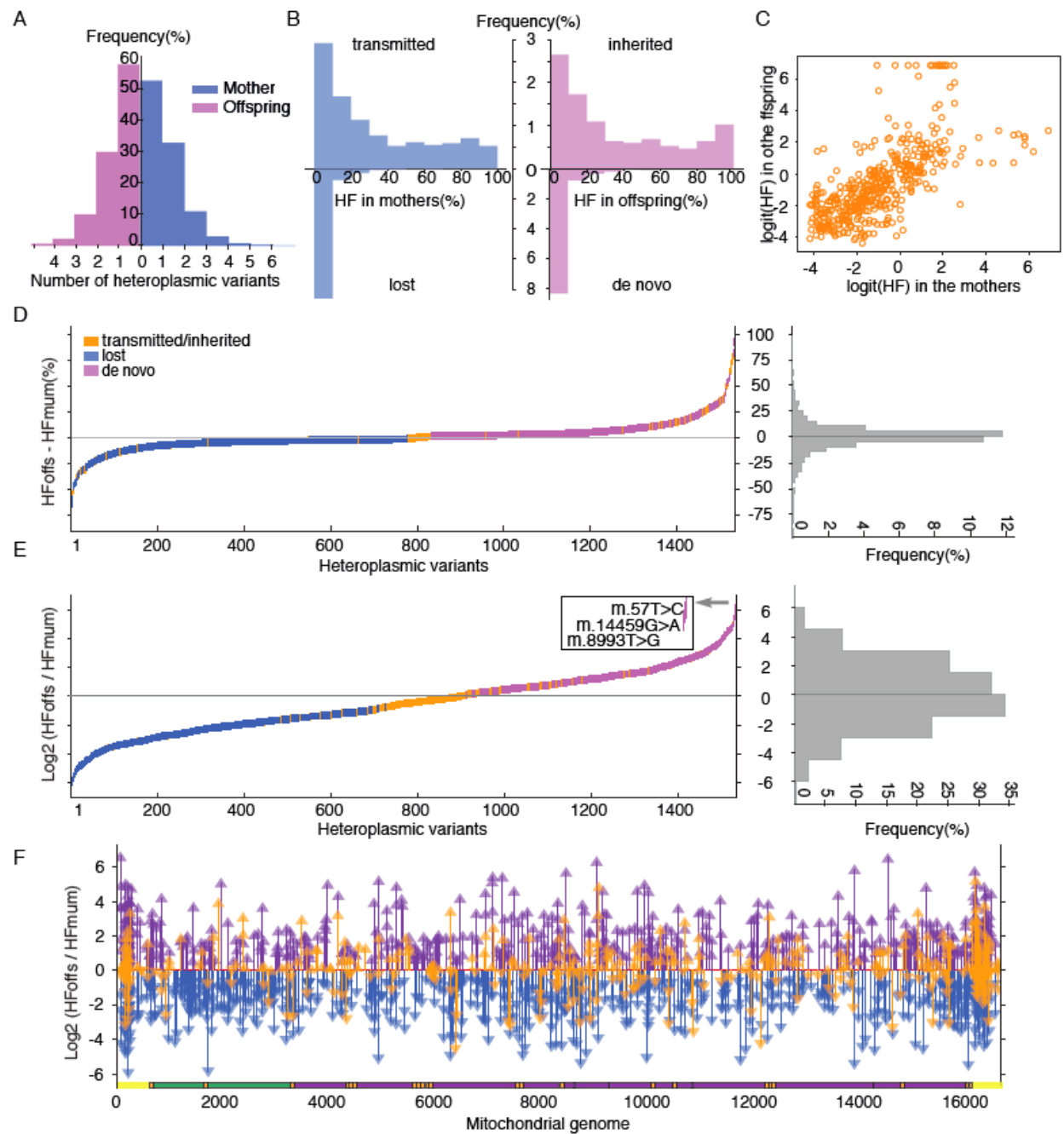


Fig 2. Transmission of heteroplasmic mtDNA variants in 1,562 mother-offspring pairs.

(A) Frequency distribution of number of heteroplasmic variants in the mothers and offspring. (B) Distribution of HF in the mothers and offspring, transmitted, inherited, lost and *de novo* are shown separately. (C) Scatter plot of logit(HF) in transmitted heteroplasmic variants between the mothers and offspring ($R^2=0.79$, $P=1.52 \times 10^{-93}$, Pearson's correlation). (D) Left, difference in the percentage shift of HF between the offspring and the corresponding mothers (HF_{offspring} - HF_{mother}) ordered by the degree of shift. Right, distribution of the difference of the percentage shift of HF between offspring and the corresponding mothers (HF_{offspring} - HF_{mother}). (E) Left, log₂ ratio of HF difference between offspring and the corresponding mothers ordered by the degree of log₂ ratio. Three increase HSs with values above 6 shown in the box. Right, distribution of log₂ ratio of HF difference between offspring and the corresponding mothers. (F) log₂ ratio of HF difference between offspring and the corresponding mothers aligned to the whole mitochondrial DNA sequence; the mtDNA regions are shown at the bottom bar in different colors (yellow - D-loop, purple – coding region, green – rRNAs and orange - tRNAs).

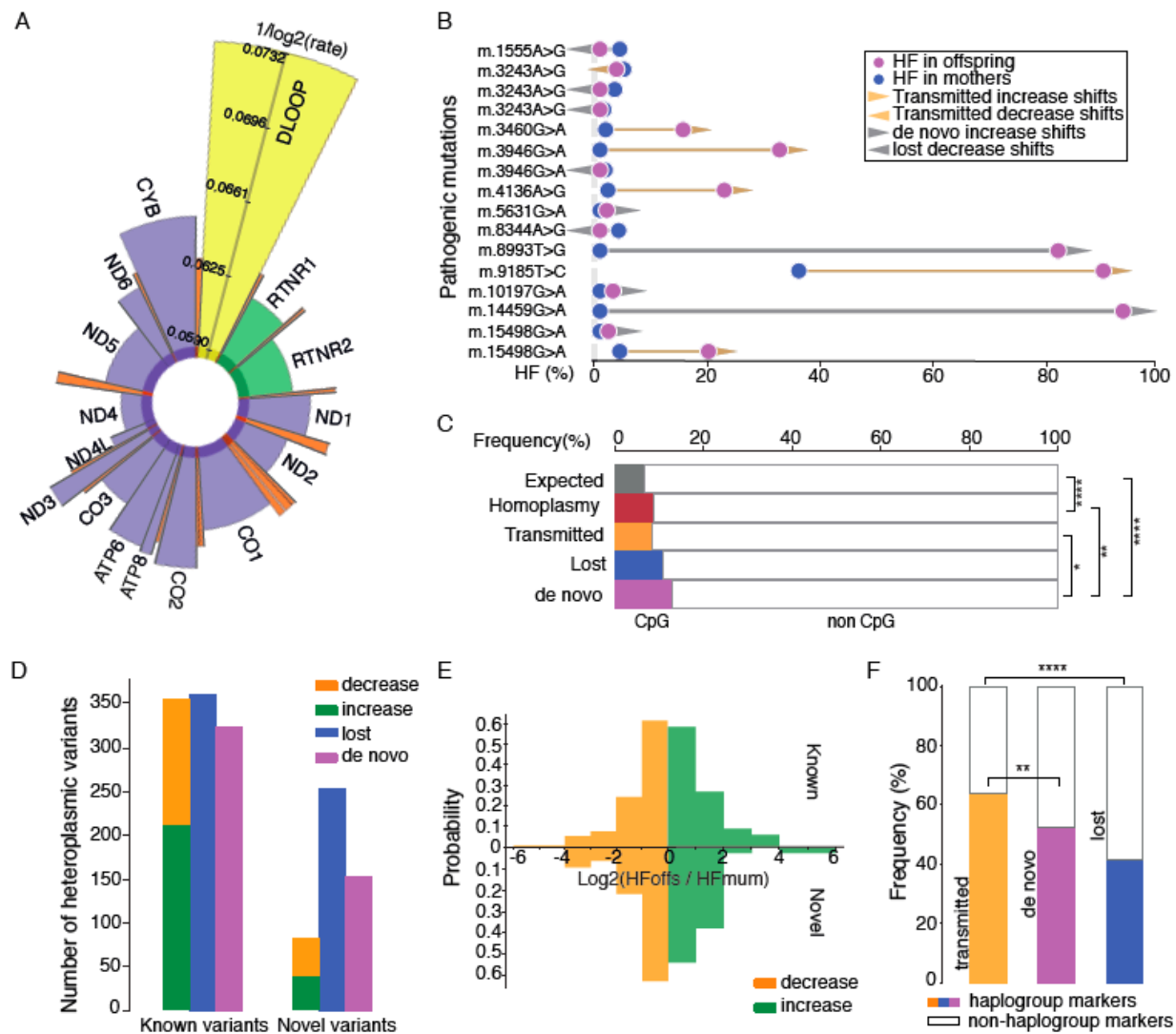


Fig 3. Characteristics of the heteroplasmic mtDNA variants in 1,562 mother-offspring pairs.

(A) Mutation rate of mtDNA genomic regions was estimated using 477 *de novo* heteroplasmic variants from 1,526 mother-offspring pairs detected at HF>1%. Vertical axes represent
5 $1/\log_2(\text{mutation rate})$ per base per mother-child transmission. mtDNA genomic regions are labeled and shown in different colors (yellow - D-loop, purple – coding region, green – rRNAs and orange - tRNAs). All tRNAs were combined to estimate the tRNA mutation rate. Note that this is the raw number of new mutations/bp/transmission detected at HF>1% in the offspring, and does not factor in the detection threshold nor segregation because current models assume
10 neutrality(13, 43), which we later show is not the case. (B) Pathogenic mutations were observed in 1,526 mother-offspring pairs. Each dot represents the HF in the mothers (blue) and the corresponding offspring (pink); the directions of arrow show increase (->) or decrease (<-) HS; the length of the arrow between each pair of points represents the change in HF (orange - transmitted heteroplasmic variants, grey - *de novo* / lost heteroplasmic variants). (C) Frequency
15 of heteroplasmic variants at CpG and non-CpG islands. Expected, homoplasmic variants, transmitted, lost and *de novo* heteroplasmic variants are shown separately. (D) Number of novel versus known heteroplasmic variants in transmitted, lost and *de novo* heteroplasmic variants, increase and decreasing HS in transmitted heteroplasmic variants are shown in different colors. (E) Distribution of HS between the offspring and the corresponding mothers in transmitted
20 known and novel heteroplasmic variants, increase and decreasing HS are shown in different colors. (F) Frequency of haplogroup defining variants in transmitted, lost and *de novo* heteroplasmic variants. The transmitted heteroplasmic variants were more likely to affect known haplogroup specific sites on the world mtDNA phylogeny than the lost and *de novo*

heteroplasmic variants ($P=7.86 \times 10^{-11}$ and $P=0.0016$ respectively, Fisher's exact test). P value < 0.05*, < 0.01**, < 0.001*** and <0.0001****.

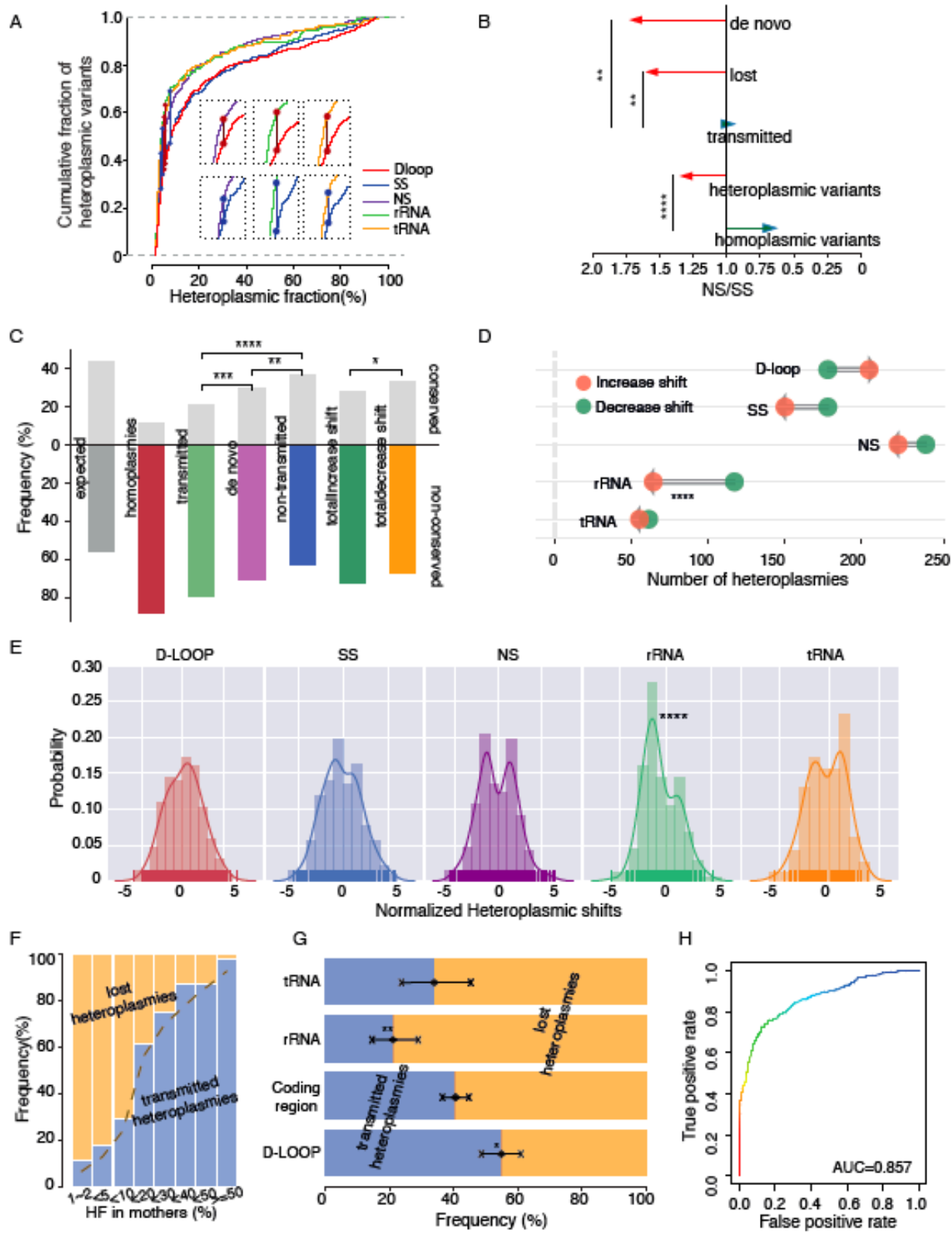


Fig 4. Evidence of selection during the transmission of mtDNA heteroplasmy in 1,526 mother-offspring pairs.

(A) Cumulative distributions of HF in the mothers and offspring within each mtDNA region.

Vertical lines between two curves show the greatest distance between Dloop / SS and NS / rRNA

5 / tRNA regions (*P*-values in **table S3**). (B) NS/SS ratio of NS and SS variants for observed

homoplasmic polymorphisms, total heteroplasmic variants, transmitted, lost and *de novo*

heteroplasmic variants. (C) Frequency of heteroplasmic variants affecting conserved and non-

conserved sites. Expected, homoplasmic variants, transmitted, lost, *de novo* heteroplasmic

variants, increase and decrease HSs are shown separately. (D) Number of heteroplasmy

10 showing an increase or decrease HF in each mtDNA region. Left-facing arrows indicate that the

number increasing was less than the number decreasing. Right-facing arrows indicate that the

number increasing was greater than the number decreasing. (E) Histograms of HS in each

mtDNA region with fitted kernel density curves. (F) Bar plot of the frequency of transmitted

heteroplasmic variants by bins of HF in the mothers. (G) Frequency of transmitted heteroplasmic

15 variants in each mtDNA region, along with 95% confidence intervals. (H) Receiver operating

characteristic (ROC) curve for the logistic regression model of transmission (Area under the

curve, AUC=0.857). *P* value < 0.05*, < 0.01**, < 0.001*** and <0.0001****.

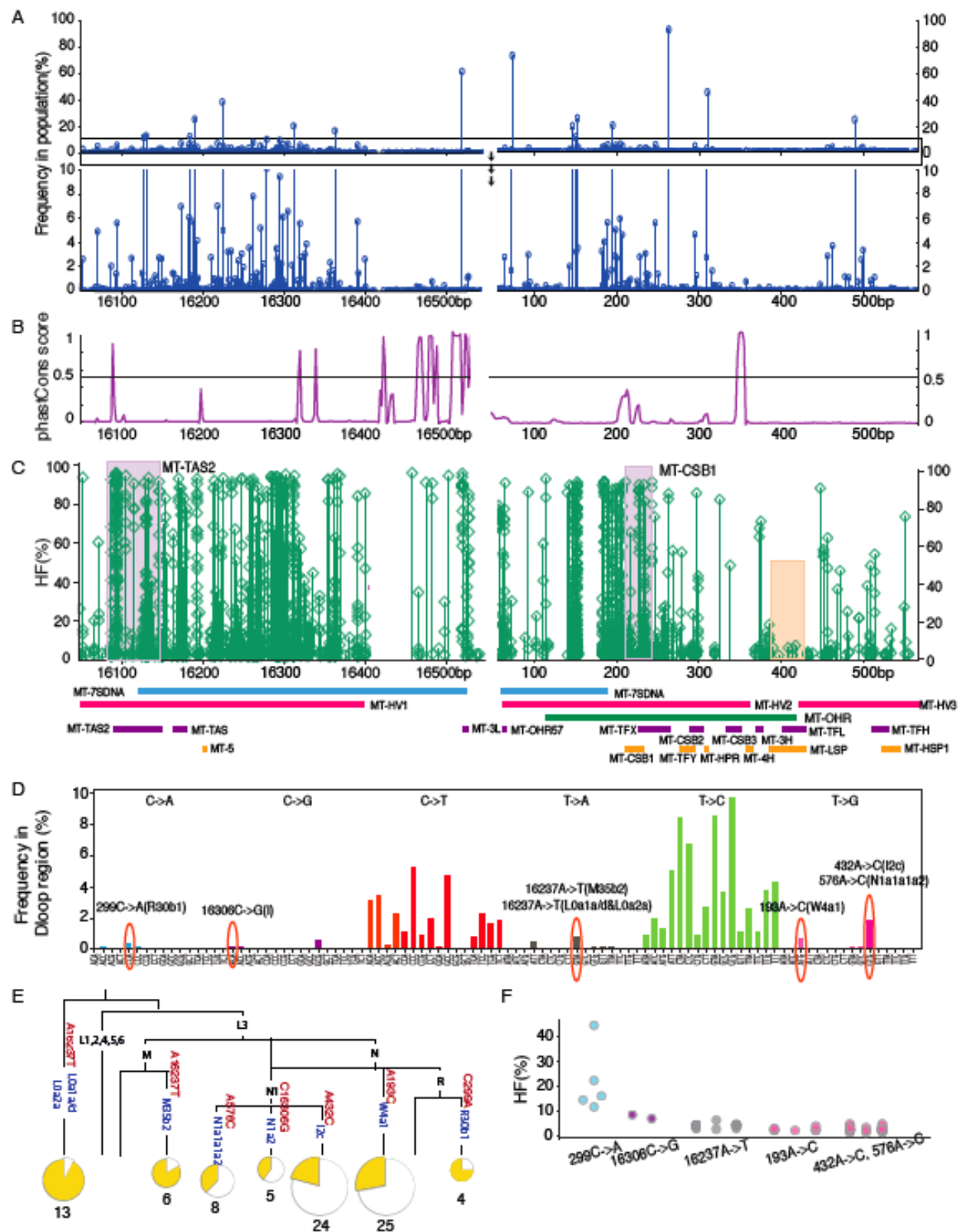


Fig 5. The distribution of heteroplasmic variants in mtDNA Dloop region.

(A) MAF of homoplasmic single nucleotide polymorphisms observed in 30,506 NCBI mtDNA sequences, with an expanded axis to show MAF<10% at the bottom. (B) Trend of PhastCons scores is shown across the mtDNA D-loop region. (C) HFs observed in 12,975 mtDNA sequences in the D-loop region. MT-TAS2 and MT-CSB1 are shadowed in light purple. MT-LSP is shadowed in light orange. Corresponding known sub-regions of the mtDNA D-loop are shown at the bottom. Key - MT-3H: mt3 H-strand control element, MT-3L: L-strand control element, MT-4H: mt4 H-strand control element, MT-7SDNA: 7S DNA, MT-CSB1: Conserved sequence block 1, MT-CSB2: Conserved sequence block 2, MT-CSB3: Conserved sequence block 3, MT-HPR: replication primer, MT-HSP1: Major H-strand promoter, MT-HV1: Hypervariable segment 1, MT-HV2: Hypervariable segment 2, MT-HV3: Hypervariable segment 3, MT-LSP: L-strand promoter, MT-OHR: H-strand origins, MT-OHR57: H-strand origin, MT-TAS: termination-associated sequence, MT-TAS2: extended termination-associated sequence, MT-TFH/MT-TFL/ MT-TFX/ MT-TFY/: mtTF1 binding site. (D) Trinucleotide mutational signature of heteroplasmic variants in the D-loop region in 12,975 mtDNA sequences. The bars representing the frequency for the six types of substitution are displayed in different colours. (E) Simplified mtDNA phylogeny tree showing 6 heteroplasmic variants (refer to main text). Variants are shown in red and haplogroups are shown in blue. The pie chart sizes are proportional to the number of samples (shown at the bottom) belonging to the corresponding haplogroup in 10,210 unrelated mtDNA sequences. The proportion of samples carrying each heteroplasmic variant within the same haplogroup is shown in yellow. (F) HF of six heteroplasmic variants shared by more than one individual who belong to the same haplogroup.

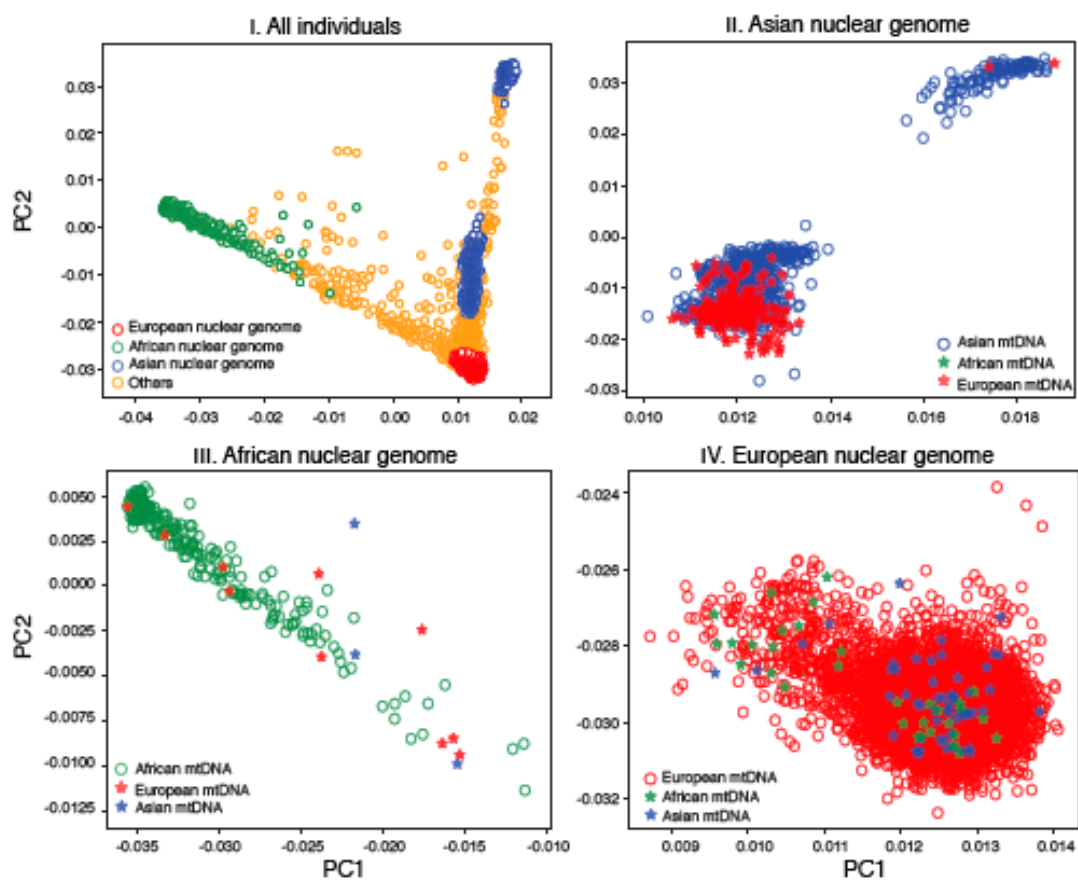
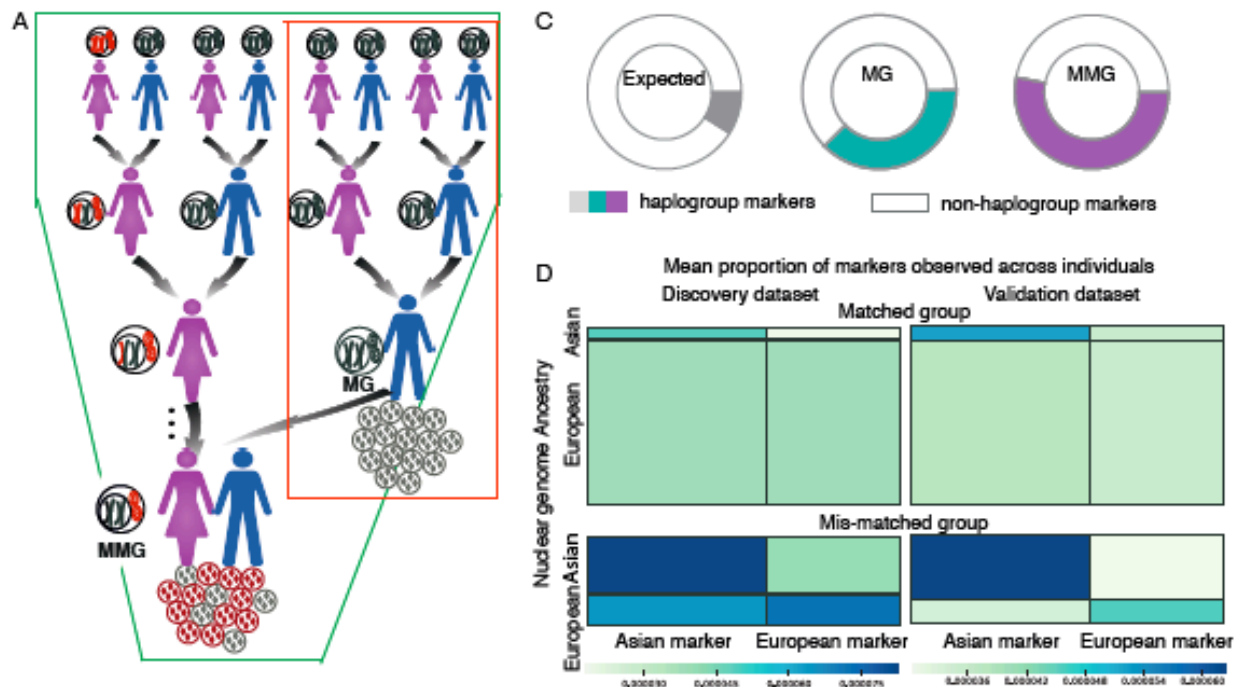


Fig 6. The characteristics of heteroplasmic variants in the nuclear ancestry and mtDNA ancestry matched and mismatched groups.

(A) Schematic showing how individuals with matched (MG, red border) and mismatched (MMG, green border) nuclear and mtDNA genomes arise over generations. Red and grey colored mtDNAs represent two different hypothetical populations. (B) I. Projection of the nuclear genotypes at common SNPs onto the two leading principal components computed with the 1000 Genomes dataset, with individuals colored by their assigned nuclear ancestry: Asian (blue), African (green), European (red) and Other (orange). The individuals colored in blue, green and red in the boxes labelled II, III and IV are shown in panels II, III and IV, respectively, where they are colored by their mitochondrial ancestries. Stars indicate that the mitochondrial ancestry does not match the nuclear ancestry. (C) Proportion of haplogroup defining variants in the matched and mismatched groups in 9,631 mtDNA sequences from unrelated individuals, along with the expected proportion shown at the left side. Distinct heteroplasmic sites were more likely to affect known haplogroup specific sites (26) than the rest of the mitochondrial genome compared to that expected by chance ($P < 2.2 \times 10^{-16}$, Fisher's exact test). This bias was stronger in the mismatched group than the matched group ($P = 0.001$, Fisher's exact test). (D) Heatmaps showing the density of observed heteroplasmic mtDNA haplogroup-defining variants in the observation dataset (left) and validation dataset (right). The matched (top) and mismatched (bottom) groups are shown separately, broken down by the nuclear ancestry of the carrier and the major haplogroup of the variants. The width of each column is proportional to the number of variants defining each of the two major haplogroups (Asian and European). Within each heatmap, the height of each row is proportional to the number of individuals having each nuclear ancestry. The density of heteroplasmic variants in each cell determines its color.