

## Mitochondrial genome variants associated with Amyotrophic Lateral Sclerosis and their haplogroup distribution

Marcelo R. S. Briones<sup>1</sup>, João H. Campos<sup>2</sup>, Renata C. Ferreira<sup>3,4</sup>, Lisa Schneper<sup>5</sup>, Ilda M. Santos<sup>2</sup>, Fernando M. Antoneli<sup>1</sup>, NYGC ALS Consortium<sup>6</sup>, James R. Broach<sup>5</sup>

<sup>1</sup>Center for Medical Bioinformatics, Escola Paulista de Medicina, Federal University of São Paulo, São Paulo, SP, Brazil, CEP 04039032.

<sup>2</sup>Graduate Program in Microbiology and Immunology, Federal University of São Paulo, São Paulo, SP, Brazil, CEP 04039032.

<sup>3</sup>Department of Neurology and Neurosurgery, Federal University of São Paulo, São Paulo, SP, Brazil, CEP 04039032.

<sup>4</sup>Bridges Genomics, M.E., São Paulo, SP, Brazil, CEP 043639040.

<sup>5</sup>Institute for Personalized Medicine, Department of Biochemistry, Pennsylvania State University College of Medicine, 500 University Drive, 17033, Hershey PA, USA.

<sup>6</sup>The New York Genome Center (NYGC) ALS Consortium members are listed in the Supplemental Material

### \*Corresponding authors:

Marcelo R. S. Briones  
Center for Medical Bioinformatics, Federal University of São Paulo  
Rua Pedro de Toledo 669 L4E, São Paulo SP CEP 04039032, Brazil.  
Tel 011 5576-4537

James R. Broach  
Institute for Personalized Medicine, Department of Biochemistry, Pennsylvania State University College of Medicine, 500 University Drive, 17033, Hershey PA, USA.

## Abstract

Amyotrophic Lateral Sclerosis (ALS) is a progressive, fatal, and non-reversible disease. It is the third most common neurodegenerative disease and the most frequent, adult-onset, motor neuron disease. Although familial and sporadic forms are observed, twin studies revealed that even sporadic forms have a significant genetic component. Variants in 55 nuclear genes have been associated with ALS and although mitochondrial dysfunction is associated with ALS, variants in mitochondrial genomes (mitogenomes) have not yet been associated with ALS. Here we conducted a genome wide association study (GWAS) in mitogenomes of 1,965 ALS patients and 2,547 controls. We identified 51 mitogenome variants with  $p$ -values  $<10^{-7}$  of which 13 variants have Odds Ratios (OR) $>1$ , in genes *RNR1*, *ND1*, *CO1*, *CO3*, *ND5*, *ND6* and *CYB*, while 38 variants have OR $<1$  in genes *RNR1*, *RNA2*, *ND1*, *ND2*, *CO2*, *ATP8*, *ATP6*, *CO3*, *ND3*, *ND4*, *ND5*, *ND6* and *CYB*. The frequency of haplogroups H, U and L, the most frequent in our ALS dataset, in the most severe, bulbar onset, versus limb, spinal and axial onsets shows that the proportion of Bulbar onset is the same in these haplogroups as compared to the other onset sites. Also, intra-haplogroup GWAS revealed unique ALS-associated variants in haplogroups L and U. Our study suggests that mitogenome variants (SNVs) could be included in routine genetic testing for ALS and that mitochondrial replacement therapy would have a potential basis for ALS treatment.

## KEYWORDS

Amyotrophic Lateral Sclerosis, mitochondrial genome, SNVs, GWAS, mitochondrial haplogroups.

## 1 | INTRODUCTION

Amyotrophic lateral sclerosis (ALS) is an incurable fatal disease with progressive degeneration of motor functions and causes enormous physical and psychological suffering in patients and their families. It is the third most common neurodegenerative disease and the most common, adult-onset, motor neuron disease<sup>1</sup>. The overall prevalence of ALS is 3.9 per 100,000 people in the United States, with the highest prevalence (17.0 per 100,000) among those aged 70 to 79 years. Men have a higher prevalence than women (4.8 vs 3.0 per 100,000), and Caucasians have a higher rate than African-Americans (4.2 vs 2.0 per 100,000)<sup>2</sup>. The worldwide prevalence ranges from 0.8 to 10.32 cases per 100,000 people<sup>3-6</sup>. Differences in regional prevalence may be explained by different methodologies, diagnostic criteria over time, and absence of suitable systems for disease notification<sup>7</sup>. Incidence of ALS is between 1.5 and 2.5 cases per 100,000 people per year<sup>8</sup>. The two most relevant forms, in epidemiological sense, are: (a) Sporadic ALS (SALS), corresponding to approximately 90% of all cases and is the most common form in the United States, and (b) Familial ALS (FALS), which corresponds to approximately 10% of all cases, occurs more than once in a family lineage (usually dominant inheritance) and accounts for a small number of cases in the United States. Different phenotypes can be observed in ALS, which depend on motor neuron pattern

involvement. Eight phenotypes are described with different participation of upper and lower motor neurons leading to bulbar onset or limb, spinal or axial onsets<sup>9</sup>.

ALS is associated with mitochondrial abnormal morphology and bioenergetics, which contribute to the denervation of motor neurons from muscles that occur in early stages of ALS<sup>10</sup>. Extensive fragmentation of mitochondria has been observed in cultured motor neuron-like NSC-34 cells overexpressing mutant *SOD1*<sup>11,12</sup>. Besides morphological alterations, mutant *SOD1* causes the loss of membrane potential and disruption of respiratory chain activity<sup>13,14</sup>. Dysregulation of electron transfer chain complexes has been observed in human ALS patients and G93A *SOD1* transgenic mice<sup>13,15</sup>. The mitochondrial respiratory chain contains components expressed by two genomes, nuclear and mitochondrial, and therefore mutations in both genomes affect cellular respiration<sup>16</sup>. The mitochondrial genome has a high mutation rate, which is 10 to 20-fold higher than the nuclear genome. It is organized in 466 to 806 nucleoids per cell, with each nucleoid containing 2 to 10 copies of mitochondrial genome<sup>17</sup>. The first pathogenic mtDNA mutations were identified in 1988. To date, over 250 pathogenic mtDNA mutations (point mutations and rearrangements) have been described. As many as 1 in 10,000 people have a clinically manifested mtDNA disease and 1 in 6,000 are at risk<sup>16</sup>. A full list of mtDNA diseases can be found at MITOMAP database<sup>18</sup>. Among most relevant diseases with reported mtDNA mutations are: dystonia, encephalomyopathy, Alzheimer's disease and Parkinson's disease, which are associated with mutations in mitochondrial encoded complex I NADH dehydrogenase subunits<sup>19</sup>.

At least nine genes are known to carry ALS-causing mutations and 28 others are implicated in ALS pathogenesis<sup>1</sup>. All genes associated with ALS so far identified are nuclear encoded and although mitochondrial dysfunction has been associated with ALS, no ALS-associated mutations in mitochondrial genomes have been described. A mutation in a nuclear gene encoding a mitochondrial protein (*CHCHD10*) is linked to a syndrome with ALS-FTD (frontotemporal dementia) features<sup>20</sup>. Twin studies have estimated ALS heritability to be 60-70%, suggesting that many patients who present with sporadic ALS (SALS) may also have an underlying genetic cause<sup>21</sup>. An initial limited study of mitochondrial DNA sequencing with thirty-eight ALS patients (23 males and 15 females) and 42 unaffected controls (18 males and 24 females) indicates that ALS patients have a higher mean number of variants in protein-coding genes, such as *ND4L*, *ND5*, *ND6*, and *ATP8* and that haplogroups Y and M7c may modulate the clinical expression of ALS<sup>22</sup>. This study, however, does not perform PCA analysis to verify population structure or perform standard GWAS statistics. Because of mitochondrial impairment in ALS and absence of mitochondrial genes among those associated with ALS, we performed a GWAS of 1,965 patients and 2,547 controls to determine whether variants in the mitochondrial genome are associated with ALS.

## 2 | METHODS

### 2.1 | Source of data

A total of 4,512 mitochondrial genomes were included in the present study main analysis. This dataset included 1,965 mitogenomes from ALS patients, being 1,930 from the NYGC ALS

Consortium and 35 sequenced in the Institute for Personalized Medicine at the Pennsylvania State University with REPLI-g Mitochondrial DNA Kit (Qiagen), library construction with Kapa Hyper Prep kit and adapters (BioO scientific) and sequencing with Illumina MiSeq 2x300. These sequences have been deposited in GenBank with accession numbers MZ458603-MZ460580 (Supplementary Table S8). Mitochondrial reads, from samples sequenced at the NYGC, were extracted from whole genome BAM files, using samtools view (v. 1.6.0). The control, non-ALS set, comprised 2,534 mitogenomes from the 1000 Genomes Project <sup>23</sup>, phase 3, plus 13 control mitogenomes sequenced in the Institute for Personalized Medicine at the Pennsylvania State University with the same method as for ALS samples (Illumina MiSeq 2x300), totaling 2,547 controls (Supplementary Table S9). The NYGC samples were prepared using either the Illumina truseq nano library prep or the Illumina PCR free library prep protocol and all sequence reads were obtained using the Illumina HiSeq X-10 platform (Supplementary Table S10).

## 2.2 | Genome assembly

Mitochondrial genomes were assembled from sequencing reads using the NCBI reference sequence *Homo sapiens* mitochondrion complete genome (NC\_012920.1). The reference NC\_012920.1 contig is 16,569 bp long and corresponds to the Revised Cambridge Reference Sequence (rCRS) belonging to the European haplogroup H2a2a1 <sup>24</sup>. Alignments of sequencing reads to the reference were done using the Burrows-Wheeler transformation method of the BWA-MEM algorithm as implemented in BWA version 0.7.17 <sup>25</sup>. The mean coverage of NYGC mitogenome assemblies is 8,000x while the mean coverage of the Penn State samples is 180x.

## 2.3 | GWAS statistical analysis

For GWAS statistical analysis, a VCF file was created with case and control samples. After initial mapping, the resulting BAM files were processed using the GATK best practices quality control <sup>26</sup> as follows: (a) all reads were sorted using SortSam algorithm, (b) duplicate reads were marked using MarkDuplicates algorithm (c) mitochondrial variants were called using the HaplotypeCaller, GenomicsDBImport and GenotypeGVCFs algorithms, as implemented in GATK version 4.0.1.0 <sup>27</sup>. The VCF file generated by this procedure comprised the 1,965 mitogenomes from ALS patients plus the 13 control mitogenomes. This VCF file was merged with the VCF file containing the 2,534 mitogenomes from the 1000 Genomes Project phase 3, using BCFTools package in SamTools <sup>28,29</sup>. The final merged VCF totaled 4,512 mitogenomes. After filtering by the standard quality control procedure as described by Li and colleagues <sup>28,29</sup>, a total of 1,946 variants remained.

The statistical analysis was performed using PLINK version 1.90b6.4 <sup>30</sup>. SNVs with a genotype failure rate >1% were removed. To keep only the common and the low frequency variants we also removed SNVs with Minimum Allele Frequency (MAF) <0.1%. The sex information was not used. A PCA was computed to correct for population stratification <sup>31</sup>. The association study was performed by logistic regression with the first 20 PCA eigenvectors as covariates. Bonferroni correction for multiple testing with the strict threshold of  $10^{-7}$  was

employed resulting in 51 variants. Among these, 43 were common variants (MAF  $\geq 1\%$ ) and 8 were low frequency variants (MAF  $>0.1\%$  and  $<1\%$ ). The  $p$ -values and Odds Ratios were calculated using PLINK. Estimation of Odds Ratios in 2x2 contingency tables with zeros was performed with Haldane-Anscombe correction<sup>32</sup>. Haplogroup assignment was done using Haplogrep2<sup>33</sup>. Variants were annotated using the NCBI Variation Viewer (<https://www.ncbi.nlm.nih.gov/variation/view/>). Genomic, clinical and functional annotations of all nucleotide changes that cause non-synonymous substitutions in human mitochondrial genes associated with ALS were annotated using MitImpact (<https://mitimpact.css-mendel.it/>). The Relative Risk was calculated from the Odds Ratios values according to the formula proposed by Zhang and collaborators<sup>34</sup>

The mitochondrial genomes of 1,978 samples sequenced at the New York Genome Center (1,965 cases and 13 controls) were scanned to detect homoplasmic or heteroplasmic sites using the Mutserve tool (version 2.0.0-rc12)<sup>35</sup>.

## 2.4 | Methods of alternative datasets analysis

In addition to the dataset described in the main text, two additional datasets were used in parallel analyses: 1. From the original dataset of 4,512 individuals, we randomly removed 860 affected individuals from haplogroups with overrepresented cases (H, J, K, T, U), to balance the haplogroup distribution, leaving a total of 3,652 individuals (1,106 cases and 2,546 controls). The same filters of the original analysis were used (removal of SNVs with genotype failure rate  $>1\%$  and MAF  $<0.1\%$ ) leaving a total of 1,934 variants; 2. We aligned 30,502 mitogenomes from GenBank as described in<sup>36</sup>. After removal of partial sequences, low quality samples (with Ns) and large gaps in 5' and 3' ends an alignment was obtained for a subset of 17,725 mitogenomes. A VCF file from this multiple alignment was generated and merged with the VCF file containing the 1,979 samples sequenced at the Pennsylvania State University Institute of Personalized Medicine, resulting in a total of 19,701 mitogenomes (1,965 cases and 17,736 controls). Using the same filters as described in the Methods section (removal of SNVs with genotype failure rate  $>1\%$  and MAF  $<0.1\%$ ), a total of 1,012 variants were obtained. The remaining statistical analysis of these two additional sets followed the same steps of the original dataset of 4,512 individuals.

## 3 | RESULTS

### 3.1 | Principal Component Analysis

The Principal Component Analysis (PCA) of the 1,965 ALS mitochondrial genomes reveal the expected populational structure with haplogroup L (present in higher frequency in populations with African background) showing significantly increased diversity as compared to haplogroups associated with European and Asian populations (Fig. 1). Figures 1B and 1C zoom in, with more detail, the "European/Asian" (e.g., H, K, J, U, M) haplogroups clustering far from L haplogroups in the same PCA. The PCA profile is maintained if the control samples (2,547) are added to the analysis (Supplementary Figure S1), and the distribution of cases and controls is shown to be quite similar (Supplementary Figure S2).

The ALS mitogenomes and control genomes were aligned, and variants with a genotyping rate  $>0.1$  or a Minimum Allele Frequency (MAF)  $<0.1\%$  were removed resulting in the retention of 1,946 variants (SNVs). The QQ-plot of all 1,946 variants show the expected slope increase as compared to the theoretical expected distribution (Fig. 2).

### 3.2 | Distribution of SNVs in the mitogenome

The Manhattan Plot of the 1,946 SNVs shows the distribution of the negative logarithm of the  $p$ -values of the case-control logistic regression for the corresponding SNVs along the mitogenome. There are 51 variants above the threshold value  $-\log_{10}(p) = 7$  (below the  $p$ -value  $10^{-7}$ ) (Fig. 3). These SNVs are significantly associated with ALS and their annotation reveals that they are scattered among almost all mitochondrial genes. For the  $p$ -value threshold of  $10^{-5}$  there are 73 SNVs above the threshold, for  $5 \times 10^{-8}$ , there are 48 SNVs above the threshold and for  $10^{-8}$ , there are 40 SNVs above the threshold. We used the threshold of  $10^{-7}$  because it is 1000-fold below the number of independent tests ( $10^4$ ) equivalent to the mitochondrial genome size. Nevertheless, with the threshold of  $10^{-7}$  we are still at a 99.9% significance level per SNV.

### 3.3 | Odds Ratios of mitogenome variants

Among these 51 ALS variants, the Odds Ratios (OR) of 13 SNVs are  $>1$  (Table 1) and 38 SNVs show OR  $<1$  (Table 2). The lower bound of confidence intervals of OR  $>1$  are always  $>1$  and in SNVs with OR  $<1$  the upper bound is always  $<1$  indicating that if the study was repeated and the range calculated each time, the true value would lie within these ranges on 95% of occasions. These higher confidence levels ascertain that the interval contains the true odds ratio and that in none of these ORs the confidence interval crosses the threshold value OR = 1. The odds ratios in Tables 1 and 2 are corrected for population structure by the logistic regression using 20 components from the PCA with cases and controls (Supplementary Figures S1 and S3). The change in magnitude and direction (from protective to risk increasing), is a direct function of the weights of the components used in the logistic regression. The Relative Risk values, calculated as described by Zhang and colleagues<sup>34</sup> to correct the adjusted odds ratios from logistic regression analyses, are included in Tables 1 and 2 for SNVs associated with ALS.

### 3.4 | Annotation of ALS SNVs

Annotation of ALS SNVs reveals that the control region (D-loop), protein coding genes and ribosomal RNA variants are associated with ALS (Table 1 and Table 2). The SNVs that increase the odds of ALS (Table 1) are in the hypervariable regions 1, 2 and 3, small subunit ribosomal ribonucleic acid (SSU rRNA) of the mitochondrial ribosome, NADs 5 and 6, Cytochrome c oxidases 1 and 3 and Cytochrome B. One nonsynonymous SNV with  $p$ -value  $<10^{-7}$  and OR  $>1$  is found in *ND6* at 14318 (Asp>Ser). Seven nonsynonymous SNVs with OR  $<1$  are observed in *ND2* (4824 Thr>Ala and 5178 Leu>Met), *ATP8* (8414 Leu>Phe), *ATP6* (8684 Thr>Ile and 8794



His>Tyr) and *ND5* (13105 Ile>Val and 13928 Ser>Thr) (Table 2). All other ALS SNVs are either synonymous or located in noncoding regions.

### 3.5 | Mitochondrial haplogroups and ALS SNVs

Distribution of haplogroups in cases and controls shows three main haplogroups, namely H, U and L (Fig. 4). Haplogroup H is associated with Western European lines of descent, haplogroup U is associated with European Mediterranean populations and haplogroup L is associated with African and Middle Eastern (L3) populations. Haplogroup H contains more cases than controls (1,016 vs 259), haplogroup U has approximately the same frequency of cases and controls (178 vs 187) and haplogroup L contains less cases than controls (95 vs 687) (Fig. 4).

Because these three haplogroups are the most frequent in the dataset we performed intra-haplogroup GWAS and the results are detailed in Table 3. In haplogroup H, all SNVs identified were also identified in the general GWAS (Tables 1 and 2). In haplogroup L no non-synonymous variants were observed, and only one variant is also present in the main GWAS (rs28358568) (Table 3). In haplogroup U three unique ALS-SNVs were identified by intra-haplogroup GWAS being 2 synonymous in genes *HV2* and *ND5* and one missense variant in *ND2* (5186 Trp>Cys) (Table 3). The small differences observed in intra-haplogroup versus inter-haplogroup analysis are due to fact that haplogroup H, that show the same SNPs in both analyses, have significantly more cases than controls whereas haplogroups L and U have more controls than cases (Fig. 4).

The frequency of haplogroups H, U and L in the most severe, Bulbar onset, *versus* limb, spinal and axial onsets was calculated. For this analysis we used 1,460 ALS cases for which data on the site of ALS onset was available. In ALS cases, 50.68% belonged to haplogroup H, 9.04% to haplogroup U and 4.73% to haplogroup L. In patients with Bulbar onset, 46.15% belonged to haplogroup H, 9.72% to haplogroup U and 4.86% to haplogroup L (Table 4). This suggests that the proportion of Bulbar onset versus other sites of onset is approximately the same in these haplogroups (Table 4).

To test if variations in the control dataset might affect the results, we produced two alternative datasets: (1) to verify haplogroup bias we generated one dataset with balanced haplogroup distribution by removal of 860 ALS samples to produce the dataset of 1,106 ALS samples, balanced for haplogroups, respective to the 2,546 control samples (Supplementary Figures S4, S5 and S6) and (2) another dataset with a larger control mitochondrial representation encompassing the 17,736 unaffected mitogenomes analyzed by Wei and collaborators<sup>36</sup> (Supplementary Figures S7, S8 and S9). Results with different control datasets reveal that the SNVs identified in association with ALS are constant especially when genes that contain SNPs that reduce the risk of ALS (Odds Ratios < 1) are compared (Supplementary Tables S2 to S5)<sup>36</sup>.

### 3.6. | Heteroplasmy

Analysis of heteroplasmy carried out on 1,978 mitogenomes (1,965 cases and 13 controls) shows that 26.7% of the sites found are heteroplasmic, and 73.3% are homoplasmic (Fig. 5)<sup>37</sup>. The distribution of variants found along the mtDNA reveals specific profiles of allelic exchange frequencies per position (Fig. 5), supported by high sequencing coverage (average of 8,000x) (Fig. 5). Homoplasmic and heteroplasmic sites are indicated as to their occurrence along the mtDNA (Fig. 5C), and the coverage, in Fig. 5D. A summary of heteroplasmy analysis is in Supplementary Table S6.

Analysis of non-synonymous SNVs by PolyPhen<sup>38</sup> indicates that 14318 (aAt/aGt) is deleterious (FatHm), 4824 (Acc/Gcc) 5178 (Cta/Ata) and 8414 (Ctc/Ttc) are possibly damaging (PolyPhen2) and deleterious (CADD and MtoolBox). SNVs 8684 (aCc/aTc), 8794 (Cac/Tac), 13105 (Atc/Gtc) and 13928 (aGc/aCc) are predicted to have either benign or neutral effects (Supplementary Table S7). However, all are listed as benign by ClinVar.

Among variants that are associated with increased risk for ALS (Table 1), only one (rs28358587) has an association described for another phenotype (Resistance to high altitude pulmonary edema, HAPE), while 13 variants that are associated with reduced risk (Table 2) have other disease associations described (Supplementary Table S8).

#### 4| DISCUSSION

Mitochondrial dysfunction is a key factor in Amyotrophic Lateral Sclerosis (ALS). Although ALS has 40-60% heritability, with 5-10% family history, the typical Mendelian pattern is not always observed. Mitochondrial DNA is a major cause of non-Mendelian, extrachromosomal inheritance, especially in muscular and neuronal diseases. Our study analyzed 1,965 mitochondrial genomes of American patients and identified 51 mutations statistically associated with the disease. These mitochondrial mutations might contribute to refinement in diagnostic testing. Furthermore, mitochondrial dysfunction is well associated with oxidative stress in ALS. This increased mutagenic mitochondrial milieu might enhance the frequency of somatic, heteroplasmic mutations that, in turn, might contribute to deviations from a Mendelian inheritance pattern. Although mitochondrial dysfunction is a central feature in ALS motor neuron degeneration, it is not clear whether increased mitogenome mutational load leads to phenomena that trigger ALS degenerative processes or if increased superoxide accumulation due to mitochondrial malfunction leads to increased mitogenome mutational loads in ALS. For example, in ALS mice expressing *SOD1* G73A mutations, Sod1 accumulates in an inactivated, misfolded form in mitochondria<sup>39</sup>. Sod1 catalyzes the reduction of superoxide into molecular oxygen. Sod1 might affect mitochondrial respiration by binding to Bcl2, promoting aggregation and interfere with superoxide removal<sup>40</sup>. Reactive oxygen species can cause DNA damage thus increasing the mutational load<sup>41</sup>. Because of specific accumulation of Sod1 and the consequent accumulation of superoxide in motor neuron mitochondria, it is plausible to infer that in these cells the probability of somatic mutations is much higher than in other tissues where inactive *SOD1* does not accumulate. Furthermore, these mitochondrial mutations could be responsible for progressive degeneration of motor neurons because they might begin to affect the respiratory chain genes and proteins. Neurons



are completely dependent on oxidative respiration for its function and a decrease in respiratory rates leads to neuron death <sup>42</sup>.

Other degenerative diseases such as Parkinson's disease, also exhibit progressive damage related to respiratory-chain-deficient neurons <sup>43</sup>. In addition, reduced respiration caused by damage to mtDNA has been shown to be related to aging and muscular denervation <sup>44</sup>. It has been shown that protein components of complex I of the mitochondrial respiratory chain were reduced or absent in a proportion of aged motor neurons, a phenomenon not observed in fetal tissue. In this same study, additional data shows that complex I-deficient neurons had reduced mitochondrial DNA content and smaller body size. The authors propose that mitochondrial dysfunction in these motor neurons could lead to the cell loss and ultimately denervation of muscle fibers <sup>44</sup>. In accordance with these data, we observed in this present study, via GWAS, that at least eight nonsynonymous mutations in mitochondrion encoded genes could affect mitochondrial function and are associated with ALS with  $p$ -values  $<10^{-7}$ . These nonsynonymous mutations led to two amino acid changes (m.4824A>G Thr>Ala, m.5178C>A Leu>Met) in *ND2*, two (m.1305A>G Ile>Val, m.13928G>C Ser>Thr) in *ND5*, one amino acid change (m.14318T>C Asn>Ser) in *ND6* (*ND2*, *ND5* and *ND6* are members of respiratory complex I), two changes (m.8684C>T Thr>Ile, m.8794C>T His>Tyr) in *ATP6* (complex V), and one change (m.8414C>T Leu>Phe) in *ATP8* (also complex V) (Table 1 and Table 2).

The total number of SNVs associated with ALS in this case-control sample is 51 for  $p$ -values  $<10^{-7}$ . This threshold is extremely low considering the size of mitogenome (16,569bp) of approximately  $1.6 \times 10^4$  and the order of SNVs tested is  $10^3$  (1,946 SNVs) which would implicate that a strict Bonferroni correction would draw the threshold around  $p$ -values below  $10^{-4}$ . The effective number of mitochondrial SNVs associated with ALS would therefore be higher. However, the threshold used here, with logistic regression, is extremely conservative to mitigate false positives. The SNVs selected are well above the expected positive association by sheer chance due to multiple tests (Fig. 2). Also, the threshold of  $10^{-7}$  would be used in case mtDNA was analyzed together with the nuclear genome in the same GWAS. Therefore the  $10^{-7}$  is used so that it is possible to infer that the 51 SNVs here selected would still be statistically associated with ALS in a whole genome analysis.

Sequences in the control set of the main association analysis (Fig. 3) are the low-coverage whole genome sequencing (lcWGS) from phase 3 of the 1000 Genomes Project <sup>45</sup>. While lcWGS data returns a low nuclear DNA coverage (usually around 1x to 2x), it still includes a high mtDNA coverage due to the much higher copy number of mtDNA, which is present at hundreds or thousands of copies per cell <sup>46</sup>. The mean mtDNA coverage for each individual ranged from 37.7x to 3,535x <sup>47</sup>. Therefore, while standard variant calling on the low-coverage nuclear DNA is challenging and prone to errors, reliable variant calling in the mitochondrial genome is feasible because the mtDNA read depth in the same sample is hundreds to thousands of times higher. In fact, lcWGS has been used to analyze mtDNA as part of the 1000 Genomes Project <sup>47</sup>, as well as phylogenetic analysis of mitogenomes <sup>48</sup>. Our present work exploits the high mtDNA read depth in lcWGS, which meets standard variant calling thresholds,

to detect SNPs associated with ALS. Also, the usual concern in using the 1000 Genomes Project sequences as controls in association studies focuses on the problem of imputation<sup>49</sup>, which was not an issue in the analysis presented here.

Different control groups were used for testing haplogroup bias and wider representation. Analyses with both supported conclusions of the main analysis. The differences between haplogroups in cases and controls indicate that mutations that simply differentiate the haplogroups are not artificially inflating the ALS associated SNVs. In intra-haplogroup analysis the SNVs associated with ALS are the same as compared to general GWAS. Because H is the most abundant haplogroup in ALS cases, at least in the population here tested, the SNVs below the  $p$ -value threshold of  $10^{-7}$  still hold. As the sample number decreases, the  $p$ -values increase and that is why in U and L fewer SNVs are positive for ALS association, most likely because the sample is smaller (Tables 1, 2 and 3). To evidence the impact of population structure in our analysis, the odds ratios obtained if the population was unstructured are depicted in Supplementary Table S1. For example, in the case of rs28358875 (mtDNA position 7196) the odds ratio without correction is OR = 0.2725 by either Chi-Square or Fisher's exact test (Supplementary Table S1), and with correction by logistic regression, OR = 136.5 (Table 1). The differences in results obtained with different control datasets are a general characteristic of GWAS, that might lead to false positives depending on the cutoff values considered. Corrections for population heterogeneity and structure mitigate these false positives<sup>50</sup>.

It is noteworthy that 30% of SNVs that increase the risk of ALS are in the control region (D-loop variable regions 1 and 2) while 16% of protective SNVs are in this region. It is therefore possible that besides direct effects in respiratory complex subunits, ALS mitochondria might be affected by variants associated with mtDNA replication.

A previous genomic study of ALS patients from mainland China reported that ALS patients have a higher mean number of variants in protein-coding genes, such as *ND4L*, *ND5*, *ND6*, and *ATP8* and that haplogroups Y and M7c may modulate the clinical expression of ALS (22). However, this conclusion was based on the analysis of whole mitochondrial genomes of only thirty-eight ALS patients and 42 unaffected controls. Moreover, the haplogroup analysis, while conducted on a larger number of individuals (585 ALS patients and 371 healthy controls), interrogated only the mitochondrial D-loop sequences for purposes of haplogroup classification. However, important substitutions that characterize several haplogroups occur outside the D-loop. Finally, this group did not report population structure by PCA in either case or controls to correct for underlying populational structure in their data. In our present study we found SNPs associated with ALS in genes *ND4*, *ND5*, *ND6*, and *ATP8* (Tables 1 and 2), however, due to the limitations of the Ni et al. study<sup>22</sup>, we cannot verify if the SNPs described in our work coincide with the variants they observed. Regarding haplogroups, both Y and M are significantly underrepresented in our set of 1,965 ALS patients (Figure 5), which precludes us from verifying that haplogroups Y and M7c modulate the clinical expression of ALS.

Our tests cannot establish, without experimental confirmation, whether the SNVs associated with ALS here described have a causal impact on ALS. For this, further experiments with cybrids should be performed and additional analysis with larger datasets with almost exact haplogroup frequencies in cases and controls should be done. This study is however the first to explicitly test the association of mitochondrial genome variants with ALS using global analysis techniques and to consider different mitochondrial haplogroups. Odds ratios tend to overestimate the risk ratios if the outcome is common in a group. Also, as our data show, the variance in African haplogroups is substantially larger than found in Asian and European haplogroups. For this reason, in clinical practice, ethnic background should be taken into consideration upon reporting the risk for ALS. Nevertheless, sufficient statistical data was presented to suggest that these variants could be included in panels for ALS testing in conjunction with nuclear genome variants. The combined analysis of nuclear and mitochondrial mutations associated with ALS could provide important information on covariance between these two genomes and mitonuclear compatibility. It is important to recall that strong conformational physical association between gene products of the mitogenome and the nuclear genome are essential for the formation of respiratory complexes I, III, IV and V. Respiratory complex II just contains subunits encoded by the nuclear genome, that are imported into the mitochondria. Although nonsynonymous mutations associated with ALS were present only in subunits of complex I, III and V, synonymous mutations associated with ALS were identified and might interact and exert effects at the level of mRNA and not simply at the level of translated products. These synonymous SNVs might affect the mRNA stability, the association with mitochondrial rRNAs and tRNAs and impose selective constraints in mitochondrial codon usage. The mitonuclear ALS associations should be studied in more detail because they could provide novel strategies for diagnosis and treatment of ALS.

#### AUTHOR CONTRIBUTIONS

**Marcelo R. S. Briones:** Conceptualization; investigation; methodology; formal analysis; data curation; writing – original draft. **João H. Campos:** Investigation; methodology; formal analysis; data curation; visualization; editing – original draft. **Renata C. Ferreira:** Conceptualization; supervision; writing – review and editing. **Lisa Schneper:** Investigation; methodology; validation; data curation; supervision; writing – original draft. **Ilda M. Santos:** Investigation; methodology; validation; data curation. **Fernando M. Antoneli:** Investigation; methodology; formal analysis; statistical analysis; writing – original draft. **James R. Broach:** Conceptualization; investigation; methodology; formal analysis; data curation; writing – original draft.

#### ACKNOWLEDGEMENTS

The authors thank Fabricio Landi for assistance with computational facilities. J.H.C. was supported by a PhD fellowship from CAPES, I.M.S. was supported by an MSc fellowship from CAPES. This work was supported by grants to M. R. S. B. (FAPESP 2014/25602-6, FAPESP 2013/07838-0 and CNPq 303912/2017-0) and NIH grants to J. R. B. (#). All NYGC ALS Consortium activities are supported by the ALS Association (ALSA, 19-SI-459) and the Tow Foundation.

#### CONFLICT OF INTEREST STATEMENT

None of the authors has any conflict of interest to disclose.

#### DATA AVAILABILITY STATEMENT

All data produced in the present study are available upon reasonable request to the corresponding authors.

#### ETHICS STATEMENT

We confirm that we have read the Journal's position on issues involved in ethical publication and affirm that this report is consistent with those guidelines.

#### REFERENCES

1. Renton AE, Chiò A, Traynor BJ. State of play in amyotrophic lateral sclerosis genetics. *Nat Neurosci*: 2014;17:17–23. <https://doi.org/10.1038/nn.3584>.
2. Mehta P. Prevalence of Amyotrophic Lateral Sclerosis — United States, 2012–2013. *MMWR Surveill Summ*: 2016;65. <https://doi.org/10.15585/mmwr.ss6508a1>.
3. Chancellor AM, Warlow CP. Adult onset motor neuron disease: worldwide mortality, incidence and distribution since 1950. *J Neurol Neurosurg Psychiatry*: 1992;55:1106–1115.
4. Román GC. Neuroepidemiology of amyotrophic lateral sclerosis: clues to aetiology and pathogenesis. *J Neurol Neurosurg Psychiatry*: 1996;61:131–137.
5. Mandrioli J, Faglioni P, Merelli E, Sola P. The epidemiology of ALS in Modena, Italy. *Neurology*: 2003;60:683–689.
6. Huisman MHB, de Jong SW, van Doormaal PTC, Weinreich SS, Schelhaas HJ, van der Kooij AJ, et al. Population based epidemiology of amyotrophic lateral sclerosis using capture-recapture methodology. *J Neurol Neurosurg Psychiatry*: 2011;82:1165–1170. <https://doi.org/10.1136/jnnp.2011.244939>.
7. Linden Junior E, Becker J, Schestatsky P, Rotta FT, Marrone CD, Gomes I. Prevalence of amyotrophic lateral sclerosis in the city of Porto Alegre, in Southern Brazil. *Arquivos de Neuro-Psiquiatria*: 2013;71:959–962. <https://doi.org/10.1590/0004-282X20130177>.
8. Logroscino G, Traynor BJ, Hardiman O, Chio' A, Couratier P, Mitchell JD, et al. Descriptive epidemiology of amyotrophic lateral sclerosis: new evidence and unsolved issues. *J Neurol Neurosurg Psychiatry*: 2008;79:6–11. <https://doi.org/10.1136/jnnp.2006.104828>.

9. Swinnen B, Robberecht W. The phenotypic variability of amyotrophic lateral sclerosis. *Nat Rev Neurol*: 2014;10:661–670. <https://doi.org/10.1038/nrneurol.2014.184>.
10. Shi P, Gal J, Kwinter DM, Liu X, Zhu H. Mitochondrial dysfunction in amyotrophic lateral sclerosis. *Biochim Biophys Acta*: 2010;1802:45–51. <https://doi.org/10.1016/j.bbadis.2009.08.012>.
11. Menzies FM, Cookson MR, Taylor RW, Turnbull DM, Chrzanowska-Lightowlers ZMA, Dong L, et al. Mitochondrial dysfunction in a cell culture model of familial amyotrophic lateral sclerosis. *Brain*: 2002;125:1522–1533.
12. Raimondi A, Mangolini A, Rizzardini M, Tartari S, Massari S, Bendotti C, et al. Cell culture models to investigate the selective vulnerability of motoneuronal mitochondria to familial ALS-linked G93ASOD1. *Eur J Neurosci*: 2006;24:387–399. <https://doi.org/10.1111/j.1460-9568.2006.04922.x>.
13. Borthwick GM, Johnson MA, Ince PG, Shaw PJ, Turnbull DM. Mitochondrial enzyme activity in amyotrophic lateral sclerosis: implications for the role of mitochondria in neuronal cell death. *Ann Neurol*: 1999;46:787–790.
14. Jung C, Higgins CMJ, Xu Z. A quantitative histochemical assay for activities of mitochondrial electron transport chain complexes in mouse spinal cord sections. *J Neurosci Methods*: 2002;114:165–172.
15. Bowling AC, Schulz JB, Brown RH, Beal MF. Superoxide dismutase activity, oxidative damage, and mitochondrial energy metabolism in familial and sporadic amyotrophic lateral sclerosis. *J Neurochem*: 1993;61:2322–2325.
16. Tuppen HAL, Blakely EL, Turnbull DM, Taylor RW. Mitochondrial DNA mutations and human disease. *Biochim Biophys Acta*: 2010;1797:113–128. <https://doi.org/10.1016/j.bbabi.2009.09.005>.
17. Chen XJ, Butow RA. The organization and inheritance of the mitochondrial genome. *Nat Rev Genet*: 2005;6:815–825. <https://doi.org/10.1038/nrg1708>.
18. Ruiz-Pesini E, Lott MT, Procaccio V, Poole JC, Brandon MC, Mishmar D, et al. An enhanced MITOMAP with a global mtDNA mutational phylogeny. *Nucleic Acids Res*: 2007;35:D823–D828. <https://doi.org/10.1093/nar/gkl927>.
19. Lott MT, Leipzig JN, Derbeneva O, Xie HM, Chalkia D, Sarmady M, et al. mtDNA Variation and Analysis Using Mitomap and Mitomaster. *Curr Protoc Bioinformatics*: 2013;44:1.23.1-26. <https://doi.org/10.1002/0471250953.bi0123s44>.
20. Bannwarth S, Ait-El-Mkadem S, Chaussonnet A, Genin EC, Lacas-Gervais S, Fragaki K, et al. A mitochondrial origin for frontotemporal dementia and amyotrophic lateral sclerosis through CHCHD10 involvement. *Brain*: 2014;137:2329–2345. <https://doi.org/10.1093/brain/awu138>.
21. Al-Chalabi A, Fang F, Hanby MF, Leigh PN, Shaw CE, Ye W, et al. An estimate of amyotrophic lateral sclerosis heritability using twin data. *J Neurol Neurosurg Psychiatry*: 2010;81:1324–1326. <https://doi.org/10.1136/jnnp.2010.207464>.

22. Ni J, Liu Z, Yuan Y, Li W, Hu Y, Liu P, *et al.* Mitochondrial genome variations are associated with amyotrophic lateral sclerosis in patients from mainland China. *J Neurol*: 2022;269:805–814. <https://doi.org/10.1007/s00415-021-10659-7>.
23. Clarke L, Fairley S, Zheng-Bradley X, Streeter I, Perry E, Lowy E, *et al.* The international Genome sample resource (IGSR): A worldwide collection of genome variation incorporating the 1000 Genomes Project data. *Nucleic Acids Research*: 2017;45:D854–D859. <https://doi.org/10.1093/nar/gkw829>.
24. Andrews RM, Kubacka I, Chinnery PF, Lightowlers RN, Turnbull DM, Howell N. Reanalysis and revision of the Cambridge reference sequence for human mitochondrial DNA. *Nat Genet*: 1999;23:147. <https://doi.org/10.1038/13779>.
25. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv:13033997 [q-bio]*: March 2013.
26. Anderson CA, Pettersson FH, Clarke GM, Cardon LR, Morris AP, Zondervan KT. Data quality control in genetic case-control association studies. *Nat Protoc*: 2010;5:1564–1573. <https://doi.org/10.1038/nprot.2010.116>.
27. Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, Del Angel G, Levy-Moonshine A, *et al.* From FastQ data to high confidence variant calls: the Genome Analysis Toolkit best practices pipeline. *Curr Protoc Bioinformatics*: 2013;43:11.10.1-33. <https://doi.org/10.1002/0471250953.bi1110s43>.
28. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics*: 2009;25:2078–2079. <https://doi.org/10.1093/bioinformatics/btp352>.
29. Li H. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics*: 2011;27:2987–2993. <https://doi.org/10.1093/bioinformatics/btr509>.
30. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, *et al.* PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet*: 2007;81:559–575. <https://doi.org/10.1086/519795>.
31. Hagen CM, Gonçalves VF, Hedley PL, Bybjerg-Grauholm J, Bækvad-Hansen M, Hansen CS, *et al.* Schizophrenia-associated mt-DNA SNPs exhibit highly variable haplogroup affiliation and nuclear ancestry: Bi-genomic dependence raises major concerns for link to disease. *PLOS ONE*: 2018;13:e0208828. <https://doi.org/10.1371/journal.pone.0208828>.
32. Ruxton GD, Neuhäuser M. Review of alternative approaches to calculation of a confidence interval for the odds ratio of a 2 × 2 contingency table. *Methods in Ecology and Evolution*: 2013;4:9–13. <https://doi.org/10.1111/j.2041-210x.2012.00250.x>.
33. Weissensteiner H, Pacher D, Kloss-Brandstätter A, Forer L, Specht G, Bandelt H-J, *et al.* HaploGrep 2: mitochondrial haplogroup classification in the era of high-throughput sequencing. *Nucleic Acids Res*: 2016;44:W58–W63. <https://doi.org/10.1093/nar/gkw233>.
34. Zhang J, Yu KF. What's the Relative Risk? A Method of Correcting the Odds Ratio in Cohort Studies of Common Outcomes. *JAMA*: 1998;280:1690–1691. <https://doi.org/10.1001/jama.280.19.1690>.



35. Weissensteiner H, Forer L, Fuchsberger C, Schöpf B, Kloss-Brandstätter A, Specht G, *et al.* mtDNA-Server: next-generation sequencing data analysis of human mitochondrial DNA in the cloud. *Nucleic Acids Research*: 2016;44:W64–W69. <https://doi.org/10.1093/nar/gkw247>.
36. Wei W, Gomez-Duran A, Hudson G, Chinnery PF. Background sequence characteristics influence the occurrence and severity of disease-causing mtDNA mutations. *PLOS Genetics*: 2017;13:e1007126. <https://doi.org/10.1371/journal.pgen.1007126>.
37. Wei W, Keogh MJ, Wilson I, Coxhead J, Ryan S, Rollinson S, *et al.* Mitochondrial DNA point mutations and relative copy number in 1363 disease and control human brains. *Acta Neuropathol Commun*: 2017;5:13. <https://doi.org/10.1186/s40478-016-0404-6>.
38. Adzhubei I, Jordan DM, Sunyaev SR. Predicting functional effect of human missense mutations using PolyPhen-2. *Curr Protoc Hum Genet*: 2013;Chapter 7:Unit7.20. <https://doi.org/10.1002/0471142905.hg0720s76>.
39. Jaarsma D, Rognoni F, van Duijn W, Verspaget HW, Haasdijk ED, Holstege JC. CuZn superoxide dismutase (SOD1) accumulates in vacuolated mitochondria in transgenic mice expressing amyotrophic lateral sclerosis-linked SOD1 mutations. *Acta Neuropathol*: 2001;102:293–305.
40. Tan W, Nanche N, Bogush A, Pedrini S, Trotti D, Pasinelli P. Small Peptides against the Mutant SOD1/Bcl-2 Toxic Mitochondrial Complex Restore Mitochondrial Function and Cell Viability in Mutant SOD1-Mediated ALS. *J Neurosci*: 2013;33:11588–11598. <https://doi.org/10.1523/JNEUROSCI.5385-12.2013>.
41. Keyer K, Gort AS, Imlay JA. Superoxide and the production of oxidative DNA damage. *J Bacteriol*: 1995;177:6782–6790.
42. Bal-Price A, Brown GC. Inflammatory neurodegeneration mediated by nitric oxide from activated glia-inhibiting neuronal respiration, causing glutamate release and excitotoxicity. *J Neurosci*: 2001;21:6480–6491.
43. Ekstrand MI, Terzioglu M, Galter D, Zhu S, Hofstetter C, Lindqvist E, *et al.* Progressive parkinsonism in mice with respiratory-chain-deficient dopamine neurons. *Proc Natl Acad Sci U S A*: 2007;104:1325–1330. <https://doi.org/10.1073/pnas.0605208103>.
44. Rygiel KA, Grady JP, Turnbull DM. Respiratory chain deficiency in aged spinal motor neurons. *Neurobiol Aging*: 2014;35:2230–2238. <https://doi.org/10.1016/j.neurobiolaging.2014.02.027>.
45. Byrska-Bishop M, Evani US, Zhao X, Basile AO, Abel HJ, Regier AA, *et al.* High-coverage whole-genome sequencing of the expanded 1000 Genomes Project cohort including 602 trios. *Cell*: 2022;185:3426–3440.e19. <https://doi.org/10.1016/j.cell.2022.08.004>.
46. Yang Z, Slone J, Wang X, Zhan J, Huang Y, Namjou B, *et al.* Validation of low-coverage whole-genome sequencing for mitochondrial DNA variants suggests mitochondrial DNA as a genetic cause of preterm birth. *Human Mutation*: 2021;42:1602–1614. <https://doi.org/10.1002/humu.24279>.
47. Durbin RM, Altshuler D, Durbin RM, Abecasis GR, Bentley DR, Chakravarti A, *et al.* A map of human genome variation from population-scale sequencing. *Nature*: 2010;467:1061–1073. <https://doi.org/10.1038/nature09534>.

48. Baeza JA. Yes, we can use it: a formal test on the accuracy of low-pass nanopore long-read sequencing for mitophylogenomics and barcoding research using the Caribbean spiny lobster *Panulirus argus*. *BMC Genomics*: 2020;21:882. <https://doi.org/10.1186/s12864-020-07292-5>.
49. Zheng-Bradley X, Flicek P. Applications of the 1000 Genomes Project resources. *Briefings in Functional Genomics*: 2017;16:163–170. <https://doi.org/10.1093/bfgp/elw027>.
50. Mukherjee S, Simon J, Bayuga S, Ludwig E, Yoo S, Orlov I, *et al.* Including additional controls from public databases improves the power of a genome-wide association study. *Hum Hered*: 2011;72:21–34. <https://doi.org/10.1159/000330149>.

**TABLE 1** Mitochondrial genome variants associated with increased ALS risk. Only variants with  $p$ -value  $<10^{-7}$  are shown. Logistic regression of 2,547 controls and 1,965 ALS cases. *MT* = mitochondrial, *HV* = hyper variable region, *CO* = Cytochrome *c* oxidases, *CYB* = Cytochrome *b*, *ND* = NADH ubiquinone oxidoreductases, *RNR1* = 12S ribosomal RNA. Nucleotide changes as shown in PLINK program. Variants annotated using NCBI Variation Viewer. (<https://www.ncbi.nlm.nih.gov/variation/view/>). 95% CI = 95% confidence interval. Risk = Relative Risk Ratio calculated as proposed by Zhang and collaborators

34

dbSNP ID	Position	Change	Gene	P-value	Odds Ratio	95% C.I.	Risk	95% CI
rs879158303	72	T>C	<i>MT-HV2</i>	$8.888 \times 10^{-8}$	3.842	2.346 - 6.292	3.75	2.32 - 6.02
rs28660704	497	C>T	<i>MT-HV3</i>	$4.783 \times 10^{-12}$	5.031	3.182 - 7.955	4.86	3.12 - 7.49
rs28358568	710	T>C	<i>MT-RNR1</i>	$1.897 \times 10^{-8}$	361.2	46.33 - 2815	281.56	44.73 - 877.06
rs28358587	3552	T>A	<i>MT-ND1</i>	$9.02 \times 10^{-27}$	$1.992 \times 10^4$	$3.255 \times 10^3$ - $1.219 \times 10^5$	35.80	35.48 - 35.85
rs28516468	6455	C>T	<i>MT-CO1</i>	$1.363 \times 10^{-10}$	63.44	17.87 - 225.2	22.15	11.88 - 29.28
rs28358875	7196	C>A	<i>MT-CO1</i>	$3.895 \times 10^{-13}$	136.5	36.19 - 514.8	24.00	16.32 - 27.42
rs878853022	9545	A>G	<i>MT-CO3</i>	$1.201 \times 10^{-9}$	39.24	12.02 - 128.1	17.35	8.81 - 24.67
rs1556424102	12405	C>T	<i>MT-ND5</i>	$7.097 \times 10^{-20}$	$3.322 \times 10^6$	$1.321 \times 10^5$ - $8.354 \times 10^7$	45.47	45.46 - 45.475
rs28359175	13263	A>G	<i>MT-ND5</i>	$3.154 \times 10^{-8}$	111.6	21.01 - 593.2	25.44	13.02 - 31
rs28357675	14318	T>C	<i>MT-ND6</i>	$2.886 \times 10^{-23}$	8969	1490 - $5.4 \times 10^4$	35.72	35.04 - 35.84
rs28357370	15487	A>T	<i>MT-CYB</i>	$6.983 \times 10^{-13}$	129.3	34.28 - 487.4	23.75	15.92 - 27.31
rs386829303	16297	T>C	<i>MT-HV1</i>	$5.838 \times 10^{-13}$	$3.994 \times 10^5$	$1.195 \times 10^4$ - $1.334 \times 10^7$	47.09	46.92 - 47.1
rs148377232	16298	T>C	<i>MT-HV1</i>	$3.139 \times 10^{-11}$	4.939	3.082 - 7.914	4.18	2.81 - 6.01

**TABLE 2** Mitochondrial genome variants associated with reduced ALS risk. Only variants with  $p$ -value  $< 10^{-7}$  are shown. Logistic regression of 2,547 controls and 1,965 ALS cases. *MT* = mitochondrial, *HV* = hyper variable region, *CO* = Cytochrome *c* oxidases, *CYB* = Cytochrome *b*, *ND* = NADH ubiquinone oxidoreductases, *ATP6* = mitochondrially encoded ATP synthase membrane subunit 6, *ATP8* = mitochondrially encoded ATP synthase membrane subunit 8, *RNR1* = 12S ribosomal RNA and *RNR2* = 16S ribosomal RNA. Nucleotide changes as shown in PLINK program. Variants annotated using NCBI Variation Viewer. (<https://www.ncbi.nlm.nih.gov/variation/view/>). CI = confidence interval. Risk = Relative Risk Ratio calculated as proposed by Zhang and collaborators<sup>34</sup>.

dbSNP ID	Position	Change	Gene	P-value	Odds Ratio	OR 95% CI	Risk	RR 95% CI
rs3937037	235	A>G	MT-HV2	4.447x10 <sup>-11</sup>	1.7x10 <sup>-3</sup>	2.55x10 <sup>-4</sup> - 1.133x10 <sup>-2</sup>	0.00	0 - 0.012
rs3901846	499	G>A	MT-HV3	3.575x10 <sup>-9</sup>	1.837x10 <sup>-2</sup>	4.837x10 <sup>-3</sup> - 6.927x10 <sup>-2</sup>	0.02	0.005 - 0.073
rs56489998	663	A>G	MT-RNR1	2.879x10 <sup>-16</sup>	2.118x10 <sup>-4</sup>	2.789x10 <sup>-5</sup> - 1.609x10 <sup>-3</sup>	0.00	0 - 0.002
rs397515731	980	T>C	MT-RNR1	1.969x10 <sup>-10</sup>	5.584x10 <sup>-2</sup>	2.296x10 <sup>-2</sup> - 1.358x10 <sup>-1</sup>	0.06	0.023 - 0.137
rs193303006	1736	A>G	MT-RNR2	3.238x10 <sup>-13</sup>	6.873x10 <sup>-4</sup>	9.684x10 <sup>-5</sup> - 4.877x10 <sup>-3</sup>	0.00	0 - 0.005
rs28358579	2352	T>C	MT-RNR2	4.92x10 <sup>-10</sup>	1.462x10 <sup>-1</sup>	7.98x10 <sup>-2</sup> - 2.68x10 <sup>-1</sup>	0.16	0.086 - 0.283
rs878907222	3741	C>T	MT-ND1	4.278x10 <sup>-12</sup>	2.942x10 <sup>-2</sup>	1.085x10 <sup>-2</sup> - 7.977x10 <sup>-2</sup>	0.03	0.011 - 0.08
rs9629042	3970	C>T	MT-ND1	2.551x10 <sup>-18</sup>	1.848x10 <sup>-6</sup>	9.536x10 <sup>-8</sup> - 3.58x10 <sup>-5</sup>	0.00	0
rs9326618	4248	T>C	MT-ND1	1.595x10 <sup>-11</sup>	1.388x10 <sup>-3</sup>	2.048x10 <sup>-4</sup> - 9.408x10 <sup>-3</sup>	0.00	0 - 0.01
rs15564229903	4824	A>G	MT-ND2	1.791x10 <sup>-16</sup>	1.961x10 <sup>-4</sup>	2.571x10 <sup>-5</sup> - 1.496x10 <sup>-3</sup>	0.00	0 - 0.002
rs200763872	4883	C>T	MT-ND2	1.233x10 <sup>-10</sup>	8.86x10 <sup>-2</sup>	4.235x10 <sup>-2</sup> - 1.854x10 <sup>-1</sup>	0.09	0.044 - 0.193
rs28357984	5178	C>A	MT-ND2	1.46x10 <sup>-10</sup>	7.716x10 <sup>-2</sup>	3.525x10 <sup>-2</sup> - 1.689x10 <sup>-1</sup>	0.08	0.037 - 0.176
rs879217723	5360	C>T	MT-ND2	4.278x10 <sup>-12</sup>	2.942x10 <sup>-2</sup>	1.085x10 <sup>-2</sup> - 7.977x10 <sup>-2</sup>	0.03	0.011 - 0.08
rs879043235	8137	C>T	MT-CO2	4.278x10 <sup>-12</sup>	2.942x10 <sup>-2</sup>	1.085x10 <sup>-2</sup> - 7.977x10 <sup>-2</sup>	0.03	0.011 - 0.08
rs28358884	8414	C>T	MT-ATP8	1.545x10 <sup>-8</sup>	1.029x10 <sup>-1</sup>	4.679x10 <sup>-2</sup> - 2.263x10 <sup>-1</sup>	0.11	0.049 - 0.233
rs201336180	8684	C>T	MT-ATP6	2.706x10 <sup>-12</sup>	2.906x10 <sup>-2</sup>	1.078x10 <sup>-2</sup> - 7.835x10 <sup>-2</sup>	0.03	0.011 - 0.079
rs2298007	8794	C>T	MT-ATP6	4.062x10 <sup>-16</sup>	2.247x10 <sup>-4</sup>	2.969x10 <sup>-5</sup> - 1.7x10 <sup>-3</sup>	0.00	0 - 0.002
rs3134801	9950	T>C	MT-CO3	1.876x10 <sup>-8</sup>	3.782x10 <sup>-2</sup>	1.208x10 <sup>-2</sup> - 1.184x10 <sup>-1</sup>	0.04	0.013 - 0.125
rs878969753	10142	C>T	MT-ND3	1.126x10 <sup>-11</sup>	3.074x10 <sup>-2</sup>	1.125x10 <sup>-2</sup> - 8.401x10 <sup>-2</sup>	0.03	0.011 - 0.084
rs41467651	10310	G>A	MT-ND3	5.501x10 <sup>-8</sup>	6.296x10 <sup>-3</sup>	1.012x10 <sup>-3</sup> - 3.916x10 <sup>-2</sup>	0.01	0.001 - 0.041
rs28358278	10400	C>T	MT-ND3	3.185x10 <sup>-53</sup>	2.007x10 <sup>-2</sup>	1.219x10 <sup>-2</sup> - 3.306x10 <sup>-2</sup>	0.03	0.016 - 0.043
rs28358283	10819	A>G	MT-ND4	5.658x10 <sup>-11</sup>	1.809x10 <sup>-1</sup>	1.085x10 <sup>-1</sup> - 3.017x10 <sup>-1</sup>	0.19	0.113 - 0.312
rs2853497	12007	G>A	MT-ND4	1.594x10 <sup>-13</sup>	1.277x10 <sup>-1</sup>	7.394x10 <sup>-2</sup> - 2.206x10 <sup>-1</sup>	0.14	0.082 - 0.241
rs878955011	12618	G>A	MT-ND5	4.832x10 <sup>-10</sup>	1.08x10 <sup>-1</sup>	5.359x10 <sup>-2</sup> - 2.177x10 <sup>-1</sup>	0.11	0.054 - 0.221
rs386420001	12882	C>T	MT-ND5	7.383x10 <sup>-8</sup>	4.087x10 <sup>-3</sup>	5.513x10 <sup>-4</sup> - 3.029x10 <sup>-2</sup>	0.00	0.001 - 0.031
rs2853501	13105	A>G	MT-ND5	1.753x10 <sup>-8</sup>	1.796x10 <sup>-1</sup>	9.883x10 <sup>-2</sup> - 3.263x10 <sup>-1</sup>	0.20	0.113 - 0.359
rs879066842	13500	T>C	MT-ND5	6.433x10 <sup>-10</sup>	9.579x10 <sup>-2</sup>	4.552x10 <sup>-2</sup> - 2.016x10 <sup>-1</sup>	0.10	0.046 - 0.204
rs28359184	13928	G>C	MT-ND5	3.003x10 <sup>-8</sup>	6.307x10 <sup>-2</sup>	2.373x10 <sup>-2</sup> - 1.676x10 <sup>-1</sup>	0.07	0.025 - 0.177
rs28357672	14212	T>C	MT-ND6	1.757x10 <sup>-10</sup>	1.39x10 <sup>-1</sup>	7.586x10 <sup>-2</sup> - 2.549x10 <sup>-1</sup>	0.15	0.079 - 0.264
rs386420019	14569	G>A	MT-ND6	2.422x10 <sup>-12</sup>	7.852x10 <sup>-2</sup>	3.854x10 <sup>-2</sup> - 1.6x10 <sup>-1</sup>	0.08	0.039 - 0.163
rs28357678	14668	C>T	MT-ND6	1.056x10 <sup>-8</sup>	1.52x10 <sup>-1</sup>	7.973x10 <sup>-2</sup> - 2.898x10 <sup>-1</sup>	0.16	0.083 - 0.299
rs28357680	14783	T>C	MT-CYB	1.036x10 <sup>-53</sup>	2.274x10 <sup>-2</sup>	1.406x10 <sup>-2</sup> - 3.677x10 <sup>-1</sup>	0.03	0.018 - 0.433
rs193302985	15043	G>A	MT-CYB	9.625x10 <sup>-59</sup>	4.309x10 <sup>-2</sup>	2.943x10 <sup>-2</sup> - 6.309x10 <sup>-2</sup>	0.06	0.039 - 0.082
rs193302991	15301	G>A	MT-CYB	2.787x10 <sup>-84</sup>	2.61x10 <sup>-2</sup>	1.808x10 <sup>-2</sup> - 3.769x10 <sup>-2</sup>	0.05	0.032 - 0.065
rs2853817	16172	T>C	MT-HV1	2.132x10 <sup>-8</sup>	3.969x10 <sup>-1</sup>	2.873x10 <sup>-1</sup> - 5.485x10 <sup>-1</sup>	0.42	0.306 - 0.571
rs35134837	16217	T>C	MT-HV1	8.59x10 <sup>-10</sup>	3.452x10 <sup>-2</sup>	1.177x10 <sup>-2</sup> - 1.012x10 <sup>-1</sup>	0.04	0.012 - 0.107
rs138126107	16261	C>T	MT-HV1	2.04x10 <sup>-9</sup>	3.492x10 <sup>-1</sup>	2.476x10 <sup>-1</sup> - 4.926x10 <sup>-1</sup>	0.36	0.256 - 0.503
rs879067317	16318	A>T	MT-HV1	1.887x10 <sup>-10</sup>	3.008x10 <sup>-2</sup>	1.024x10 <sup>-2</sup> - 8.841x10 <sup>-2</sup>	0.03	0.01 - 0.089

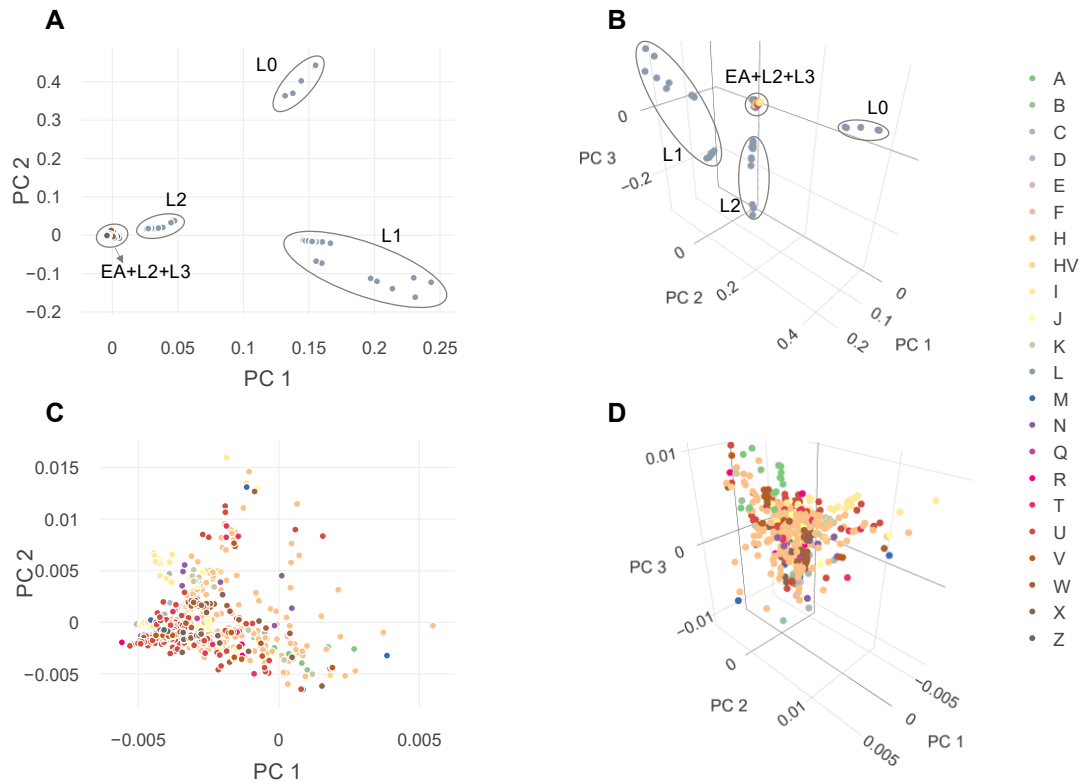
**TABLE 3** Intra haplogroup GWAS of SNVs associated with ALS. Only variants with  $p$ -value  $< 10^{-5}$  are shown. The  $p$ -values were calculated by the  $X^2$  test. *MT* = mitochondrial, *HV2* = hyper variable region 2, *HV3* = hyper variable region 3, *ND* = NADH ubiquinone oxidoreductases, *MT-TT* = mitochondrially encoded tRNA threonine, *RNR1* = 12S ribosomal RNA and *RNR2* = 16S ribosomal RNA. Variants were annotated using the MITOMAP (<https://www.mitomap.org/MITOMAP>) numbers confirmed in dbSNP build 152 ([www.ncbi.nlm.nih.gov/snp/](http://www.ncbi.nlm.nih.gov/snp/)). C.I. = 95% confidence interval of the Odds Ratio. Hap. = Haplogroup. (\*) = the odds ratios and C.I. of these variants were calculated with Haldane-Anscombe correction. (\*\*) ALS-associated variants also detected in the analysis of 4,512 samples dataset (Table 1).

Hap.	dbSNP ID	Position	Change	Gene	P-value	Odds Ratio	C.I.
L	-	76	C>T	<i>MT-HV2</i>	$1.175 \times 10^{-6}$	30.15	3.334 – 272.7
L	rs879040416	198	C>T	<i>MT-HV2</i>	$4.605 \times 10^{-19}$	99.18	12.73 – 772.5
L	-	*517	del A	<i>MT-HV3</i>	$3.06 \times 10^{-6}$	52.02	20.33 – 133.12
L	rs28358568	**710	T>C	<i>MT-RNR1</i>	$2.518 \times 10^{-22}$	118.6	15.39 – 913.5
L	rs28619217	2755	A>G	<i>MT-RNR2</i>	$4.766 \times 10^{-7}$	15.37	3.778 – 62.54
L	rs41423746	*3348	A>G	<i>MT-ND1</i>	$3.06 \times 10^{-6}$	52.02	20.33 – 133.12
L	rs1556423778	10196	C>T	<i>MT-ND3</i>	$3.175 \times 10^{-6}$	30.15	3.334 – 272.7
L	rs1569484396	*10793	C>T	<i>MT-ND4</i>	$3.06 \times 10^{-6}$	52.02	20.33 – 133.12
L	rs1603224091	13209	C>T	<i>MT-ND5</i>	$1.417 \times 10^{-6}$	19.03	3.638 – 99.52
L	-	13608	T>C	<i>MT-ND5</i>	$3.175 \times 10^{-6}$	30.15	3.334 – 272.7
L	-	*15940	del T	<i>MT-TT</i>	$8.239 \times 10^{-36}$	396.81	162.97 – 966.17
U	rs368463610	*234	A>G	<i>MT-HV2</i>	$7.194 \times 10^{-6}$	$2.229 \times 10^{-1}$	$9.402 \times 10^{-3}$ - $5.558 \times 10^{-2}$
U	rs878939965	*5186	A>T	<i>MT-ND2</i>	$6.1112 \times 10^{-6}$	$1.96 \times 10^{-2}$	$8.07 \times 10^{-3}$ – $4.77 \times 10^{-1}$
U	rs1556424235	13194	G>A	<i>MT-ND5</i>	$6.112 \times 10^{-6}$	$4.028 \times 10^{-2}$	$5.38 \times 10^{-3}$ – $3.017 \times 10^{-1}$

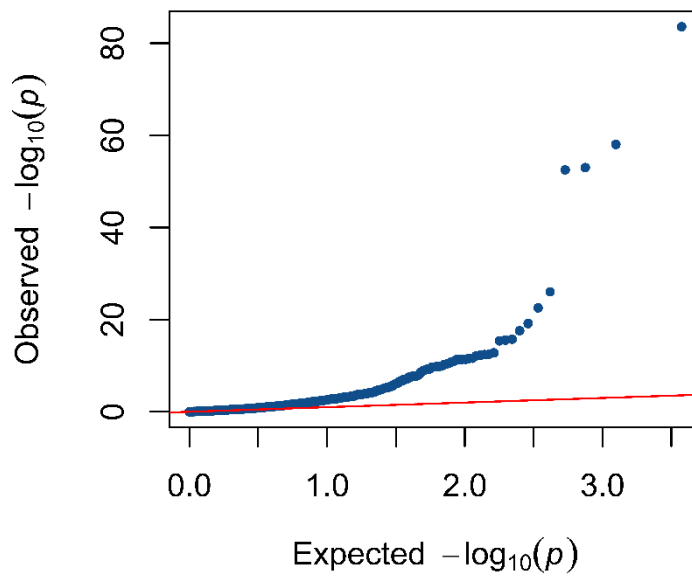
**Table 4.** Frequency of bulbar ALS cases by haplogroup. The % in the sample corresponds to the total of cases per haplogroups the total of the sample. A total of 1460 samples with clinical data were evaluated, 247 of which had a bulbar onset.

Samples/Type	Haplogroup			remaining haplogroups
	H	L	U	
Total in bulbar	114	12	24	97
Total in sample	740	69	132	519
% in bulbar	46.15	4.86	9.72	39.27
% in sample	50.68	4.73	9.04	35.55

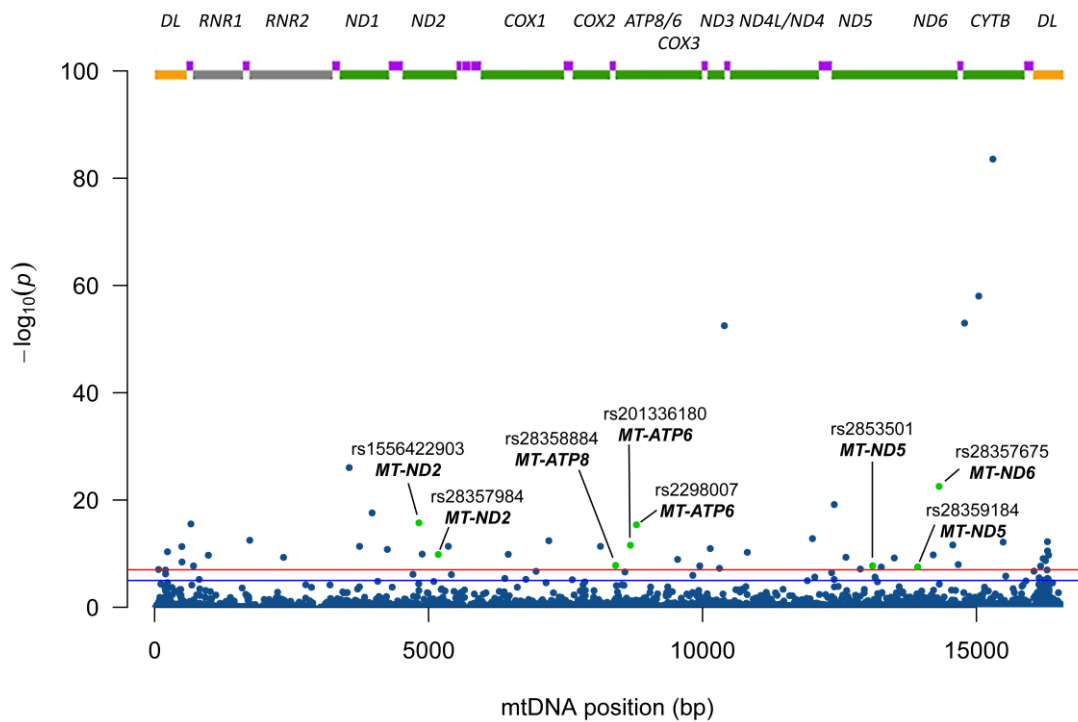




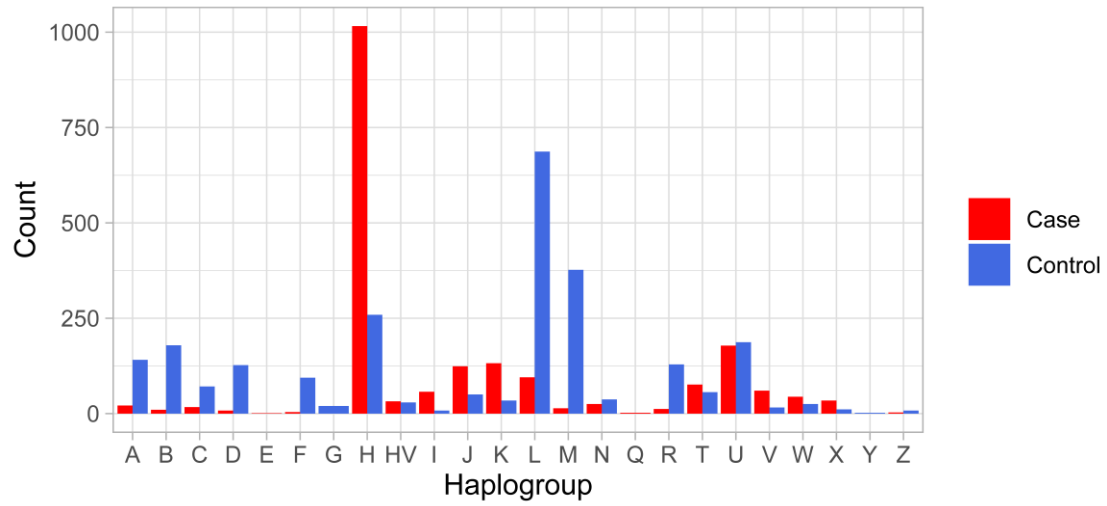
**FIGURE 1** Principal Component Analysis (PCA) of 1,965 ALS cases in this study. The color scale indicates the mitochondrial haplogroups. (A) the 2D PCA, (B) the corresponding 3D PCA with a third component, (C) a closer look at the non-African cluster and (D) the detailed 3D zooms in of the European-Asian haplogroups in the ALS samples. African (L) and Eurasian (EA) haplogroups are highlighted (ovals) in A and B.



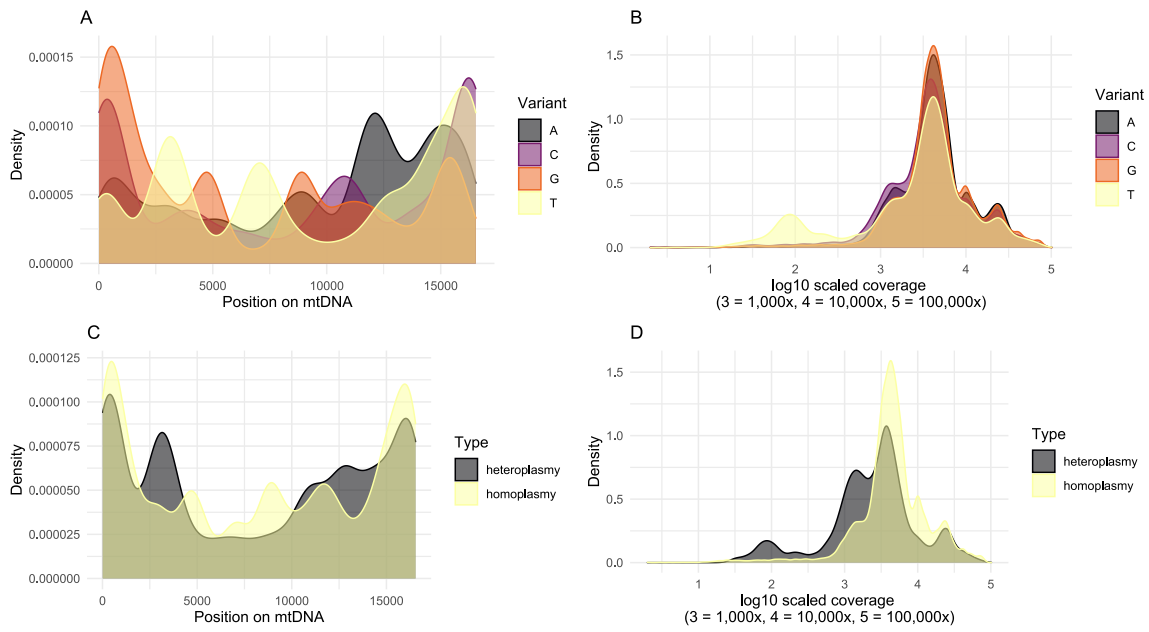
**FIGURE 2** The QQ-Plot with logistic regression of the 1,946 SNVs selected after removal of SNVs with genotype failure rate  $>0.1$  and Minimum Allele Frequency (MAF)  $<0.1\%$ . The red line is the expected random association values. This corresponds to the analysis of 4,512 samples, being 1,965 cases and 2,547 controls.



**FIGURE 3** Manhattan plot of the 1,946 SNVs mapped to the mitochondrial genome. In the abscissa are the mitochondrial positions numbered according to the revised Cambridge Reference Sequence (rCRS)<sup>24</sup> and in the ordinate the  $-\log_{10}(p)$ -values. The horizontal blue line indicates the  $10^{-5}$  threshold and the red horizontal line the  $10^{-7}$  threshold. The gene map is depicted at the top with yellow boxes representing the D-loop, green boxes the protein coding genes, gray boxes the rRNA genes and purple boxes the tRNA genes. Nonsynonymous substitutions are depicted as green dots and annotated.



**FIGURE 4** Plot of the mitochondrial haplogroup frequencies in 1,965 cases and 2,547 controls. The haplogroups are depicted in the abscissa and the number of samples in each category are depicted in the ordinate.



**FIGURE 5** Heteroplasmy analysis (1,978 samples). Density plots show the distribution of variants by position (A), distribution of sequencing coverage by variant (B), distribution of variant type by position (C) and distribution of sequencing coverage by variant type (D).