

1     **The First complete Zoroastrian-Parsi Mitochondria Reference Genome: Implications of**  
2                   **mitochondrial signatures in an endogamous, non-smoking population**

3  
4  
5     **Authors Names and Affiliations:**

6     Villoo Morawala Patell\*<sup>1,2,3</sup>, Naseer Pasha<sup>1&2</sup>, Kashyap Krishnasamy<sup>1&2</sup>, Bharti Mittal<sup>1&2</sup>,  
7     Chellappa Gopalakrishnan<sup>1&2</sup>, Raja Mugasimangalam<sup>1&4</sup>, Naveen Sharma<sup>1&2</sup>, Arati-Khanna  
8     Gupta<sup>1</sup>, Perviz Bhote-Patell<sup>1</sup>, Sudha Rao<sup>1&4</sup>, Renuka Jain<sup>1&2</sup>, The Avestagenome Project<sup>®</sup>

9  
10  
11                   <sup>1</sup>*Avesthagen Limited, Bangalore, India*

12                   <sup>2</sup>*The Avestagenome Project<sup>®</sup> International Pvt Ltd, Bangalore, Karnataka, India-*  
13                   <sup>560005</sup>

14                   <sup>3</sup>*AGENOME LLC, USA*

15                   <sup>4</sup>*Genotypic Technologies Private Limited, Bangalore 560094*

16  
17  
18  
19     \**Corresponding Author:*

20     Address correspondence to Villoo Morawala Patell, THE dry lab, Avesthagen Limited

21     Yolee Grande, 2nd Floor, Pottery Road, Richard's Town, Bangalore, 560005, Karnataka, India,

22     Email: viloo@avesthagen.com

23

24 **Abstract:**

25 The present-day Zoroastrian-Parsis have roots in ancient pastoralist migrations from circumpolar  
26 regions<sup>1</sup> leading to their settlement on the Eurasian Steppes<sup>2</sup> and later, as Indo Iranians in the  
27 Fertile Crescent<sup>3</sup>. From then, the Achaemenids (550 - 331 BC), and later the Sassanids (224 BC -  
28 642 AD) established the mighty Persian Empires<sup>2</sup>. The Arab invasion of Persia in 642 AD  
29 necessitated the migration of Zoroastrians from Pars to India where they settled as Parsis and  
30 practiced their faith, Zoroastrianism. Endogamy became a dogma, and the community has  
31 maintained the practice since their arrival in India. Fire is the medium of worship<sup>4</sup> as it is  
32 considered pure and sacrosanct; Social ostracism practiced against smokers resulted in a non-  
33 smoking community, thus forming a unique basis for our study.

34 In order to gain a clearer understanding of the historically recorded migration of the Zoroastrian-  
35 Parsis, decipher their phylogenetic relationships and understand disease association to their  
36 individual mitochondrial genomes, we generated the first complete *de novo* Zoroastrian-Parsi  
37 Mitochondrial Reference Genome, AGENOME-ZPMS-HV2a-1. Phylogenetic analysis of  
38 additional 100 Parsi mitochondrial genome sequences, showed their distribution into 7 major  
39 haplogroups and 25 sub-haplogroups and a largely Persian origin for the Parsi community. We  
40 have generated individual reference genomes for each major haplogroup and assembled the  
41 Zoroastrian Parsi Mitochondrial Consensus Genome (AGENOME-ZPMC V1.0) for the first  
42 time in the world.

43 We report 420 variants, specifically 12 unique mitochondrial variants in the 100 mitochondrial  
44 genome sequences compared with the revised Cambridge Reference Sequence (rCRS) standard.

45 Disease association mapping showed 217 unique variants linked to longevity and 41 longevity  
46 associated disease phenotypes across most haplogroups. Our results indicate none of the variants  
47 are linked to lung cancer. Mutational signatures, C>A, G>T transitions<sup>36</sup>, linked to tobacco  
48 carcinogens were found at extremely low frequencies in the Zoroastrian-Parsi cohort.

49 Our analysis of gene-coding, tRNA and the D-Loop regions revealed haplogroup specific disease  
50 associations for Parkinson's, Alzheimer's, Cancers, and Rare diseases.

51 These disease signatures investigated in the backdrop of generations of endogamy, in the rapidly  
52 declining, endangered Zoroastrian-Parsi community of India, provides exceptional universal  
53 opportunity to understand and mitigate disease.

54 **Keywords:** Mitochondria, Haplotypes, Phylogeny, Human migration, Endogamous, Non-  
55 smoking, Longevity, Cancers, Neurodegenerative disorders, Rare Diseases, t-RNA, D-loop  
56 variants, Population genetics, Unique mitochondrial variants, Zoroastrian Parsis, Persia, Iran,  
57 India, Precision healthcare.

58 **Abbreviations:** mt DNA-Mitochondrial DNA ; rCRS-revised Cambridge Reference Sequence ;  
59 NGS Next Generation Sequencing ; ZPMS- Zoroastrian Parsi Mitochondrial Sequence ; ZPMRG-  
60 Zoroastrian Parsi Mitochondrial Reference Genome ; ZPMCG- Zoroastrian Parsi Mitochondrial  
61 Consensus Genome ; AD- Alzheimers Disease ; PD- Parkinsons Disease

62

63 **Introduction :**

64

65 ***The Burden of History- Travelogue of the Zoroastrian Mitochondrion***

66 Zoroastrian-Parsis of India are followers of the ancient prophet Zarathushtra, claimed by the Greek  
67 historian Herodotus to have been born circa 6,450 BC<sup>1</sup>. Zarathushtra, advocated the first known  
68 monotheistic concept of one supreme intelligence termed Ahura Mazda - ‘Majestic Creator’<sup>2</sup>.

69 The ancient homeland of the present-day Zoroastrian-Parsis finds mention in their sacred Avestan  
70 text *Vendidad*, and the location indicated is the North Polar Arctic region<sup>3</sup>. Sanskrit scholar B G  
71 Tilak’s study *Arctic Home in the Vedas* is also corroborated by Bennet, suggesting that the Indo-  
72 European culture originated in the Hyperborean Regions of Northern Siberia and the islands of the  
73 circumpolar regions<sup>4</sup>.

74 Around 12000 years ago, this region suffered a natural calamity and became ice-clad<sup>1</sup> necessitating  
75 southward migrations of these pastoralist inhabitants, and by 4,000 BC the Indo Europeans took  
76 over the Eurasian Steppe<sup>7</sup>.

77 From the late second to early first millennium BC, the Indo Europeans, mostly on the basis of  
78 religious worship, split with the Indo Aryans who moved further south and crossed the Hindu  
79 Kush, while the Indo Iranians (Medes, Persians, and Parthians) began populating the western  
80 portion of the Iranian plateau, close to the Alborz and Zagros Mountains and northern  
81 Mesopotamia to Southeast Anatolia, in what is called the Fertile Crescent where significant  
82 innovation in agriculture occurred<sup>8</sup>.



83

84 Image of *Taq-e Bostan* means "Arch of the Garden" or "Arch made by stone" , a site with a series of large rock reliefs from the era of the  
85 Sassanid Empire of Persia (Iran), carved around the 4th century CE. Image courtesy of Irandestination.com

86 In 550 BC, Cyrus the Great overthrew the leading Median rule and conquered the Kingdom of  
87 Lydia and the Babylonian Empire after which he established the Persian Zoroastrian Achaemenid  
88 Empire (the First Persian Empire), while his successors Dariush I (522-485 BC) Xerxes I,  
89 Artaxerxes and others extended its borders to encompass several countries across three continents,  
90 namely Europe, Africa and Asia. A second Zoroastrian dynasty of Sassanian Kings followed, who  
91 ruled Persia starting with Ardashir I (224 BC). It was the Golden age of the Persian Empire. During  
92 the time of Zoroastrian Achaemenid and Sassanid empires, Persia became a global hub of culture,  
93 religion, science, art, gender equality, and technology.

94 The Persians under Yezdecard III were defeated by the Arabs in two decisive battles – (Qadisiyah-  
95 636 AD and Nahavand – 642 AD) resulting in the fall of the Zoroastrian Persian Empire.

96 It was almost a hundred years later in the 8<sup>th</sup> century that a few boatloads of Zoroastrians left Paars  
97 and Khorasan from the port of Hormuz to sail south towards India. The boats first touched shore  
98 on Diu island on the west coast of India where the refugees stayed for around 19 years. The  
99 environment being non-conducive to progress, they once again set sail and arrived in Sanjan,  
100 Gujarat. Vijayaditya of the Chalukya dynasty (aka Jadi Rana) the ruler, hesitated to give refuge,  
101 but on being explained the principles of Zoroastrianism and observing the similarities with the  
102 Vedic religion, the Parsis were given refuge.



103

104

Map showing early migration of Parsis from Iran. Image courtesy Microsoft Encarta. Reference Library 2005

105 Endogamy became the norm to preserve their identity, and for the last 1300 years the community  
106 has maintained this practice<sup>9,10</sup>. Fire being the purest of all elements is considered sacred by  
107 Zoroastrian-Parsis. Strict measures are employed to maintain the purity of fire, hence the strict  
108 social ostracism practiced against smokers in the community.

109 Today, the Zoroastrian-Parsis, are a small community of <52000 in India (2011 Census, Govt of  
110 India). We present the genetic data of the conserved Zoroastrian-Parsi mitochondrion,  
111 encapsulated in resilience of thousands of years of magnificent history: of struggles and  
112 overcoming them; of building something out of nothingness; of achievement gained with ethical  
113 standards; and philanthropy.

114 In recent decades, the analysis of the variability of maternally inherited mitochondrial DNA  
115 (mtDNA) has been commonly used to reconstruct the history of ethnic groups, especially with  
116 respect to maternal inheritance. The lack of genetic recombination in mtDNA, results in the  
117 accumulation of maternally inherited single nucleotide polymorphisms (variants). The  
118 accumulation of Variants along maternally inherited lineages results in phylogenetically traceable  
119 haplotypes<sup>15</sup> which can be used to follow maternal genealogies both historically and  
120 geographically. This approach has provided insightful findings into the origins and disease  
121 etiologies associated with another well documented endogamous European community: the  
122 Icelandic people (rev in <sup>13</sup>).

123

124 Human mtDNA (mitochondrial DNA) is a double stranded, circular (16,569 kb) genome of  
125 bacterial origin<sup>16</sup> primarily encoding vital subunits of the energy generating oxidative  
126 phosphorylation and electron transport chain (ETC) pathway that generates Adenosine Tri-  
127 Phosphate (ATP), the primary energy substrate of the eukaryotic cell. In addition, 22 tRNAs and 2  
128 rRNAs are also encoded by the mtDNA<sup>17</sup>.

129 In this study, our first aim was to gain a clear understanding and genetic impact of the historically  
130 recorded migration of the Zoroastrian-Parsis from Persia to India, and to link socio-cultural,  
131 ritualistic practices followed within the community over several millenia manifesting in genetic  
132 outcomes. To shed further light on the impact of migrations followed by integration into  
133 communities, where ritual and social practices are strictly followed within communities and  
134 between communities resulting in specific traceable signatures.

135

136 Secondly, we have attempted to elucidate the genetic basis of commonly occurring diseases in this  
137 endogamous community. To address these issues, we generated the first complete *de novo*  
138 Zoroastrian-Parsi Mitochondrial Genome (AGENOME-ZPMS-HV2a-1) and used it to arrive at  
139 the mitochondrial haplotype specific Reference Genomes from a hundred Zoroastrian-Parsi  
140 individuals. Our study for the first time, has assembled the Zoroastrian Parsi Mitochondrial  
141 Consensus Genome (AGENOME-ZPMCG V 1.0) thereby creating the first Mitochondrial  
142 Consensus Genome for the Zoroastrian-Parsi community.

143

144 Our phylogenetic analysis confirmed that present day Zoroastrian-Parsis are closely related to  
145 Persians, and like most endogamous communities have comparatively lower genetic diversity and  
146 tend to be predisposed to several inherited genetic disorders<sup>5,11</sup>. They also possess longevity as a  
147 trait and are a long-living community<sup>6</sup> with lower incidences of lung cancer<sup>12</sup>. The study of the  
148 genealogic history of a close-knit community like the Parsis provides an unique opportunity, to  
149 understand the link between disease and social behaviour, thus providing the direction for  
150 population genetics as a basis for personalized healthcare.

151

## 152 **Materials and Methods**

### 153 **Sample collection and ethics statement**

154 One hundred healthy nonsmoking Parsi volunteers residing in the cities of Hyderabad-  
155 Secunderabad and Bangalore, India were invited to attend blood collection camps at the  
156 Zoroastrian centers in their respective cities under the auspices of The Avestagenome Project<sup>TM</sup>.  
157 Each adult participant (>18 years) underwent height and weight measurements and answered an  
158 extensive questionnaire designed to capture their medical, dietary and life history. All subjects  
159 provided written informed consent for the collection of samples and subsequent analysis. All  
160 health-related data collected from the cohort questionnaire were secured in The Avestagenome  
161 Project<sup>TM</sup> database to ensure data privacy. All procedures performed in this study involving human  
162 participants were in accordance with the ethical standards of the institution (Avesthagen Limited,  
163 Bangalore, India) and in line with the 1964 Helsinki declaration and its later amendments. This  
164 study has been approved by the Avesthagen Ethics Committee (BLAG-CSP-033).

165

### 166 **Genomic DNA extraction**

167 Genomic DNA from the buffy coat of peripheral blood was extracted using the Qiagen Whole  
168 Blood and Tissue Genomic DNA Extraction kit (cat. #69504). Extracted DNA samples were

169 assessed for quality using the Agilent Tape Station and quantified using the Qubit™ dsDNA BR  
170 Assay kit (cat. #Q32850) with the Qubit 2.0® fluorometer (Life Technologies™). Purified DNA  
171 was subjected to both long-read (Nanopore GridION-X5 sequencer, Oxford Nanopore  
172 Technologies, Oxford, UK) and short-read (Illumina sequencer) for sequencing.

173

#### 174 **Library preparation for sequencing on the Nanopore platform**

175 Libraries of long reads from genomic DNA were generated using standard protocols from Oxford  
176 Nanopore Technology (ONT) using the SQK-LSK109 ligation sequencing kit. Briefly, 1.5 µg of  
177 high-molecular-weight genomic DNA was subjected to end repair using the NEBNext Ultra II End  
178 Repair kit (NEB, cat. #E7445) and purified using 1x AmPure beads (Beckman Coulter Life  
179 Sciences, cat. #A63880). Sequencing adaptors were ligated using NEB Quick T4 DNA ligase (cat.  
180 #M0202S) and purified using 0.6x AmPure beads. The final libraries were eluted in 15 µl of elution  
181 buffer. Sequencing was performed on a GridION X5 sequencer (Oxford Nanopore Technologies,  
182 Oxford, UK) using a SpotON R9.4 flow cell (FLO-MIN106) in a 48-hr sequencing protocol.  
183 Nanopore raw reads (fast5 format) were base called (fastq5 format) using Guppy v2.3.4 software.  
184 Samples were run on two flow cells and generated a dataset of ~14 GB.

185

#### 186 **Library preparation and sequencing on the Illumina platform**

187 Genomic DNA samples were quantified using the Qubit fluorometer. For each sample, 100 ng of  
188 DNA was fragmented to an average size of 350 bp by ultrasonication (Covaris ME220  
189 ultrasonicator). DNA sequencing libraries were prepared using dual-index adapters with the  
190 TruSeq Nano DNA Library Prep kit (Illumina) as per the manufacturer's protocol. The amplified  
191 libraries were checked on Tape Station (Agilent Technologies), quantified by real-time PCR using  
192 the KAPA Library Quantification kit (Roche) with the QuantStudio-7flex Real-Time PCR system  
193 (Thermo). Equimolar pools of sequencing libraries were sequenced using S4 flow cells in a  
194 Novaseq 6000 sequencer (Illumina) to generate 2 x 150-bp sequencing reads for 30x genome  
195 coverage per sample.

196

#### 197 **Generation of the *de novo* Parsi mitochondrial genome (AGENOME-ZPMS-HV2a-1)**

198 a) Retrieval of mitochondrial reads from whole-genome sequencing (WGS) data:

199 A total of 16 GB of raw data (.fasta) was generated from a GridION-X5 Nanopore sequencer for  
200 AGENOME-ZPMS-HV2a-1 from WGS. About 320 million paired-end raw reads were generated  
201 for AGENOME-ZPMS-HV2a-1 by Illumina sequencing.

202 Long Nanopore reads (. fastaq5) were generated from the GridION-X5 samples. The high-quality  
203 reads were filtered (PHRED score =>20) and trimmed for adapters using Porechop (v0.2.3). The  
204 high-quality reads were then aligned to the human mitochondrial reference (rCRS) NC\_12920.1  
205 using Minimap2 software. The aligned SAM file was then converted to a BAM file using  
206 SAMtools. The paired aligned reads from the BAM file were extracted using Picard tools (v1.102).

207

208 The short Illumina high-quality reads were filtered (PHRED score  $\Rightarrow$ 30). The adapters were  
209 trimmed using Trimgalore (v0.4.4) for both forward and reverse reads, respectively. The filtered  
210 reads were then aligned against a human mitochondrial reference (rCRS<sup>21</sup>) using the Bowtie2  
211 (v2.2.5) aligner with default parameters. The mapped SAM file was converted to a BAM file using  
212 SAMtools, and the mapped paired reads were extracted using Picard tools (v1.102).

213

214 b) *De novo* mitochondrial genome assembly

215 Mapped reads were used for *de novo* hybrid assembly using the Maryland Super-Read Celera  
216 Assembler (MaSuRCA-3.2.8) tool. The configuration file from the MaSuRCA tool was edited by  
217 adding appropriate Illumina and Nanopore read files. The MaSuRCA tool uses a hybrid approach  
218 that has the computational efficiency of the de Bruijn graph methods and the flexibility of overlap-  
219 based assembly strategies. It significantly improves assemblies when the original data are  
220 augmented with long reads. AGENOME-ZPMS-HV2a-1 was generated by realigning the mapped  
221 mitochondrial reads from Illumina as well as Nanopore data with the initial assembly.

222

223 **Confirmation of Variants in the *de novo* Parsi mitochondrial genome using Sanger**  
224 **sequencing**

225 To validate the *de novo* Parsi mitochondrial sequence, AGENOME-ZPMS-HV2a-1, selected  
226 variants were identified and subjected to PCR amplification. Genomic DNA (20 ng) was PCR  
227 amplified using LongAmpTaq 2X master mix (NEB). The PCR amplicons of select regions were  
228 subjected to Sanger sequencing and BLAST analysis to confirm the presence of eight Variants  
229 using primers listed in Supplemental Table 1.

230

231 **Generation of the Zoroastrian Parsi mitochondrial consensus genome (AGENOME-**  
232 **ZPMCG-V1.0) and Parsi haplogroup-specific consensus sequences**

233 a) **Retrieving mitochondrial reads from 100 Parsi whole-genome sequences**

234 The whole-genome data from 100 Parsi samples were processed for quality assessment. The  
235 adapters were removed using the Trimgalore 0.4.4 tool for paired end reads (R1 and R2), and sites  
236 with PHRED scores less than 30 and reads shorter than 20 bp in length were removed. The  
237 processed Illumina reads were aligned against a human mitochondrial reference sequence (rCRS<sup>18</sup>,  
238 NC\_012920.1) using the Bowtie 2 (version 2.4.1) aligner with default parameters. Mapped reads  
239 were further used for the *de novo* assembly using SPAdes (version 3.11.1) and Velvet and IVA  
240 (version 1.0.8). Comparison of the assembly and statistics were obtained using Quast (version  
241 5.0.2). The assembled scaffolds were subjected to BLASTn against the NCBI non-redundant  
242 nucleotide database for validation.

243 b) **Variant calling and haplogroup classification**

244 Sequencing reads were mapped to the human mitochondrial genome (rCRS<sup>21</sup>) assembly using the  
245 MEM algorithm of the Burrows–Wheeler aligner (version 0.7.17-r1188) with default parameters.  
246 Variants were called using SAMtools (version 1.3.1) to transpose the mapped data in a sorted  
247 BAM file and calculate the Bayesian prior probability. Next, Bcftools (version 1.10.2) was used



248 to calculate the prior probability distribution to obtain the actual genotype of the variants detected.  
249 The classification and haplogroup assignment were performed for each of the 100 Parsi mtDNAs  
250 after variant calling and after mapping reference and alternate alleles to the standard haplogroups  
251 obtained from MITOMAP (**Appendix 4**).

### 252 **c) Haplogroup-based consensus sequence**

253 Ninety-seven of 100 full-length Parsi mtDNA sequences were segregated based on haplogroups  
254 and separately aligned using the MUSCLE program to obtain the multiple sequence alignments.  
255 The Zoroastrian Parsi Mitochondrial Reference Genome (ZPMRG) and the Parsi haplogroup-  
256 specific consensus sequences were generated after calculation of the ATGC base frequency by  
257 comparison of the nucleotides in an alignment column to all other nucleotides in the same column  
258 called for other samples at the same position. The highest frequency (%) was taken to build seven  
259 Parsi haplogroup ZPMRGs and the seven Parsi haplogroup-specific consensus sequences.

260

### 261 **Phylogeny build and analysis**

262 Ninety-seven of 100 full-length Parsi mtDNA sequences generated as described above were  
263 compared with 100 randomly chosen Indian mtDNA sequences derived from NCBI Genbank  
264 under the accession codes FJ383174.1-FJ 383814.1<sup>22</sup>, DQ246811.1-DQ246833.1<sup>23</sup>, and  
265 KY824818.1-KY825084.1<sup>24</sup> and from previously published data on 352 complete Iranian mtDNA  
266 sequences<sup>25</sup>. All mtDNA sequences were aligned using MUSCLE software<sup>26</sup> using the “maxiters  
267 2” and “diags 1” options, followed by manual verification using BioEdit (version 7.0.0). Following  
268 alignment, the neighbor-joining method, implemented in MEGAX<sup>27</sup>, was employed to reconstruct  
269 the haplotype-based phylogeny. The neighbor-joining method was used because it is more efficient  
270 for large data sets<sup>28</sup>.

271

### 272 **Variant disease analysis**

273 One hundred Parsi mitochondria sequences extracted from the WGS were uploaded into the  
274 VarDiG<sup>®</sup>-R search engine (<https://vardigrviz.genomatics.life/vardig-r-viz/>) on AmazonWeb  
275 Services. VarDiG<sup>®</sup>-R, developed by Genomatics Private Ltd, connects variants, disease, and  
276 genes in the human genome. Currently, the VarDiG<sup>®</sup>-R knowledgebase contains manually curated  
277 information on 330,000+ variants, >20 K genes covering >4500 phenotypes, including nuclear and  
278 mitochondrial regions for 150,000+ published articles from 388+ journals. Variants obtained from  
279 Parsi mitochondria were mapped against all the published variants in VarDiG<sup>®</sup>-R. Associations  
280 with putative diseases was ascertained for each variant through VarDIG<sup>®</sup>-R.

281

282 Seventeen tRNA SNP sites were identified in the 100 Parsi mitochondrial SNP data. The PON-  
283 mt-tRNA database<sup>42</sup> was downloaded to annotate the tRNA Variants for their impact and disease  
284 associations. This database employs a posterior probability-based method for classification of  
285 mitochondrial tRNA variations. PON-mt-tRNA integrates the machine learning-based probability  
286 of pathogenicity and the evidence-based likelihood of pathogenicity to predict the posterior

287 probability of pathogenicity. In the absence of evidence, it classifies the variations based on the  
288 machine learning-based probability of pathogenicity.

289

290 For annotation of disease pathways associated with Variants, we employed MitImpact  
291 (<https://mitimpact.css-mendel.it/>) to predict the functional impact of the nonsynonymous Variants  
292 on their pathogenicity. This database is a collection of nonsynonymous mitochondrial Variants  
293 and their functional impact according to various databases, including SIFT, Polyphen, Clinvar,  
294 Mutationtester, dbSNP, APOGEE, and others. The disease associations, functional classifications,  
295 and engagement in different pathways were determined using the DAVID and UNIPROT  
296 annotation tools.

297

### 298 **Haplogroup and disease linkage**

299 Principal component analysis (PCA) was performed to visualize the linkage of the haplogroup  
300 with disease. XLSTAT (Addinsoft 2020, New York, USA. <https://www.xlstat.com>) was used for  
301 statistical and data analysis, including PCA.

302

### 303 **Data Accessibility:**

304 The GenBank (<http://www.ncbi.nlm.nih.gov/genbank>) accession numbers for the 105 novel  
305 complete mtDNA sequences (97 ZPMS, 7 ZPMRG and 1 ZPMCG) reported in this paper are  
306 MT506242-MT506346. The raw reads for 97 ZPMS mitochondrial genome sequences have been  
307 deposited with BioProject ID: PRJNA636291. The SRA accession numbers for the 97 ZMPS:  
308 SRR11888826-SRR11888922.

309

## 310 **Results**

### 311 **Assembly of the first complete Zoroastrian Parsi mitochondrial sequence, AGENOME-** 312 **ZPMS-HV2a-1**

313 The first complete *de novo* non-smoking Zoroastrian Parsi mitochondrial sequence, AGENOME-  
314 ZPMS-HV2a-1, was assembled from a healthy Parsi female sample by combining the sequence  
315 data generated from two next-generation sequencing (NGS) platforms using a protocol, as outlined  
316 in Materials and Methods. Our approach combines the sequencing depth and accuracy of short-  
317 read technology (Illumina) with the coverage of long-read technology (Nanopore). QC parameters  
318 for mitochondrial reads, mitochondrial coverage, and X-coverage were found to be optimal, as  
319 seen in **Supplementary Figure 1**. The hybrid Parsi mitochondrial genome was assembled as a  
320 single contig of 16.6 kb (with 99.82% sequence identity), resulting in the consensus sequence for  
321 the *de novo* Parsi mitochondrial genome with 99.84% sequence identity to the revised Cambridge  
322 Reference Sequence (rCRS<sup>21</sup>).

323

### 324 **Identification of 28 unique Variants in AGENOME-ZPMS-HV2a-1**

325 The variants identified from both the Illumina and Nanopore data were considered to be significant  
326 for this *de novo* Zoroastrian Parsi mitochondrial genome, henceforth referred to as AGENOME-  
327 ZPMS-HV2a-1.

328

329 A total of 28 significant variants (i.e., variants) were identified by BLAST alignment between the  
330 Parsi mitochondrial hybrid assembly and the rCRS<sup>21</sup> (**Figure 1, Table 1**). To confirm the  
331 authenticity of the identified variants, we selected a total of 7 identified variants from the D-loop  
332 region and one SNP from the *COI* gene (m.C7028T, A375A) and subjected them to Sanger  
333 sequencing using primers. All 8 predicted variants were verified and confirmed for their presence  
334 in the consensus Parsi mitochondrial genome (**Figure 2**).

335

336 The majority (n=11) of the variants identified in the AGENOME-ZPMS-HV2a-1 were found in  
337 the hypervariable regions (HVRI and HVRII) of the D-loop. Of the remaining 17 variants, eight  
338 were found to represent synonymous variants, while four were in genes for 12S, 16S-rRNA (n=3)  
339 and tRNA (n=1) (**Figure 1**). The remaining 5 nonsynonymous variants were located in the genes  
340 for *ATPase6* (m8860G>A), *COIII* (m.9336 A>G), *ND4* (m.11016 G>A), and two in the *CytB* gene  
341 (m15326 A>G and m15792 T>C, (**Table 1**). Except for the *ATPase6* gene variant, which has been  
342 found to be associated with hypertrophic cardiomyopathy in Iranian individuals<sup>29</sup>, no associations  
343 were found in the published literature for these gene variants, and they need to be further  
344 investigated.

345

346 Given that the Zoroastrian Parsis are known to have originated in Persia and have practiced  
347 endogamy since their arrival on the Indian subcontinent, we wished to determine the mitochondrial  
348 haplogroup associated with the first complete Zoroastrian Parsi mitochondrial genome. We  
349 therefore compared the variants associated with ZPMS-HV2a-1 to standard haplogroups obtained  
350 from MITOMAP and determined the haplogroup to be HV2a (**Figure 1**). This haplogroup is  
351 known to have originated in Iran<sup>25</sup>, suggesting Persian origins for this Parsi individual, based on  
352 maternal inheritance patterns.

353

### 354 **Seven major haplogroups identified in 100 Zoroastrian Parsi individuals**

355 Keeping in mind the endogamous nature of the Indian Parsis and to understand the extent of the  
356 diversity of the mitochondrial haplogroups in this population, we analyzed mitochondrial genomes  
357 from 100 consenting Parsi individuals. Our study had an equal representation of both genders, and  
358 60% of the subjects were of age 30–59 (mean age 50±1.6) (**Figure 3**). Complete analysis of the  
359 variants in the 100 Parsi samples identified a total of 420 unique Variants (**Figure 4, Appendix**  
360 **1**). QC analysis of the 100 mitochondrial genomes sequenced were found to be optimal:  
361 PHRED>30 (**Supplementary Figure 2**). Variant distribution in the coding region normalized to  
362 gene length showed the *ND6* gene has the highest number of variants (**Supplementary Figure 3**).  
363 The 100 Zoroastrian Parsi mitochondrial genomes were subjected to haplogroup analysis using  
364 haplogroup specific variant assignment matrix from MITOMAP (**Appendix 4**). The haplogroup  
365 assignment based on variants classified the genomes into seven principal haplogroups (HV, U, T,  
366 M, A, F, and Z) and 25 sub-haplogroups were also identified within the principal haplogroups  
367 (**Figure 5**). The variant count across all sub-haplogroups ranged between 14-64 (**Figure 6A**).

368 Analysis of the sub-haplogroups demonstrated that HV2a was the single largest representative sub-  
369 haplogroup within the Parsi population (n=14, n=9 females, n=5 males, (**Figure 6B**), that includes  
370 the AGENOME-ZPMS-HV2a-1.

371  
372 The sub-haplogroup HV2a (n=14 subjects) contained 28 variants observed in the AGENOME-  
373 ZPMS-HV2a-1 are common across all 14 subjects. In total, the HV2a sub-haplogroup had 38  
374 variants, with the highest number in the HVR II region (n=8). Coding region mutations constituted  
375 20/38 variants, with equal distribution between synonymous (n=10) and non-synonymous  
376 substitutions observed for this sub-haplogroup (n=10). Among the coding regions, the largest  
377 number of Variants was found in the gene encoding *COI* (n=6, **Supplementary Figure 4A**). Four  
378 *COI* Variants distributed across all of the 14 subjects in the HV2a sub-haplogroup (m.6104 C>T,  
379 m.6179 G>A, m.7028 C>T, and m.7193 T>C) constitute synonymous mutations (amino acid  
380 change: F67F, M92M, A375A, and F430F, respectively). Two Variants (m.7080 T>C and m.7146  
381 A>G), found to occur in one subject each in the sub-haplogroup HV2a, were nonsynonymous  
382 substitutions (F393L and T415A, respectively). Further analysis of rare Variants (occurring only  
383 in single subjects or n<8/14) showed their presence in the 16S-RNR2 gene (m.1883 G>A and  
384 m.1888 G>A), as well as the *COII*, *COIII* (m.8203 C>T and m.9540 T>C), and HVR I (m.16153  
385 G>A and 16274 G>A) genes, which were synonymous substitutions in these coding genes, while  
386 we found nonsynonymous substitutions in the *COII* (m.7650 C>T; T22I), *ND5* (m.12358 C>T;  
387 T8A), and *CYTB* (m.14954 A>G; T70A) genes in our analysis. We found a variant in the gene  
388 encoding for tRNA[R] at m.10410 T>C (n=14 subjects), but no mutations were observed in the D-  
389 loop region for the entire group under analysis.

390  
391 The sub-haplogroup HV12b (n=1 subject) contained 17 Variants. HVR II harbors four Variants,  
392 while the coding genes together contain six Variants that encode three synonymous and three  
393 nonsynonymous substitutions. We observed Variants encoding nonsynonymous substitutions in  
394 this sub-haplogroup in *ATPase6* (m.8860 A>G; T112A), *ND5* (m.13889 G>A; C518Y), and *CYTB*  
395 (m.15326 A>G; T194A). Three Variants were found in 12S-RNR1, two Variants in 16S-RNR2.  
396 In the non-coding regions 5 variants were observed in HVR II, 1 in HVR I and 1 in the D-loop  
397 region (m.16519 T>C). No Variants were observed in the genes coding for tRNAs in the HV12b  
398 sub-haplogroup.

399  
400 The 21 subjects analyzed that fell into the U haplogroup consisted of four sub-haplogroups U1a  
401 (n=1), U4b (n=11), U2e (n=3), and U7a (n=6). The U1a sub-haplogroup contained 44 Variants  
402 distributed across 19 positions in the mitochondrial genome. Twenty-one Variants were observed  
403 in the coding region (17 synonymous, 4 nonsynonymous). *ND5*, containing a coding region,  
404 contains six Variants, the most for any position within the U1a haplogroup. All *ND5* Variants  
405 coded for synonymous substitutions, while nonsynonymous substitutions were observed for *ND2*  
406 (m.4659 G>A; A64T), *ATPase6* (m.8860 A>G; T112A), and *CYTB* (m.14766 C>T; T7I and  
407 m.15326 A>G; A190T). 21/44 variants fell within coding genes, while the rest were distributed

408 across HVR I (n=4 Variants), HVR II (n=3 Variants), HVR III (n=5 Variants), 12S-RNR1 (n=2  
409 Variants), 16S-RNR2 (n=4 Variants), the D-loop region (n=1 SNP), and control regions (n=2  
410 Variants). Two Variants were found in regions coding for tRNA[D] and tRNA[L:CUN].

411  
412 The U4b sub-haplogroup is the most common sub-haplogroup among the U haplogroup in our  
413 analysis. In all, 64 Variants were observed for the U4b sub-haplogroup, with most of the variants  
414 (n=20) found in the gene encoding 16S-RNR2 (**Supplementary Figure 4B**). Twenty-one Variants  
415 were found in coding regions (14 synonymous and 7 nonsynonymous substitutions), with the  
416 highest number seen in the gene coding for *COI* (n=6 Variants). Five of six Variants coded for  
417 synonymous substitutions, while m.6366 G>A coded for a nonsynonymous substitution (V155I).  
418 Three Variants were found in the gene encoding *CYTB* and were distributed across all subjects  
419 (n=11) in the U4b sub-haplogroup. All three encoded nonsynonymous substitutions, m14766 C>T  
420 (T7I), m.15326 A>G (T194A), and m.15693 T>C (M316T), and need to be further investigated.  
421 Four tRNA mutations were observed in this sub-haplogroup and one mutation in the D-loop region.

422  
423 A total of 52 variants were observed across all samples in the U7a subgroup (**Supplementary**  
424 **Figure 4B**). Twenty-seven Variants were found in noncoding regions, 12S-RNR1, 16S-RNR2,  
425 and the D-loop region. Twenty-five Variants were found in the coding region (17 synonymous and  
426 8 nonsynonymous substitutions), with 17/25 distributed among the ND genes coding for *ND1-6*.  
427 *ND5* (n=6 Variants) encodes five synonymous mutations, with a nonsynonymous mutation  
428 observed at m.14110 T>C (F592L, in 4/6 subjects).

429  
430 A total of 55 Variants was observed for U2e, with the majority (n=33 Variants) falling in the  
431 noncoding regions (HVR I-III and D-loop) and the 12S-RNR1, 16S-RNR2, and tRNA genes.  
432 Twenty-two Variants fell within the coding region (15 synonymous and 7 nonsynonymous  
433 substitutions), of which 8 fell in the ND gene complex (four *ND2*, four *ND5*) and four in the *CYTB*  
434 gene. While all the Variants in the *ND2* and *ND4* genes are synonymous substitutions, all the  
435 Variants in the *CYTB* gene encoded nonsynonymous mutations (m.14766 C>T; T7I in 3/3 subjects,  
436 m.15326 A>G; T194A in 3/3 subjects; m.14831 G>A; A29T and m.15479 C>T; F245L, both in  
437 1/3 subjects).

438  
439 Five subjects in our analysis (n=100) fell within the T haplogroup. We found four sub-haplogroups  
440 within this haplogroup (T1a, 2 subjects: T2b, T2i, and T2g, with 1 subject each). Our analysis  
441 indicated a total of 39 Variants (**Supplementary Figure 4C**) for T1a, with 21/39 Variants found  
442 in noncoding regions, including 12S-rRNA, 16S-rRNA, tRNAs, and control regions, including the  
443 D-loop. Eighteen Variants were observed in the coding region, with the greatest number occurring  
444 in the *CYTB* gene (n=5 Variants). Three Variants within the *CYTB* gene coded for nonsynonymous  
445 mutations, including m.14776 C>T, m.14905 G>A, and m.15452 C>A, coding for T7I, T194A,  
446 and L236I substitutions, respectively.

447

448 The T2b, T2g, and T2i sub-haplogroups contained 35, 42, and 34 Variants, respectively, in total.  
449 We found that *CYTB* contained the majority of the Variants found in the coding regions in these  
450 sub-haplogroups, except for the T2i group in which the *CYTB* Variants (n=5) constituted the  
451 majority of the Variants found in coding and noncoding regions of the genome. Two Variants,  
452 m.14766 C>T and m.15326 A>G, seen in all three groups code for nonsynonymous substitutions,  
453 and m.15452 C>A was seen in T2g and T2i and codes for a nonsynonymous mutation. Single  
454 mutations were seen for m.15497 G>A and m.14798 T>C and code for nonsynonymous  
455 substitutions and need further investigation.

456  
457 The A haplogroup in our study consists of the sub-haplogroup A2v (n=3 subjects). The subjects in  
458 the A2v sub-haplogroup had a total of 17 Variants (**Supplementary Figure 4D**) distributed across  
459 the mitochondrial genome. Twelve of seventeen Variants were found in the noncoding regions  
460 (HVR I, II) and in the 12S rRNA and 16S rRNA genes. Five Variants were distributed in the  
461 coding region across *ND2* (m.4769 A>G and m.6095 A>G), *ATPase6* (m.8860 A>G), *ND4*  
462 (m.11881 C>T), and *CYTB* (m.15326 A>G). Two nonsynonymous substitutions were observed in  
463 the *ATPase6* and *CYTB* genes that need further investigation.

464  
465 F1g (n=1 subject) is a sub-haplogroup, along with Z1a (n=1 subject). A total of 33 and 32 Variants,  
466 respectively, were identified in these groups. Nine *CYTB* Variants were observed in total for both  
467 groups. Two encoded nonsynonymous substitutions, m.14766 C>T (T7I) and m.15326 A>G  
468 (T194A), while the seven other Variants resulted in synonymous mutations. Variants for *ND4L*  
469 are seen only across Z1a and F1g, with the m.10609 T>C SNP in F1g resulting in a  
470 nonsynonymous shift (M47T), while the Z1a SNP resulted in a synonymous substitution  
471 (**Supplementary Figure 4D**).

472  
473 The M haplogroup (n=52 subjects) consists of 12 sub-haplogroups, the most number for a  
474 haplogroup in our study (**Supplementary Figure 4E**). M30d is the sub-haplogroups with the  
475 highest number of subjects in the M haplogroup (n=11 subjects). Fifty-one Variants were identified  
476 in this sub-haplogroup in total, of which 28 Variants were seen in the noncoding regions (HVR I,  
477 II, III), the D-loop region, and the 12S-RNR1 and 16S-RNR2 genes. The remaining 23 Variants  
478 were part of the coding region within *CYTB* (n=8 Variants) and *ND4* (n=5 Variants) and formed a  
479 majority. Nine of thirteen Variants in *CYTB* and *ND4* code for synonymous substitutions, while  
480 four Variants in *CYTB* resulted in nonsynonymous substitutions (m.14766 C>T; T7I, m.15218  
481 A>G; T158A, m.15326 G>A; T194A, and m.15420 G>A; A229T).

482  
483 M39b (n=10 subjects) is one of the largest sub-haplogroups, and a total of 59 Variants were seen  
484 for this sub-haplogroup. The noncoding regions, 12S, 16S, and control regions, together constitute  
485 33/59 of the Variants. Of the remaining 26 Variants, the 5 Variants in the *CYTB* complex constitute  
486 the greatest number, while the ND gene complex accounts for 12 Variants (2 *ND1*, 1 *ND2*, 2 *ND3*,

487 2 *ND4*, 3 *ND5*, and 2 *ND6*). Of the nine remaining Variants, six are seen in the *COI*, *II*, and *III*  
488 genes (two each), while three Variants are found in the *ATPase6* gene.

489  
490 The M2 sub-haplogroup consists of M2a (n=2 subjects) and M2b (n=1 subject). A total of 110  
491 Variants was observed in total for M2a and M2b (**Supplementary Figure 4E**). In M2a, 23/53  
492 Variants occurred in noncoding regions (HVR I, II, III), the 12S-RNR1 and 16S-RNR2 genes, the  
493 control region (OL), and the D-loop region. Thirty Variants occurred in the coding regions, making  
494 this one of the sub-haplogroups in which Variants in the coding region outnumber the Variants in  
495 the noncoding region. *CYTB* harbors seven Variants, followed by three Variants in *ND4* and three  
496 Variants in *ATPase8*, *ATPase6*, and *COI*. A total of 55 Variants was observed for M2b, in which  
497 31/55 Variants occurred in the noncoding regions. Twenty-four Variants were observed in genes  
498 coding for *COI*, *III*; *ND1,2,3,4,5*; *ATPase6,8*; and *CYTB*. The six Variants in *CYTB* constitute the  
499 greatest number of Variants in the coding region. The M2a/b sub-haplogroup is also conspicuous  
500 by the presence of Variants in the *ATPase8* gene, which is not observed in any sub-haplogroup  
501 besides U4b. The complete distribution of the Variants across all the sub-haplogroups is presented  
502 in **Table 2**.

503  
504 The M3a sub-haplogroup (n=8 subjects) consists of 38 variants, with 12/38 variants in the HVR I,  
505 II, III, D-loop regions (**Supplementary Figure 4E**). 19/38 variants were observed in the protein  
506 coding regions, with the most variants in this region occurring in *CYTB* (n=5). We found 15 coding  
507 for synonymous substitutions and 5 for non-synonymous variants (Supplementary Figure 4E)

508  
509 M52b sub haplogroup (n=9 subjects) contained a total of 90 variants. 29/90 variants were observed  
510 in HVR I, II, III and the D-loop (**Supplementary Figure 4E**). 31 variants were observed for  
511 protein coding genes. *CYTB* (n=9 variants) contains the most variants for this region. 2 variants  
512 were found in t-RNA coding genes. 22 variants coded for synonymous substitutions while 9  
513 variants coded for non-synonymous substitutions.

514  
515 M24a subhaplogroup (n=8 subjects) contains a total of 48 variants, 12/48 are seen in HVR I, II,  
516 III and D-loop (**Supplementary Figure 4E**). 22/48 are found in protein encoding genes with the  
517 most on *CYTB* (n=5 variants). 13 synonymous variants and 7 non-synonymous variants are seen  
518 in this sub-haplogroup. The rest of the variants are seen in 12S, 16S-rRNA. No variants for t-RNA  
519 genes were observed in this sub-haplogroup.

520  
521 M27b (n=1 subject) has a total of 41 variants (**Supplementary Figure 4E**). 16/41 are seen in HVR  
522 I, II, III and the D-loop. 22/41 variants are seen in protein encoding genes with the highest variant  
523 count in *CYTB* (n=6 variants). 14 synonymous and 8 non-synonymous variants are observed for  
524 this sub-haplogroup and 1 variant for t-RNA coding gene.

525

526 M4a (n=1 subject) contains a total of 40 variants. 15/40 variants are seen in the non-coding regions  
527 of HVRI, II, III and D-loop (**Supplementary Figure 4E**). 21 variants are seen in the protein  
528 coding region with *CYTB* gene (n=5 variants) containing the highest variant count. Like M27b,  
529 M4a contains 14 synonymous and 7 non-synonymous variants and 1 variant on the t-RNA coding  
530 gene.

531  
532 A total of 45 variants was seen in M5a sub-haplogroup (n=2 subjects) (**Supplementary Figure**  
533 **4E**). 19/45 was seen in protein coding genes with *CYTB* (n=7 variants) representing the highest  
534 variants in the protein coding region. 13 variants code for synonymous substitutions while 6 code  
535 for non-synonymous variants. 1 variant is observed for a t-RNA coding gene.

536  
537 M35b sub-haplogroup (1 subject) contains a total of 40 variants (**Supplementary Figure 4E**).  
538 15/40 variants are seen in HVR I, II, III and D-loop and 20/40 variants are found in protein  
539 encoding regions with the most variants observed in *CYTB* gene (n=5 variants). 14 code for  
540 synonymous substitution while 7 code for non-synonymous substitutions. 1 variant is observed for  
541 a t-RNA coding gene.

542  
543 M33a sub-haplogroup (n=1 subject) contains 39 variants (**Supplementary Figure 4E**). 15/39  
544 variants are observed in HVR I, II, III and D-loop, 19/39 variants are seen in the protein coding  
545 region, with the highest count seen for *CYTB* (n=5 variants) for this region. 12 are synonymous  
546 and 7 are non-synonymous substitutions. 1 variant for t-RNA coding gene is also observed in this  
547 sub-haplogroup. This haplogroup is unique amongst the 25 sub-haplogroups owing to the presence  
548 of a variant (m.8562 C>T) at *ATPase6/8* gene.

#### 549 550 **Phylogenetic analysis of the Parsi mitochondrial haplotypes with those of Iranians and** 551 **Indians**

552 To further investigate the substructure of the major haplogroups identified in the Parsi cohort, a  
553 comparative analysis of haplotypes from 452 complete mtDNA sequences, including 352  
554 Iranians<sup>25</sup> and 100 Indian mitochondrial genome sequences, was undertaken. The rationale for  
555 selection of these two populations centered around the ancestral migration patterns of the Parsis of  
556 India<sup>30</sup>. This grouping also complements the model of the Parsi origin stemming from the ancient  
557 Iranian plateau<sup>31</sup>.

558  
559 Analysis of the haplogroups identified in the Parsis compared with the Iranians, of whom the  
560 Persians (n=180) and the Qashqais (n=112) were the most frequent representatives, demonstrated  
561 that a) all seven Parsi haplogroups were found within the Iranian haplogroup set and b) a marked  
562 lack of haplogroup diversity was observed in the Parsi dataset (n=7 principal haplogroups)  
563 compared with the Persians and Qashqais (n=14 principal haplogroups, **Figure 7A, B**). The reason  
564 for this lack of haplotype diversity likely lies in the practice of endogamy, which has been strictly  
565 adhered to in the Parsi community for centuries, following their arrival from the Iranian plateau.



566 Contemporary populations of Iranians in the Iranian plateau represent diverse haplogroupings,  
567 possibly due to admixture following political upheavals in the region after the departure of Parsis  
568 from ancient Iran around 745 AD<sup>31</sup>.

569  
570 The presence of the predominantly Eurasian mtDNA haplotypes HV, T, and U in our study cohort  
571 was remarkable, given that Parsis have resided on the Indian subcontinent for over 1200 years.  
572 While the majority of Parsis with M haplogroups can be linked to Persian descent, 2 sub-  
573 haplogroups (M2a, M2b) and 1 subject from M30d (n=4 subjects in total) were found to be related  
574 to relic tribes of Indian origin within the M haplogroups in our analyses.

575  
576 A detailed phylogenetic clustering of the Parsis to establish more precise ethnic relationships was  
577 next undertaken. Our analysis revealed that the Parsis predominantly clustered with populations  
578 from Iran (Persians and people of Persian descent, **Figure 8A, 8E**), and the most common HV  
579 group showed that all Parsis in the HV2a tree (n=14) clustered with Persians and Qashqais  
580 (neighbour-joining tree weight > 0.72/72% (**Figure 8A and Table 3**), while the single Parsi in the  
581 HV12b (n=1) haplotype demonstrated a strong association with other Iranian ethnicities, including  
582 the Khorasani and Mazandarani, in addition to the Qashqai and Persians (**Table 3**).

583  
584 A total of 20 Parsi individuals in the U macro-haplogroup were found to fall into four subclades,  
585 U7a (n=6), U2e (n=3), U4b (n=10), and U1a (n=1), with the highest representation in U4b and  
586 U7a (**Figure 8B**). Phylogenetic analysis demonstrated that the Parsis in the U haplogroup cluster  
587 with the Persians most frequently, while a few cluster with Kurds, Armenians, Mazandarani,  
588 Azeris, and Khorasanis, who all claim descent from Mesopotamia and the older Persian empire  
589 (<https://journals.openedition.org/asiacentrale/480>). Among the U haplogroup, U4b and U7a (the  
590 dominant branch of U7) haplotypes are distributed throughout the Near East and South Asia<sup>24</sup> with  
591 subclades specific to Central Asia in the Volga-Ural region<sup>33</sup>, Mediterranean, and Southeast  
592 Europe, with lower frequencies in populations around the Baltic Sea, such as in Latvians and Tver  
593 Karelians<sup>33</sup>. Haplogroup U2 harbors frequency and diversity peaks in South Asia, whereas its U2d  
594 and U2e subclades are confined to the Near East and Europe<sup>24</sup>.

595  
596 The T haplogroup in the Parsi cohort was found to consist of T1a, T2g, T2i, and T2b, with an even  
597 distribution of samples across the subgroups (n=2, 1, 1, 1, respectively). Similar to the haplogroups  
598 HV and U, the Persians and Qashqais form the largest ethnic denomination associated with the  
599 Parsis with respect to the T haplogroups (>60%, **Figure 8C**). Five Parsi individuals of the  
600 haplogroups A2v (n=3), F1g (n=1), and Z1a (n=1) were observed to be phylogenetically related to  
601 Persian, Kurd, Turkmen, and Iranian ethnicities, further attesting to their origin in the Iranian  
602 plateau (**Figure 8C**). The T haplogroup is also well distributed in Eastern and Northern Europe,  
603 as well as in the Indus Valley and the Arabian Peninsula. Younger T subclades are reported to  
604 have expanded into Europe and Central Asia during the Neolithic transition<sup>34</sup>

605

606 Unlike the HV, U, and T haplogroups, within which Parsi's cluster closely with Persians, Parsis  
607 harboring the M haplogroup appear to demonstrate more diversity in their mitochondrial genomes.  
608 This study showed the following breakdown: 8/12 M sub-haplogroups of the 29 Parsi M  
609 haplotypes (M24a [n= 8], M33a [n=1], M5a [n=2], M4a [n=1]), M3a [n=7], M52b [n=8], M27b  
610 [n=1], and M35b [n=1]) clustered with the Persians, Qashqais, Azeris of Iranian ethnicity, and  
611 others of Persian descent (**Figure 8D, Table 3**). Only two sub-haplogroups in our study (M2a and  
612 M2b [n=21], M30d [n=1], (**Figure 8D**) clustered more closely with relic tribes of Indian origin.  
613 Our phylogenetic analyses further showed that 19 Parsi individuals belonging to the M30d (n=10)  
614 and M39d (n=9) haplogroups did not cluster either with Indian or Iranian ethnic groups (**Figure**  
615 **8D**) but remained clustered within their own subgroups.

616  
617 Outgroup sampling is of primary importance in phylogenetic analyses, affecting ingroup  
618 relationships and, in placing the root, polarizing characters. Accordingly, we used AGENOME-  
619 OUTGROUP-Y2b to root the phylogenetic tree. AGENOME-OUTGROUP-Y2b did not associate  
620 with the Zoroastrian-Parsis, Indians and Iranians attesting to the robustness of the method  
621 employed for phylogenetic analysis (**Figure 8E**, black line)

622  
623 **Assembly of the Zoroastrian Parsi mitochondrial consensus genome (AGENOME-ZPMRG-**  
624 **V1.0) and Parsi haplogroup-specific reference sequences**

625 The Parsis of India are a nonsmoking, long-living community despite the prevalence of many  
626 genetic disease manifestations. This prompted us to generate a Parsi-specific mitochondrial  
627 consensus genome to better understand the nuances of disease and wellness in this unique  
628 community. Considering this goal, we classified the Parsi mitochondrial genome based on the  
629 seven identified major haplogroups, HV, M, U, T, A, F, and Z. The haplogroup-specific Parsi  
630 mitochondrial sequences were aligned, and a consensus call for each nucleotide was made based  
631 on the maximal frequency of a base called at each position in the mtDNA genome sequence  
632 (**Appendix 2**).

633  
634 Using this approach, we derived the Zoroastrian Parsi mitochondrial reference sequences for each  
635 haplogroup: AGENOME-ZPMRG-HV-V1.0 (n=15 sequences), AGENOME-ZPMRG-U-V1.0  
636 (n=20 sequences), AGENOME-ZPMRG-T-V1.0 (n=5 sequences), AGENOME-ZPMRG-M-V1.0  
637 (n=52 sequences), AGENOME-ZPMRG-A2v-V1.0, AGENOME-ZPMRG-F1a-V1.0, and  
638 AGENOME-ZPMRG-Z-V1.0 (**Table 4**). Additionally, using all 100 Parsi mitochondrial genomes  
639 generated in this study (see Materials and Methods), we built the first standard Zoroastrian Parsi  
640 mitochondrial consensus genome (AGENOME-ZPMCG-V1.0). The consensus Parsi mtDNA  
641 sequence was found to have 31 unique Variants (**Table 5**), of which five Variants (A263G, A750G,  
642 A1438G, A4769G, and A15326G) were found to be common to the reference sequences of all  
643 seven haplogroups considered (**Table 5**). While the number of Variants unique to each of the seven  
644 haplogroups ranged from 11 to 33, haplogroup M did not appear to have any unique Variants when  
645 compared with the overall consensus sequence, AGENOME-ZPMRG-V1.0. The utility of this

646 newly generated reference standard could be found in the accurate mitochondrial-based analyses  
647 involving the global Zoroastrian Parsi population as well as for individuals of Western Asian,  
648 Indo-European and Indian origin.

649

### 650 **Disease-specific associations of mtDNA variants predict the prevalence of commonly** 651 **occurring diseases in the non-smoking Parsi cohort**

652 As demonstrated in this paper (**Figure 7B**), the practice of intermarriage has likely restricted the  
653 genetic diversity of the Parsis, as measured by the paucity of haplogroups in our cohort compared  
654 with the Persian and Qashqai populations, possibly contributing to a number of autosomal  
655 recessive and other genetic diseases. In previous studies, Parsis were found to be  
656 disproportionately affected with certain diseases, such as prostate and breast cancers<sup>5,11</sup>,  
657 Parkinsons disease (PD), and Alzheimers disease (AD). However, the Parsis are also considered  
658 to be a long-living community<sup>6</sup> with lower incidences of lung cancer<sup>12</sup>.

659

660 In order to determine whether diseases known to be prevalent in the Parsi community could in  
661 fact be predicted by association using the collective mitochondrial variants discovered in this  
662 study, we first analyzed variants identified in tRNA genes in the mitochondrial genome that have  
663 previously been implicated in rare and degenerative diseases. We found a total of 17 tRNA-  
664 associated variants, with a pathogenic variant (amino acid change: G1644A) implicated  
665 significantly in LS/HCM/MELAS, a genetically inherited mitochondrial disease<sup>47</sup>. We also found  
666 a total of six tRNA mutations associated with non-syndromic hearing loss, hypertension,  
667 breast/prostate cancer risk, and progressive encephalopathies in the analysis of our 100 Parsi  
668 individuals (**Table 6**).

669

670 Further analysis of the nucleotide transitions and transversions that constitute the 420 variants  
671 revealed that the mutational signatures (C>A and G>T) found in tobacco smoke-derived cancers<sup>36</sup>  
672 were found at an extremely low frequency (<6% compared to other mutational signatures) on both  
673 the H and L strands of the mitochondrial genomes of the Parsi population (**Figure 9**), who refrain  
674 from smoking due to their religious and social habits.

675

### 676 **Variant analysis**

677 Furthermore, we found that the 420 variants analysed were associated with 41 diseases. SNP  
678 disease-association analysis revealed that Parkinson's disease is highly associated with our  
679 variants (178 Variants, **Supplementary Figure 5**). Other neurodegenerative diseases, rare  
680 diseases of mitochondrial origin, and cardiovascular and metabolic diseases associated with the  
681 variants in our study were also predicted (**Supplementary Figure 5**).

682

683 While a predisposition to 41 diseases were spread across 25 sub-haplogroups, many diseases were  
684 found to be recurring across haplogroups, totalling 188 diseases (**Figure 10A**). Haplogroup U4b  
685 harbored 15 diseases associations, while the majority of M and T groups had five diseases (Figure

686 6B). Some of the mitochondrial rare diseases, such as mitochondrial encephalomyopathies,  
687 MELAS syndrome and cytochrome c oxidase deficiency were found to be associated with M2a  
688 and U1a, U4b and M2b sub-haplogroups respectively (**Figure 10B**).

689

### 690 **Haplogroup and disease linkage**

691 Since the 420 variants identified fell into 25 sub-haplogroups contributing to 41 diseases and  
692 conditions, Principal component analysis (PCA) showed the grouping of variants and haplogroups  
693 (**Figure 11**). Alzheimers disease, breast cancer, cardiomyopathies, and Parkinsons disease were  
694 represented in all the 25 sub-haplogroups (**Appendix 3**), and longevity was represented in 23 sub-  
695 haplogroups, with the exception of HV12b and U1a groups. Our tRNA pathogenicity analysis  
696 showed that the variability in tRNA was highest in the U, T, and M haplogroups compared with  
697 other haplogroups (**Table 6**).

698

### 699 **Analysis of Variants in tRNA genes and the D-loop region in the mitochondrial genome**

700 While most of the variants in mtDNA genome sequences do not affect mitochondrial function,  
701 unlike synonymous/neutral variants, nonsynonymous/non-neutral variants may have functional  
702 consequences, and their effect on mitochondrial metabolism may be strongly deleterious, mildly  
703 deleterious, or even beneficial. We thus analysed, a SNP dataset obtained from 100 Parsi subjects  
704 for nonsynonymous mutations and identified 63 such Variants located within different  
705 mitochondrial genes (**Figure 12**). Twenty of sixty-three variants were found in genes encoding  
706 *CYTB* (n=13) and *ND2* (n=7), followed by *ND5* and *ND1*. Disease-association analysis showed  
707 that these genes were implicated in the onset of neurodegenerative conditions like AD, PD, cancers  
708 of colorectal and prostate origin, metabolic diseases such as type 2 diabetes, and rare diseases such  
709 as LHON (*CYTB* and *ND2*), (**Figure 13, Figure 14**). Variants implicated in longevity were  
710 observed in our study and distributed across the *ND2* gene (**Figure 10B**). As observed earlier, we  
711 found no association of the nonsynonymous variants in our data set that linked to lung cancer or a  
712 risk of lung cancer.

713

714 To understand the mitochondrial pathways affected by the variants in our study, we annotated the  
715 pathways associated with Variants with DAVID and UNIPROT and found that the major genes  
716 *CYTB* and *ND2* were implicated in pathways that include the mitochondrial respiratory complex  
717 (*COI/COII/COIII/COIV*), OXPHOS, and metabolic pathways implicated in mitochondrial  
718 bioenergetics. Critical disease-related pathways in Parkinsons disease, Alzheimers disease, and  
719 cardiac muscle contraction were also associated with *CYTB*- and *ND2*-specific Variants, which  
720 possibly explains the high incidence of these disease in the Zoroastrian-Parsi population (**Figure**  
721 **15**).

722

723 A total of 87 variants, including 6 unique variants, were observed in the D-loop region across all  
724 25 sub-haplogroups (n=100 subjects, **Table 2**). 74/100 Parsis in our study, were found to have the  
725 polymorphism m.16519 T>C that is associated with chronic kidney disease<sup>43</sup>, an increased risk

726 for Huntingtons disease, migraine headache, and cyclic vomiting syndrome<sup>44</sup> and schizophrenia  
727 and bipolar disorder<sup>48</sup> a. While six subjects of the M52 sub-haplogroup were found to have  
728 m.16525 A>G. The rest of the variants were found at m.16390 G>A (n=4 subjects) and m.16399  
729 A>G, m.16401 C>T, and m.16497 A>G (all with n=1 subject each). Taken together, these results  
730 warrant a deeper investigation into the D-loop variants in the Zoroastrian-Parsi community.

731

### 732 **Identification of unique, unreported variants from the 100 Parsi-Zoroastrian mitogenome** 733 **analysis**

734 We performed a comparative analysis of the 420 variants in the Zoroastrian-Parsi community with  
735 MITOMASTER<sup>45</sup>, a database that contains all known pathogenic mtDNA mutations and common  
736 haplogroup polymorphisms, to identify unique Variants in our population, that are not reported  
737 previously. Our analysis showed the presence of 12 unique Variants distributed across 27 subjects  
738 that were not observed in MITOMASTER and additionally in the VarDIG<sup>®</sup>-R disease association  
739 dataset (**Figure 16**). These unique variants were observed across different gene loci. 12S-rRNA  
740 (2 variants), 16S-rRNA (5 Variants), 1 each at ND1, COII, COIII, ND4 and ND6. The SNP  
741 haplogroup association showed that they fell into 4 major haplogroups and 13 sub haplogroups;  
742 HV2a=1, M24a=4, M2a=1, M30d=3, M35b=1, M39b=2, M3a=1, M4a=1, M52b=4, M5a=1,  
743 T2b=1, U4b=6, U7a=1. Of the 12 variants identified, no disease associations were observed for on  
744 analysis with MITOMASTER and VarDIG<sup>®</sup>-R and needs to be further investigated.

745

### 746 **Discussion**

747

748 The first *de novo* Parsi mitochondrial genome, AGENOME-ZPMS-HV2a-1 (Genbank accession:  
749 MT506314) from a healthy, non-smoking female of haplogroup HV2a when compared with the  
750 revised Cambridge Reference Standard (rCRS) showed 28 unique variants. Upon extending our  
751 mitochondrial genome analyses to an additional 99 Parsi individuals, we found that 94 individuals  
752 separate into four major mitochondrial haplogroups, HV, U, T, and M, while 5 individuals belong  
753 to the rarer haplogroups A, F, and Z. The largest sub-haplogroup was found to be HV2a (n=14).

754

755 Due to the strict endogamy practiced by the Zoroastrian Parsis, their maternally inherited  
756 mitochondrial lineage seems to have remained aligned with those of their ancestors in Old Persia,  
757 prior to 642 AD. On comparison of the major mitochondrial haplogroups in our Parsi cohort with  
758 352 Iranian<sup>25</sup> and 100 Indian mitochondrial genomes, we observed that the Zoroastrian-Parsi  
759 genomes are phylogenetically related to the Persians and Qashqais<sup>25</sup> in HV, T, U, F, A and Z  
760 haplogroups, those associated with peopling of western Europe, Central Asia and the Iranian  
761 plateau.

762

763 The haplogroup HV2, dated at 36–42 kya, most likely arose in Persia between the time of the first  
764 settlement by modern humans and the last glacial melt, and the subclade HV2a has a demonstrated  
765 Persian ancestry. HV12b, a branch of the HV12 clade, is one of the oldest HV subclades and has

766 been found in western Iran, India, and sporadically as far as Central and Southeast Asia. It has  
767 strong associations with the Qashqais, who are Turkic-speaking nomadic pastoralists of southern  
768 Iran and who previously resided in the Iranian region of the South Caucasus<sup>32,35</sup>. The presence of  
769 these predominantly Eurasian mtDNA haplotypes, HV, T, U, F, A and Z in our Parsi cohort attests  
770 to their practice of endogamy, given that Parsis have resided on the Indian subcontinent for over  
771 1300 years.

772  
773 Despite the large grouping of the M haplogroup (the largest haplogroup in the Indian subcontinent  
774 <sup>34</sup>) in our Parsi cohort, phylogenetic analysis showed that 47/51 Parsis belonging to the M  
775 haplogroups in our study, cluster with the Persians, suggesting Persian descent, with a small  
776 minority of Parsis found to be related to relic tribes of India. This observation suggests minimal  
777 gene flow from indigenous Indian females into the Parsi gene pool, as had been previously  
778 proposed<sup>30</sup>.

779  
780 Phylogenetic analysis also revealed that two Parsi M sub-haplogroups, M30d and M39b formed a  
781 unique cluster that needs further resolution.

782  
783 We further present the first complete Zoroastrian Parsi Mitochondrial Consensus Genome  
784 (AGENOME-ZPMC G V1.0), built from the mitochondrial genomes of 100 non-smoking, Parsi  
785 individuals representing seven mitochondrial haplogroups. The need for the generation of such an  
786 ethnic-specific consensus genome, specifically for the Parsis, is self-evident for studies involving  
787 comparative analyses, designed to precisely understand patterns of maternally inherited  
788 mitochondrial DNA and aid in reconstructing the history and prevalent disease associations in this  
789 unique community.

790  
791 We found that *CYTB* gene contained the maximum number of variants ( $n \geq 5$ ) in the coding region  
792 of haplogroup M, besides having maximal representation in F1g, T, and HV12b. Haplogroups U,  
793 A2v, and Z1a showed dominance for the ND complex genes *ND5* and *ND2*, while the *COI* genes  
794 were the most highly represented in HV2a and U4b. Variants in the *CYTB* gene are associated with  
795 Alzheimers disease, diabetes mellitus, cognitive ability, breast cancer, hearing loss, and  
796 asthenozoospermia and associated with changes in metabolic pathways, cardiac contraction and  
797 rare diseases such as Huntington's disease, whereas the *ND2* and *ND5* variants were associated  
798 with prostate, ovarian cancer, rare mitochondrial neuronal diseases, such as LHON,  
799 cardiomyopathy, Alzheimers disease and Parkinsons disease.

800  
801 tRNA disease-association analysis in our study showed that these genes were implicated in the  
802 onset of neurodegenerative conditions, such as Alzheimers disease, Parkinsons disease, cancers of  
803 colorectal and prostate origin, metabolic diseases, such as type 2 diabetes, and rare diseases, such  
804 as LHON (*CYTB* and *ND2*). The D-loop SNP analysis showed the prevalence (74/100 subjects) of  
805 the m.16519 T>C polymorphism, which has been implicated in chronic kidney disease<sup>43</sup>, an

806 increased risk for Huntingtons disease, migraine headache, and cyclic vomiting syndrome<sup>44</sup>. Taken  
807 together, these results warrant a deeper investigation into the tRNA and the D-loop variants in the  
808 Zoroastrian-Parsi community.

809 Consanguineous marriages amongst the Parsis has given rise to longevity<sup>11</sup> and number of  
810 associated diseases, including colon, prostate, breast cancers<sup>9,10</sup>, Parkinson's disease (PD),  
811 Alzheimer's disease (AD), and lower incidences of lung cancer<sup>12</sup>. Interrogation of the 420  
812 variants across seven haplogroups in the Parsi cohort using the VarDIG<sup>®</sup>-R database revealed that  
813 Parkinson's disease, known to be prevalent in the Parsi community<sup>36</sup>, was predicted to have the  
814 highest prevalence with 178 of the 420 variants represented. Not surprisingly, longevity, which  
815 often co-occurs with Parkinsons disease, was also predicted to be highly prevalent in the Parsi  
816 cohort, but with a notable absence in the U1 sub-haplogroup, an interesting observation that  
817 warrants further investigation.

818  
819 Analysis of additional disease associations revealed that Alzheimer's disease (also related to  
820 ageing), breast cancer, and cardiomyopathies<sup>38,39,40</sup>, were predicted to be associated with all the 25  
821 Parsi sub-haplogroups. Additionally, the observed low birth rate among Parsi could be predicted  
822 from the presence of variants associated with asthenozoospermia<sup>37</sup>; a condition associated with  
823 reduced sperm motility.

824  
825 It is noteworthy that previously published studies demonstrating lower rates of lung cancer  
826 amongst the Parsis<sup>41</sup>, appears to hold true, given that no haplogroup in the Parsi cohort  
827 demonstrated a predicted predisposition to lung cancer. The lack of mitochondrial signatures for  
828 lung cancer in all haplogroups examined in this non-smoking Parsi population coupled with the  
829 low frequency of mutational signatures for tobacco smoke-derived cancers was not surprising,  
830 particularly since the Zoroastrian-Parsi venerate fire, and smoking would lie in gross violation of  
831 that sacred tenet. Other diseases predicted to occur at high frequencies in our analyses await further  
832 investigation in the Parsi community.

833  
834 The Parsi haplogroup specific variant-disease association analysis has shed predictive light on the  
835 association of mitochondrial variants linked to longevity, neurodegenerative diseases, cancers of  
836 the colon, breast and prostate and low birth rate, among others; diseases that have been well  
837 documented to occur in the Parsi community. The Parsis thus represent a small but unique, non-  
838 smoking community where genomic disease signatures, both mitochondrial and nuclear, can be  
839 investigated in the backdrop of generations of endogamy thus providing exceptional opportunities  
840 to understand and mitigate disease.

## 841 842 **Conclusion**

843 We have generated the first *de novo* Zoroastrian Parsi Mitochondrial Reference Sequence  
844 (AGENOME- ZPMS-HV2a-1) and the Zoroastrian Parsi Mitochondrial Consensus Genome  
845 (AGENOME-ZPMC V1.0). This newly generated reference standards will contribute in the

846 analysis of mitochondrial genomes of not only the Zoroastrian Parsi population but also other  
847 populations. We have also provided evidence that the Zoroastrian Parsis of India, through centuries  
848 of endogamy, have retained their Persian genetic heritage, distinct traits of longevity and associated  
849 diseases. We have shown the rôle of social habits in genetic signatures exemplified by the lack of  
850 mitochondrial variants associated with lung cancer.

851  
852 In sum, The Parsi haplogroup specific variant-disease association analysis has shed predictive light  
853 on the association of mitochondrial variants linked to longevity, neurodegenerative diseases,  
854 cancers of the colon, breast and prostate and low birth rate, among others; diseases that have been  
855 well documented to occur in the Parsi community. The Parsis thus represent a small but unique,  
856 non-smoking community where genomic disease signatures, both mitochondrial and nuclear, can  
857 be investigated in the backdrop of generations of endogamy thus providing exceptional  
858 opportunities to understand and mitigate disease.

859

## 860 **References**

- 861 1. Mistry, R. K. *Glimpses of Parsi history, Insights Into The Zarathustrian Religion*, p.20.
- 862 2. Nariman, R. F. *The Inner Fire – Faith, Choice, and Modern Day Living in*  
863 *Zoroastrianism*, p. 20-21
- 864 3. The Vendidad: The Zoroastrian Book Of The Law Paperback – September 10, 2010. I, 1-2  
865 & II, 5. Charles. F. Horne. ISBN-10: 1162910089; ISBN-13: 978-1162910086. Kessinger  
866 Publishing, LLC (September 10, 2010)
- 867 4. Bennet, J. G. The Hyperborean Origin of the Indo-European Culture, *Journal Systematics*.  
868 *J Syst.* **1**, (1963).
- 869 5. Jussawalla, D. J., Yeole, B. B. & Natekar, M. V. Histological and epidemiological  
870 features of breast cancer in different religious groups in greater bombay. *J. Surg. Oncol.*  
871 (1981) doi:10.1002/jso.2930180309.
- 872 6. Jussawalla, D. J. The persistence of differences in cancer incidence at various anatomical  
873 sites 1300 years after immigration. *Recent Results Cancer Res.* (1975) doi:10.1007/978-3-  
874 642-80880-7\_22.
- 875 7. Anthony, DW, (2007), *The Horse, The Wheel, And Language. How Bronze-Age Riders*  
876 *from the Eurasian Steppes Shaped the Modern World*, Princeton University Press. p. 9.
- 877 8. Alizadeh, A. The Rise of the Highland Elamite State in Southwestern Iran. *Current. Curr*  
878 *Anthropol.* **51**, 353–383 (2010).
- 879 9. Shroff Z, C. M. The potential impact of intermarriage on the population decline of the  
880 Parsis of Mumbai, India. *Demogr Res.* **25**, 545–564 (2011).
- 881 10. Karkal, M. Marriage among Parsis. *Demogr. India* **4**, 128 (1975).
- 882 11. Barnabas-Sohi, N. *et al.* Breast carcinoma in a high-risk population: Structural alterations  
883 in neu, int-2, and p-53 genes. *Breast Dis.* (1993).
- 884 12. Jussawalla, D. J. & Jain, D. K. Lung cancer in Greater Bombay: Correlations with religion



- 885 and smoking habits. *Br. J. Cancer* (1979) doi:10.1038/bjc.1979.199.
- 886 13. Helgason, A., Sigurðardóttir, S., Gulcher, J. R., Ward, R. & Stefánsson, K. mtDNA and  
887 the origin of the Icelanders: Deciphering signals of recent population history. *Am. J. Hum.*  
888 *Genet.* (2000) doi:10.1086/302816.
- 889 14. Wallace, D. C. Mitochondrial DNA Variation in Human Radiation and Disease. *Cell*  
890 (2015) doi:10.1016/j.cell.2015.08.067.
- 891 15. Wallace, D. C., Brown, M. D. & Lott, M. T. Mitochondrial DNA variation in human  
892 evolution and disease. *Gene* (1999) doi:10.1016/S0378-1119(99)00295-4.
- 893 16. Roger, A. J., Muñoz-Gómez, S. A. & Kamikawa, R. The Origin and Diversification of  
894 Mitochondria. *Current Biology* (2017) doi:10.1016/j.cub.2017.09.015.
- 895 17. Garcia, I., Jones, E., Ramos, M., Innis-Whitehouse, W. & Gilkerson, R. The little big  
896 genome: The organization of mitochondrial DNA. *Front. Biosci. - Landmark* (2017)  
897 doi:10.2741/4511.
- 898 18. Stewart, J. B. & Chinnery, P. F. The dynamics of mitochondrial DNA heteroplasmy:  
899 Implications for human health and disease. *Nature Reviews Genetics* (2015)  
900 doi:10.1038/nrg3966.
- 901 19. Bussard, K. M. & Siracusa, L. D. Understanding mitochondrial polymorphisms in  
902 cancer. *Cancer Research* (2017) doi:10.1158/0008-5472.CAN-17-1939.
- 903 20. Alston, C. L., Rocha, M. C., Lax, N. Z., Turnbull, D. M. & Taylor, R. W. The genetics  
904 and pathology of mitochondrial disease. *J. Pathol.* **241**, 236–250 (2017).
- 905 21. Andrews, R. M. *et al.* Reanalysis and revision of the cambridge reference sequence for  
906 human mitochondrial DNA [5]. *Nature Genetics* (1999) doi:10.1038/13779.
- 907 22. Chandrasekar, A. *et al.* Updating phylogeny of mitochondrial DNA macrohaplogroup m  
908 in India: dispersal of modern human in South Asian corridor. *PLoS One* **4**, e7447–e7447  
909 (2009).
- 910 23. Rajkumar, R., Banerjee, J., Gunturi, H. B., Trivedi, R. & Kashyap, V. K. Phylogeny and  
911 antiquity of M macrohaplogroup inferred from complete mt DNA sequence of Indian  
912 specific lineages. *BMC Evol. Biol.* **5**, 26 (2005).
- 913 24. Sahakyan, H. *et al.* Origin and spread of human mitochondrial DNA haplogroup U7. *Sci.*  
914 *Rep.* **7**, 46044 (2017).
- 915 25. Derenko, M. *et al.* Complete mitochondrial DNA diversity in Iranians. *PLoS One* (2013)  
916 doi:10.1371/journal.pone.0080673.
- 917 26. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high  
918 throughput. *Nucleic Acids Res.* **32**, 1792–1797 (2004).
- 919 27. Kumar, S., Stecher, G., Li, M., Knyaz, C. & Tamura, K. MEGA X: Molecular  
920 Evolutionary Genetics Analysis across Computing Platforms. *Mol. Biol. Evol.* **35**, 1547–  
921 1549 (2018).
- 922 28. Tamura, K., Nei, M. & Kumar, S. Prospects for inferring very large phylogenies by using  
923 the neighbor-joining method. *Proc. Natl. Acad. Sci. U. S. A.* **101**, 11030–11035 (2004).
- 924 29. Houshmand, M. *et al.* Is 8860 variation a rare polymorphism or associated as a secondary

- 925 effect in HCM disease? *Arch. Med. Sci.* (2011) doi:10.5114/aoms.2011.22074.
- 926 30. Chaubey, G. *et al.* ‘Like sugar in milk’: Reconstructing the genetic history of the Parsi  
927 population. *Genome Biol.* (2017) doi:10.1186/s13059-017-1244-9.
- 928 31. López, S. *et al.* The Genetic Legacy of Zoroastrianism in Iran and India: Insights into  
929 Population Structure, Gene Flow, and Selection. *Am. J. Hum. Genet.* (2017)  
930 doi:10.1016/j.ajhg.2017.07.013.
- 931 32. Quintana-Murci, L. *et al.* Where west meets east: the complex mtDNA landscape of the  
932 southwest and Central Asian corridor. *Am. J. Hum. Genet.* **74**, 827–845 (2004).
- 933 33. Shamoon-Pour, M., Li, M. & Merriwether, D. A. Rare human mitochondrial HV lineages  
934 spread from the Near East and Caucasus during post-LGM and Neolithic expansions. *Sci.*  
935 *Rep.* **9**, 14751 (2019).
- 936 34. Farjadian, S. *et al.* Discordant Patterns of mtDNA and Ethno-Linguistic Variation in 14  
937 Iranian Ethnic Groups. *Hum. Hered.* **72**, 73–84 (2011).
- 938 35. Thangaraj, K. *et al.* In situ origin of deep rooting lineages of mitochondrial  
939 Macrohaplogroup ‘M’ in India. *BMC Genomics* **7**, 151 (2006).
- 940 36. Alexandrov LB, Ju YS, Haase K, et al. Mutational signatures associated with tobacco  
941 smoking in human cancer. *Science*. 2016;354(6312):618 - 622.
- 942 37. E. Ruiz-Pesini, A.C. Lapeña, C. Díez, E. Alvarez, J.A. Enríquez, M.J. López-Pérez  
943 Seminal quality correlates with mitochondrial functionality. *Clin. Chim.*  
944 *Acta.*, 300 (2000), p. 97 105.
- 945 38. Fang, H., Shen, L., Chen, T. et al. Cancer type-specific modulation of mitochondrial  
946 haplogroups in breast, colorectal and thyroid cancer. *BMC Cancer* **10**, 421 (2010).
- 947 39. Van der Walt JM, Dementieva YA, Martin ER, Scott WK, Nicodemus KK, Kroner CC,  
948 Welsh-Bohmer KA, Saunders AM, Roses AD, Small GW, Schmechel DE, Murali  
949 Doraiswamy P, Gilbert JR, Haines JL, Vance JM, Pericak-Vance MA. Analysis of  
950 European mitochondrial haplogroups with Alzheimer disease risk. *Neurosci Lett.* 2004  
951 Jul 15; 365(1):28-32.
- 952 40. van Oven M, Kayser M *Hum Mutat.* Updated comprehensive phylogenetic tree of global  
953 human mitochondrial DNA variation. 2009 Feb; 30(2): E386-94.
- 954 41. Balkrishna Bhika Yeole, AP Kurkure, SH Advani, Sunny Lizzy; An Assessment of  
955 Cancer Incidence Patterns in Parsi and Non Parsi Populations, Greater Mumbai. *Asian*  
956 *Pacific Journal of Cancer Prevention*, Vol 2, 2001; 293-298
- 957 42. Niroula A, Vihinen M. PON-mt-tRNA: a multifactorial probability-based method for  
958 classification of mitochondrial tRNA variations. *Nucleic Acids Res.* 2016;44(5):2020-  
959 2027. doi:10.1093/nar/gkw046
- 960 43. Chen JB, Yang YH, Lee WC, et al. Sequence-based polymorphisms in the mitochondrial  
961 D-loop and potential SNP predictors for chronic dialysis. *PLoS One.* 2012;7(7):e41125.  
962 doi:10.1371/journal.pone.0041125
- 963 44. Zaki EA, Freilinger T, Klopstock T, et al. Two common mitochondrial DNA  
964 polymorphisms are highly associated with migraine headache and cyclic vomiting

- 965 syndrome. *Cephalalgia*. 2009;29(7):719-728. doi:10.1111/j.1468-2982.2008.01793.x
- 966 45. Brandon MC, Ruiz-Pesini E, Mishmar D, et al. MITOMASTER: a bioinformatics tool for
- 967 the analysis of mitochondrial DNA sequences. *Hum Mutat*. 2009;30(1):1-6.
- 968 doi:10.1002/humu.20801.
- 969 46. Narasimhan VM, Patterson N, Moorjani P, et al. The formation of human populations in
- 970 South and Central Asia. *Science*. 2019;365(6457):eaat7487. doi:10.1126/science.aat7487.
- 971 47. Menotti F, Brega A, Diegoli M, Grasso M, Modena MG, Arbustini E. A novel mtDNA
- 972 point mutation in tRNA(Val) is associated with hypertrophic cardiomyopathy and
- 973 MELAS. *Ital Heart J*. 2004;5(6):460-465.
- 974 48. Schulmann A, Ryu E, Goncalves V, et al. Novel Complex Interactions between
- 975 Mitochondrial and Nuclear DNA in Schizophrenia and Bipolar Disorder. *Mol*
- 976 *Neuropsychiatry*. 2019;5(1):13 - 27. doi:10.1159/000495658

977

### 978 **Contributions**

979 VMP conceptualized and guided the experiments; NP, CG analysed the sequences, bioinformatics

980 analysis, interpreted the results; RM, SR, NS co-ordinated wet-lab work flows and data analysis; NP,

981 BM, KK analysed data, plotted graphs and figures; VMP, AKG, PB, KK, RJ drafted the manuscript

982 with inputs from RM, SR, NS, CG. All authors reviewed the manuscript

983

984 All authors researched data for the article, made substantial contribution to discussion of content,

985 and wrote, reviewed, and edited the manuscript before submission.

986

### 987 **Sources of funding**

988 The project was funded by a grant (GG-0005), Cancer risk in smoking subjects assessed by next

989 generation sequencing profile of circulating free DNA and RNA) by the Foundation for a Smoke-

990 Free World, New York, USA.

991

### 992 **Competing interests**

993 The authors declare no competing interests

994

### 995 **Acknowledgements:**

996 We thank the Foundation for Smoke Free World who is advancing global progress in smoking

997 cessation and harm reduction for funding this project and Dr.Derek Yach for his support to this

998 important project.

999 We would like to express our thanks to Prof. Dr. Partha. P. Mazumdar, National Institute of Bio

1000 Medical Genetics (NIBMG), Kolkata and Prof. Dr. Kumarasamy Thangaraj, Center for Cellular

1001 and Molecular Biology (CCMB), Hyderabad; for their leadership, expertise and coordination of

1002 sequencing our samples.

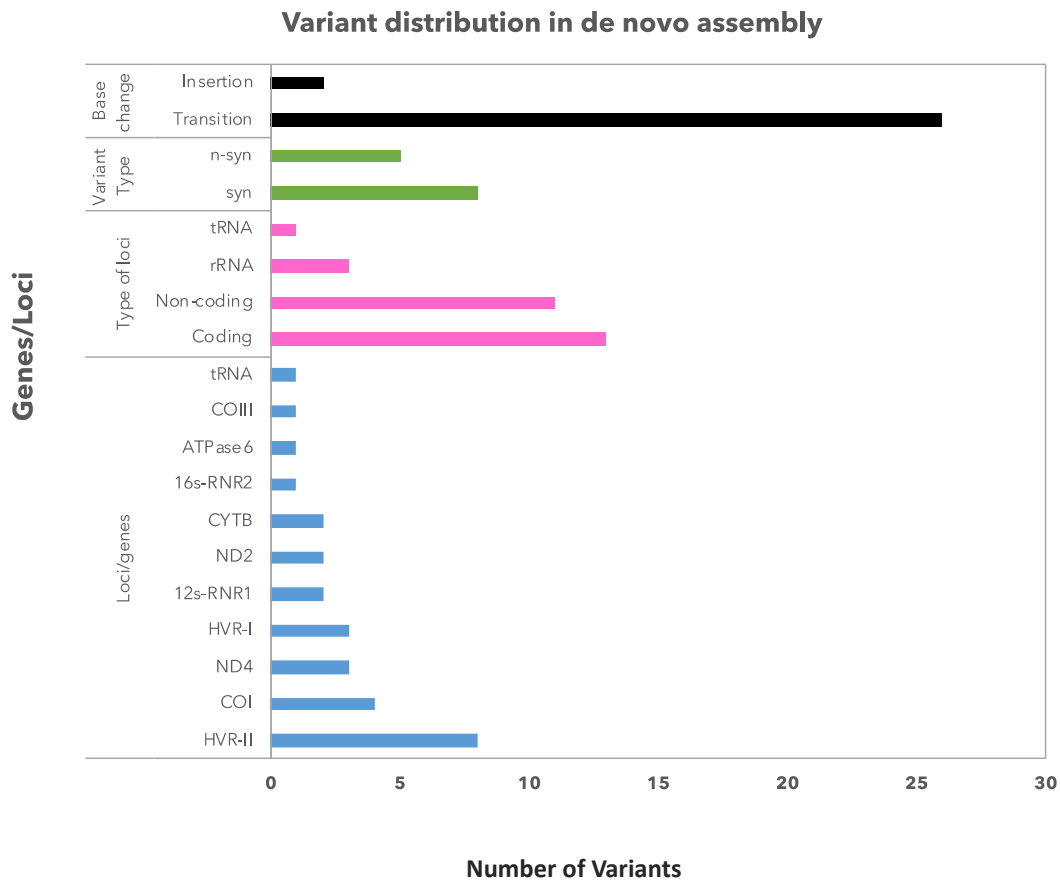
1003 We would like to express our deep gratitude to the Zoroastrian Parsi community of India for their

1004 patience and cooperation at multiple instances.

- 1005 Parzor Foundation for its partnership in outreach among the Zoroastrian-Parsi community in India  
1006 and across the globe.  
1007 Dr.Sami Gazder, Kouser Sonnekhan and the The Avestagenome Project™ project team of yore.

1008 **Main Figures**

1009  
1010 **Figure 1 : Identification of 28 variants in the de novo Parsi mitochondrial**  
1011 **genome, AGENOME-ZPMS-HV2a-1**  
1012



1013  
1014  
1015 **Figure 1: Distribution and classification of Variants in the AGENOME-ZPMS-HV2a-1.**  
1016 Representative histogram showing the base change, variant type, type of loci and distribution of  
1017 variants across genes in the *de novo* mitochondrial genome AGENOME-ZPMS-HV2a-1

1018  
1019  
1020  
1021  
1022

1023 **Figure 2 : Validation of variants in the AGENOME-ZPMS-HV2a-1 by Sanger**  
 1024 **sequencing**  
 1025

1. rCRS	CCTGACTGGCATTGTATTAGCAAACCTCATCACTAGACATCGTACTACACGACACGCTACTA	7018
2. AGENOME-ZPMS-HV2a-1	CCTGACTGGCATTGTATTAGCAAACCTCATCACTAGACATCGTACTACACGACACGCTACTA	7018
3. SANGER-SEQUENCED	CCTGACTGGCATTGTATTAGCAAACCTCATCACTAGACATCGTACTACACGACACGCTACTA	434
	*****	
1. rCRS	CGTTGTAGCCCACTTCCACTATGTCTTATCAATAGGAGCTGTATTTGCCATCATAGGAGG	7078
2. AGENOME-ZPMS-HV2a-1	CGTTGTAGCTCACTTCCACTATGTCTTATCAATAGGAGCTGTATTTGCCATCATAGGAGG	7078
3. SANGER-SEQUENCED	CGTTGTAGCTCACTTCCACTATGTCTTATCAATAGGAGCTGTATTTGCCATCATAGGAGG	494
	*****	
1. rCRS	CAACCGCTATGTATTTCCGTACATTACTGCCAGCCACCATGAATATTGTACGGTACCATAA	16138
2. AGENOME-ZPMS-HV2a-1	CAACCGCTATGTATTTCCGTACATTACTGCCAGCCACCATGAATATTGTACGGTACCATAA	16138
3. SANGER-SEQUENCED	CAACCGCTATGTATTTCCGTACATTACTGCCAGCCACCATGAATATTGTACGGTACCATAA	126
	*****	
1. rCRS	ATACTTGACCACCTATGTACATAAAAAACCAATCCACATCAAAAACCCCTCCCCATGCT	16198
2. AGENOME-ZPMS-HV2a-1	ATACTTGACCACCTATGTACATAAAAAACCAATCCACATCAAAAACCCCTCCCCATGCT	16198
3. SANGER-SEQUENCED	ATACTTGACCACCTATGTACATAAAAAACCAATCCACATCAAAAACCCCTCCCCATGCT	186
	*****	
1. rCRS	TACAAGCAAGTACAGCAATCLACCCCTCAACTATCACACATCAACTGCAACTCCAAAGCCA	16258
2. AGENOME-ZPMS-HV2a-1	TACAAGCAAGTACAGCAATCLACCCCTCAACTATCACACATCAACTGCAACTCCAAAGCCA	16258
3. SANGER-SEQUENCED	TACAAGCAAGTACAGCAATCLACCCCTCAACTATCACACATCAACTGCAACTCCAAAGCCA	246
	*****	
1. rCRS	CCCCTCACCCTAGGATACCAACAAACCTACCCACCCTTAACAGTACATGCTACATAAA	16318
2. AGENOME-ZPMS-HV2a-1	CCCCTCACCCTAGGATACCAACAAACCTACCCACCCTTAACAGTACATGCTACATAAA	16318
3. SANGER-SEQUENCED	CCCCTCACCCTAGGATACCAACAAACCTACCCACCCTTAACAGTACATGCTACATAAA	306
	*****	
	GCCATTTACCGTACATAGCACATTACAGTCAAAATCCCTTCTCGTCCCCATGGATGACCCC	16378
	GCCATTTACCGTACATAGCACATTACAGTCAAAATCCCTTCTCGTCCCCATGGATGACCCC	16378
	GCCATTTACCGTACATAGCACATTACAGTCAAAATCCCTTCTCGTCCCCATGGATGACCCC	366
	*****	

1026

1. SANGER-SEQUENCED	CAGGACATCCCGATGGTGCAGCCGCTATTAAGGTTTCGTTTGTTC AACGATTAAGTCCT	181
2. AGENOME-ZPMS-HV2a-1	CAGGACATCCCGATGGTGCAGCCGCTATTAAGGTTTCGTTTGTTC AACGATTAAGTCCT	3060
3. rCRS	CAGGACATCCCGATGGTGCAGCCGCTATTAAGGTTTCGTTTGTTC AACGATTAAGTCCT *****	3058
1. SANGER-SEQUENCED	ACGTGATCTGAGTT CAGACCCGGAGTAATCCAGGTCGGTTCTATCTAC-TTCAAATTCCT	240
2. AGENOME-ZPMS-HV2a-1	ACGTGATCTGAGTT CAGACCCGGAGTAATCCAGGTCGGTTCTATCTAC-TTCAAATTCCT	3119
3. rCRS	ACGTGATCTGAGTT CAGACCCGGAGTAATCCAGGTCGGTTCTATCTACNTTCAAATTCCT *****	3118
1. SANGER-SEQUENCED	CCCTGTACGAAAGGACAAGAGAAATAAGGCCACTTCACAAGCGCCTTCCCCGTA AAT	300
2. AGENOME-ZPMS-HV2a-1	CCCTGTACGAAAGGACAAGAGAAATAAGGCCACTTCACAAGCGCCTTCCCCGTA AAT	3179
3. rCRS	CCCTGTACGAAAGGACAAGAGAAATAAGGCCACTTCACAAGCGCCTTCCCCGTA AAT *****	3178
1. SANGER-SEQUENCED	GATATCATCTCAACTTAGTATTATACCCACACCCACCCAAGAACAGGGTTTGTTAAGATG	360
2. AGENOME-ZPMS-HV2a-1	GATATCATCTCAACTTAGTATTATACCCACACCCACCCAAGAACAGGGTTTGTTAAGATG	3239
3. rCRS	GATATCATCTCAACTTAGTATTATACCCACACCCACCCAAGAACAGGGTTTGTTAAGATG *****	3238
1. SANGER-SEQUENCED	-----CCCCA	5
2. AGENOME-ZPMS-HV2a-1	ACAATTGAATGTCTGCACAGCCGCTTTCACACAGACATCATAACAAAAAATTTCCACCA	300
3. rCRS	ACAATTGAATGTCTGCACAGCCACTTTCACACAGACATCATAACAAAAAATTTCCACCA * **	300
1. SANGER-SEQUENCED	AACCCCCCTC CCCCCTTCTGGCCACAGCACTTAAACACATCTCTGCCAAACCCCAA	65
2. AGENOME-ZPMS-HV2a-1	AACCCCCCTC CCCCCTTCTGGCCACAGCACTTAAACACATCTCTGCCAAACCCCAA	360
3. rCRS	AACCCCC--CTCCCCGCTTCTGGCCACAGCACTTAAACACATCTCTGCCAAACCCCAA *****	358
1. SANGER-SEQUENCED	AAACAAGAACCCTAACACCAGCCTAACAGATTCAAATTTATCTTTGGCGGTAGCACT	125
2. AGENOME-ZPMS-HV2a-1	AAACAAGAACCCTAACACCAGCCTAACAGATT-----	395
3. rCRS	AAACAAGAACCCTAACACCAGCCTAACAGATT----- *****	393

1027

1028

1029 **Figure 2 : Confirmation of variants identified with next-generation sequencing (NGS) data**

1030 **and confirmation by Sanger sequencing.** Sequences obtained from desired regions were

1031 analyzed for presence of variants/Variants. Low quality bases were trimmed from both ends of the

1032 sequences and used for alignment with the reference Mitochondrial Genome (rCRS). A total of 13

1033 variants/Variants from D-loop and internal region of mitochondrial genome were verified.

1034

1035

1036

1037

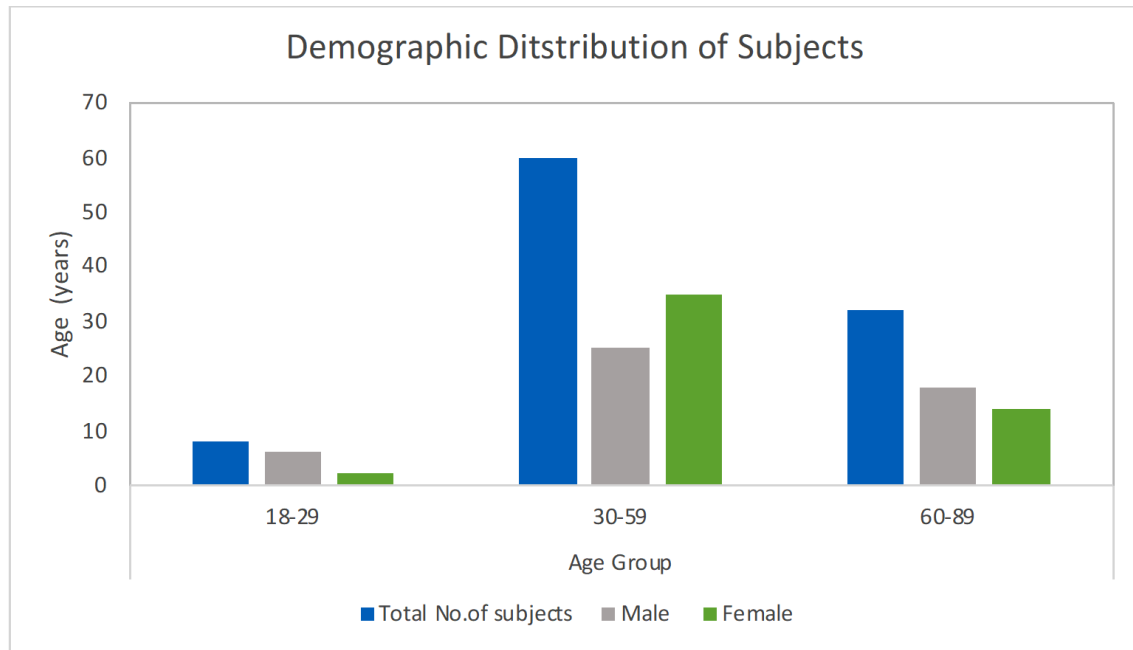
1038

1039

1040 **Figure 3 : Representation of Males and Females in the 100 Zoroastrian-Parsi**  
1041 **whole mitogenome study**

1042

1043



1044

1045 **Figure 3: Distribution of 100 Parsi subjects.** Distribution of the subjects classified based on  
1046 gender and age. The bars on the histogram depict further segmentation of the total number of  
1047 subjects, Male and Female numbers according to their age range.

1048

1049

1050

1051

1052

1053

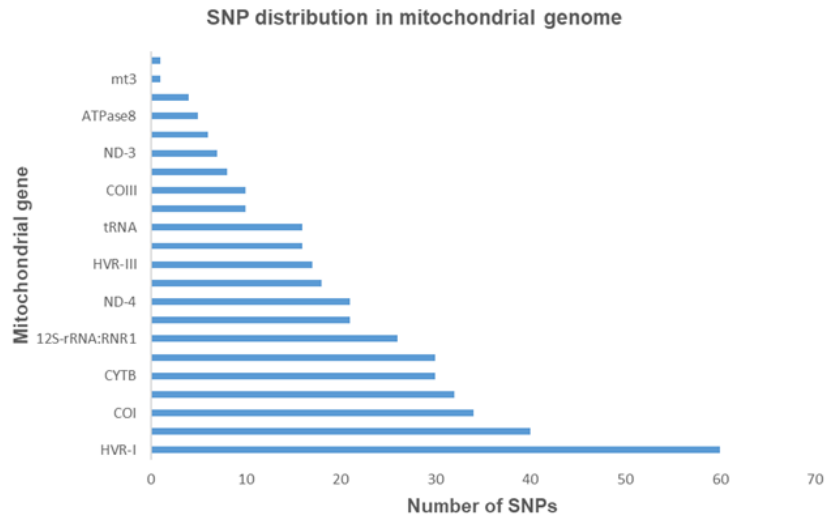
1054

1055

1056

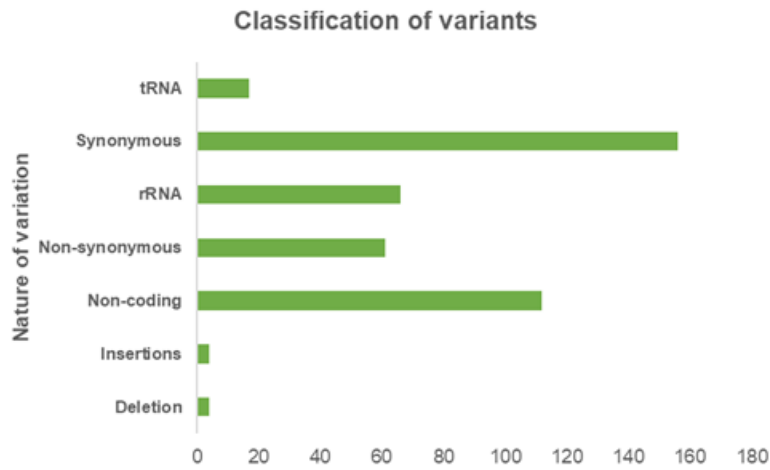


1057 **Figure 4 : Distribution of 420 variants across gene loci in the 100 Zoroastrian-**  
1058 **Parsi whole mitogenomes**  
1059



1060

1061



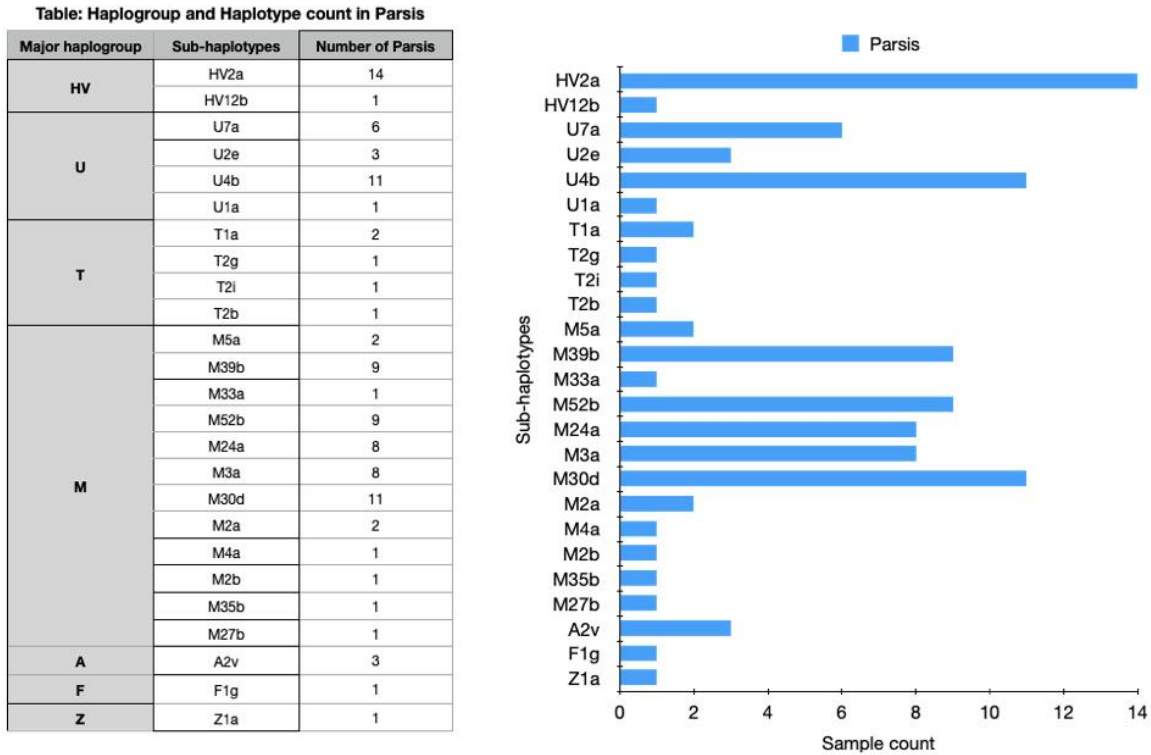
1062

1063

1064 **Figure 4 : Annotation and distribution of 420 variants across 100 Parsi complete**  
1065 **mitogenomes**

1066

1067 **Figure 5 : Identification of 25 sub-haplogroups in the 100 Zoroastrian-Parsi**  
 1068 **study group**  
 1069

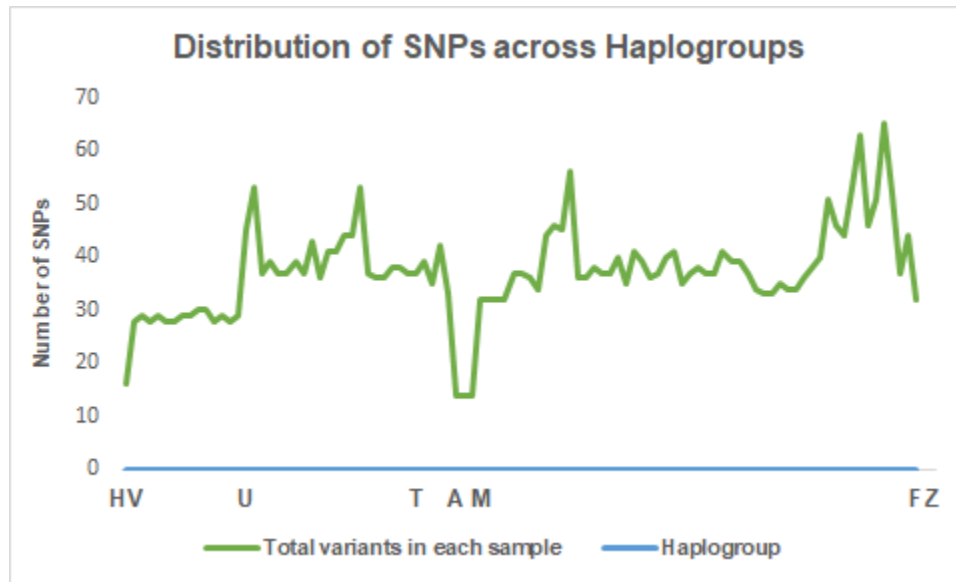


1070

1071 **Figure 5 : Distribution of Parsis across major haplogroups and sub-haplogroups.** The table  
 1072 and the histogram shows the distribution of 100 Parsi subjects across 7 major haplogroups and 25  
 1073 sub-haplogroups  
 1074

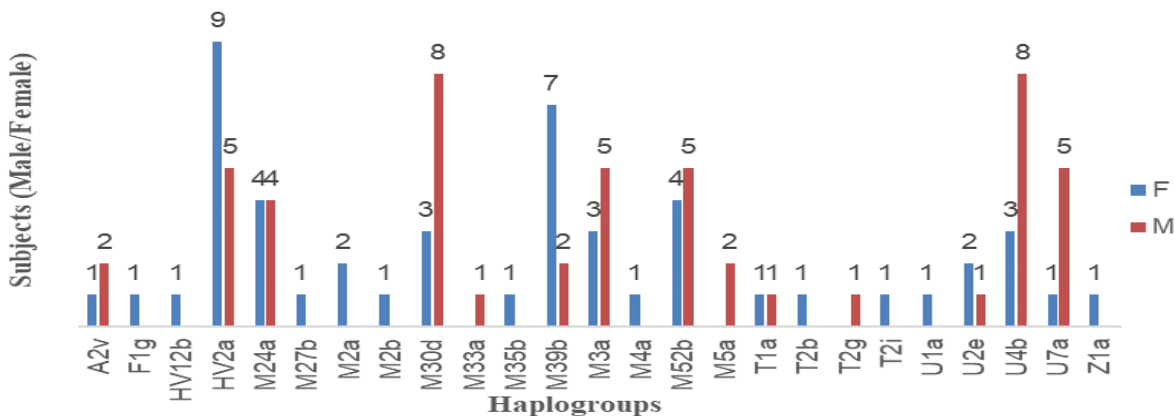
1075 **Figure 6 : Distribution of variants across haplogroups and demographic**  
 1076 **classification of the 100 Parsi study group**  
 1077

1078 **A**



1079

1080 **B**



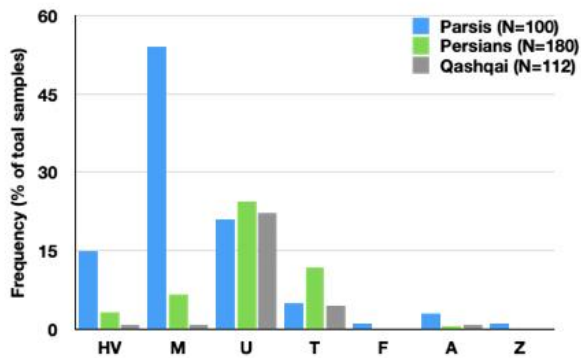
1081

1082 **Figure 6: Distribution across the 100 Zoroastrian-Parsi subjects.** (A) Representative graph  
 1083 depicting the distribution of SNP's count across the 7 major haplogroups (B) Graph depicts the  
 1084 distribution of the subjects classified based on gender across 25 sub-haplogroups  
 1085

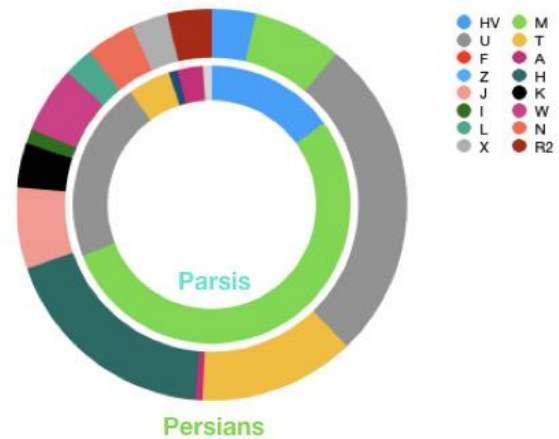
1086

1087 **Figure 7 : Lack of haplogroup diversity in the Parsi cohort suggesting**  
1088 **endogamy**  
1089

1090 A



B



1091

1092 **Figure 7: Comparative analysis of Major haplogroup distribution in the Parsis and**  
1093 **populations of Iranian ethnicities (Persians, Qashqais)** (A) Histogram depicting the analysis of  
1094 The 7 Major haplogroups across Parsis (n=100), Persians (n=180) and Qashqais (n=112) (B)  
1095 Representative figure showing the diversity of major haplogroups in the Parsis and the Persians

1096

1097

1098

1099

1100

1101

1102

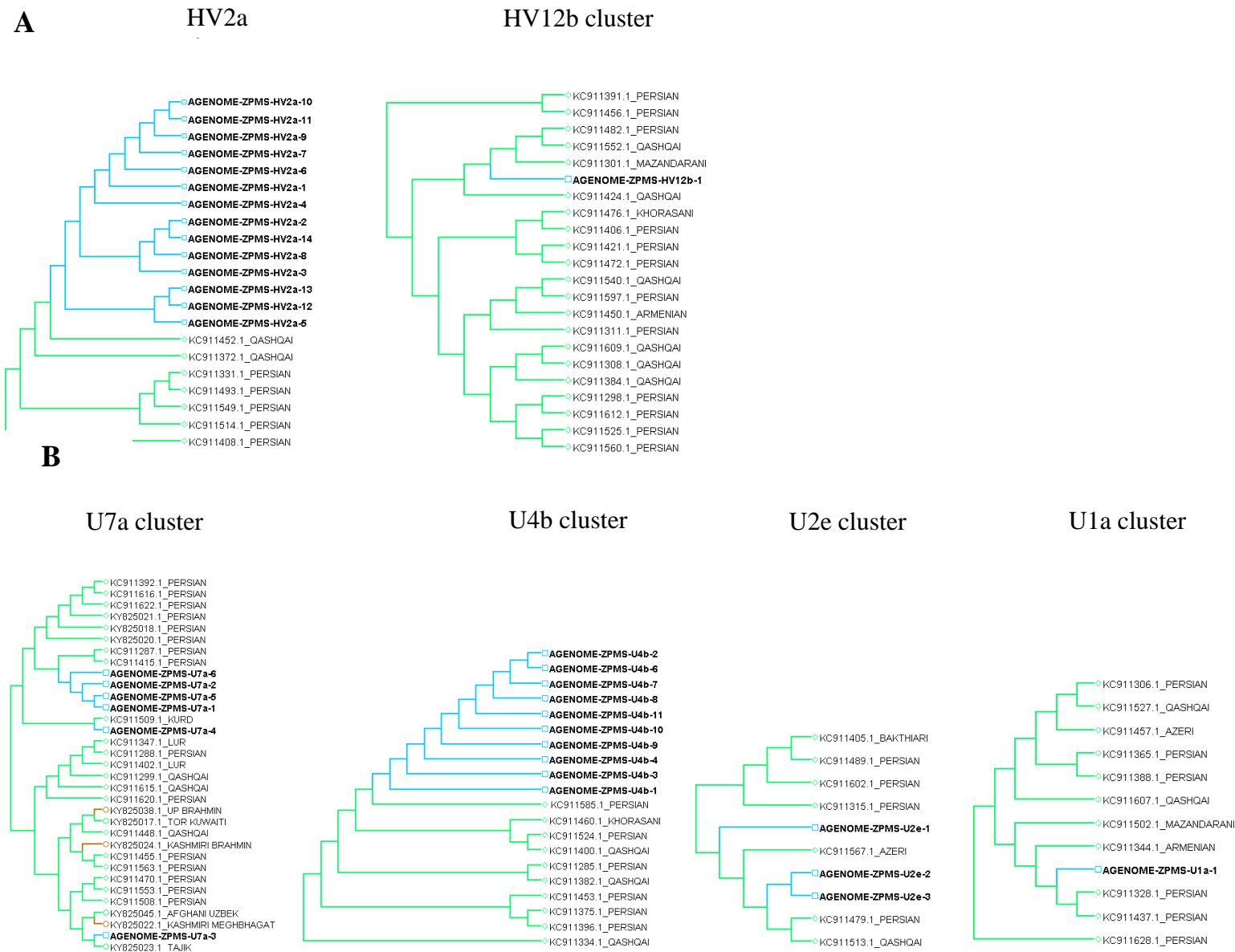
1103

1104

1105

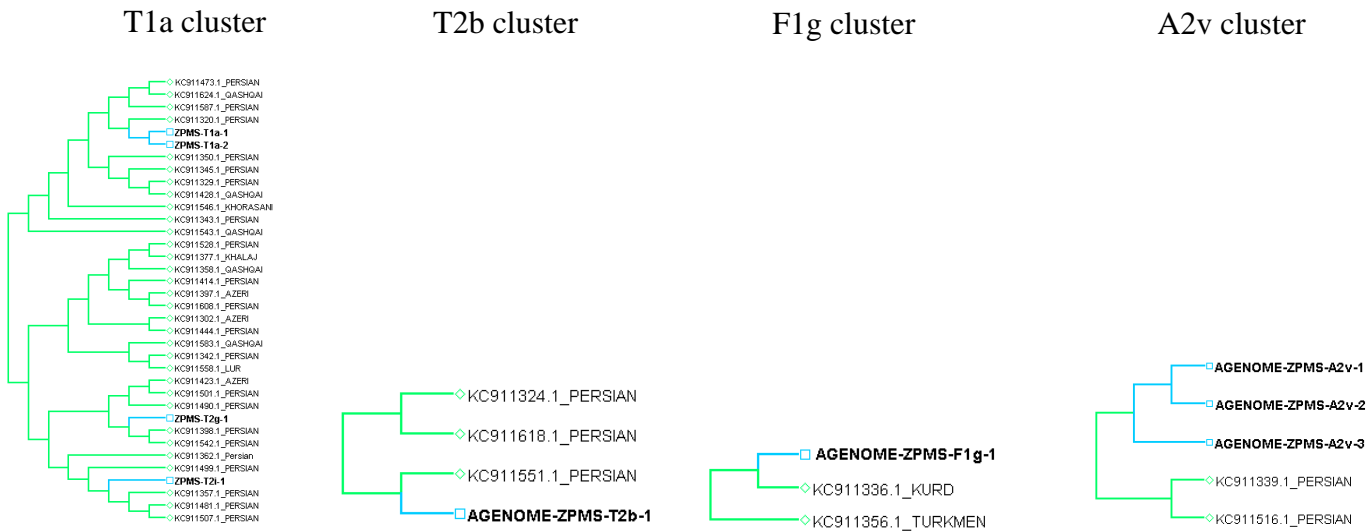
1106

1107



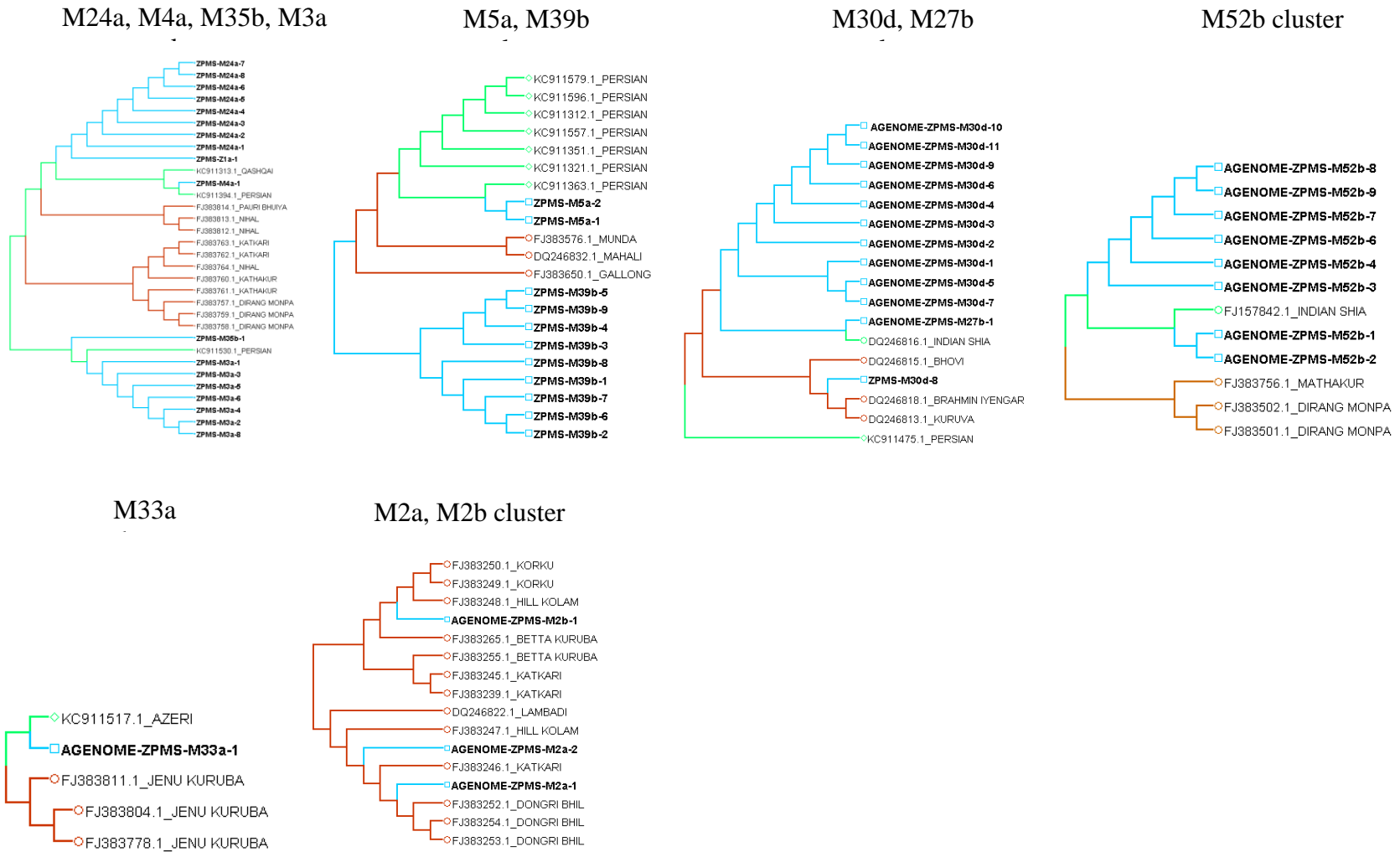
**Figure 8: Phylogenetic analysis depicting individual sub-haplogroup clusters of 97 Parsis, 352 Iranian and 100 relic tribes of Indian origin (A) Representative cladograms of the HV sub-haplogroup (B) Representative cladograms of the U sub-haplogroup**

C



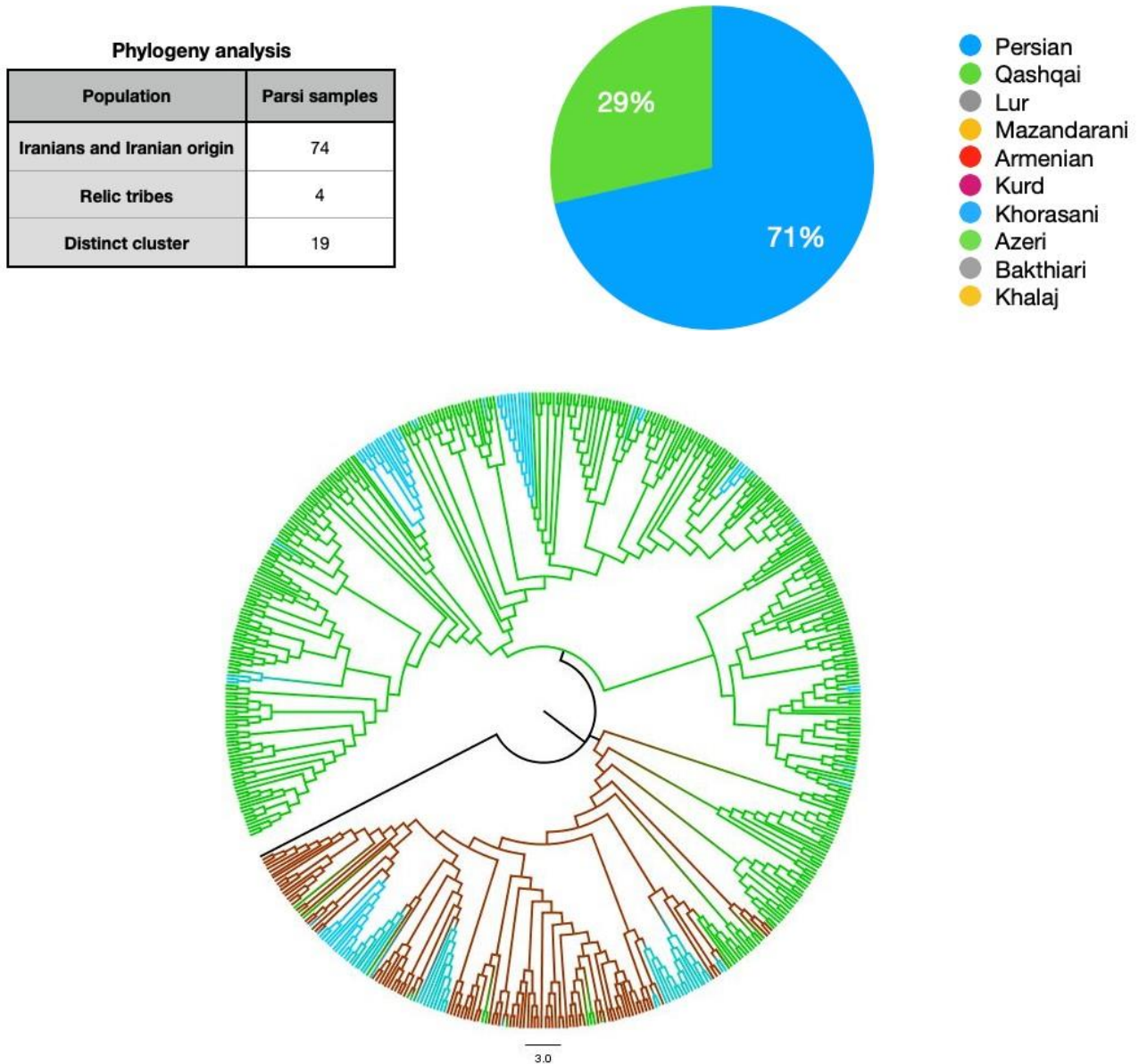
**Figure 8: Phylogenetic analysis depicting individual sub-haplogroup clusters of 97 Parsis, 352 Iranian and 100 relic tribes of Indian origin (C) Representative cladograms of the T, F and A sub-haplogroup**

D



**Figure 8: Phylogenetic analysis depicting individual sub-haplogroup clusters of 97 Parsis, 352 Iranian and 100 relic tribes of Indian origin (A-D) Representative cladograms of the each sub-haplogroup**

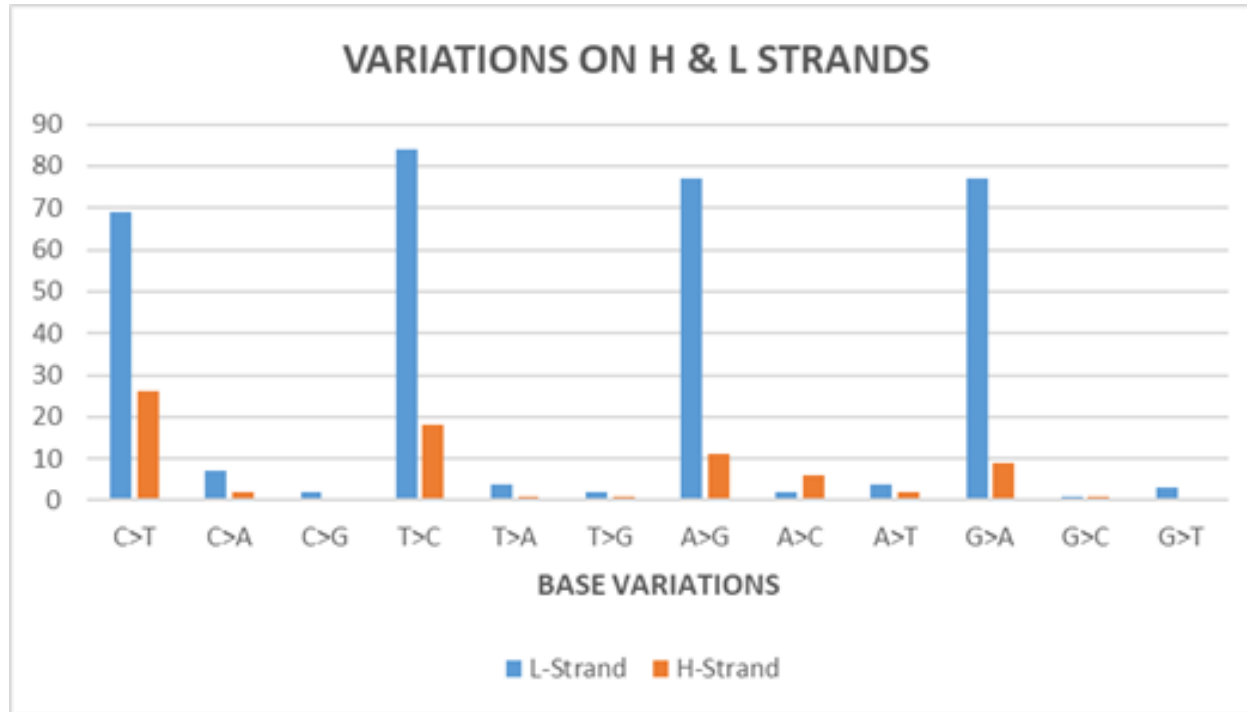
**E**



**Figure 8E:** Table indicates the number of Zoroastrian-Parsis who cluster with Persians or people of Persian origin, relic tribes of Indian origin. Pie chart indicates the percentage of clustering of the HV2a Zoroastrian-Parsis in the phylogenetic clustering analysis. Circular dendrogram of the complete Phylogenetic clustering analysis of Parsis (**Blue clades**) with Iranian mitogenomes (**Green clades**) and Indian mitogenomes (**Brown clades**). Outgroup is indicated by the black line.



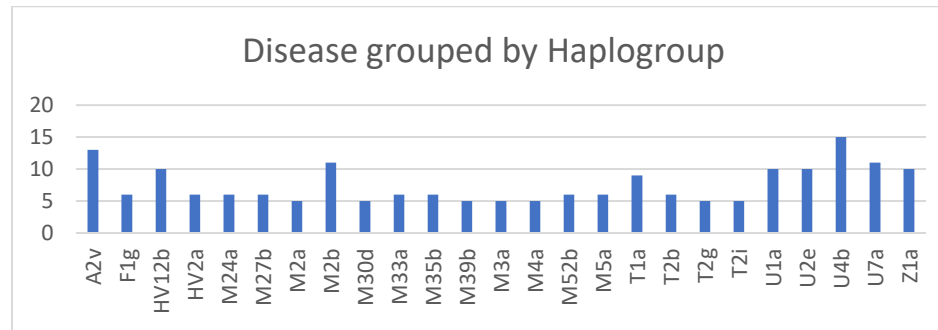
**Figure 9: Lack of smoking induced mutational signatures in the Parsi cohort**



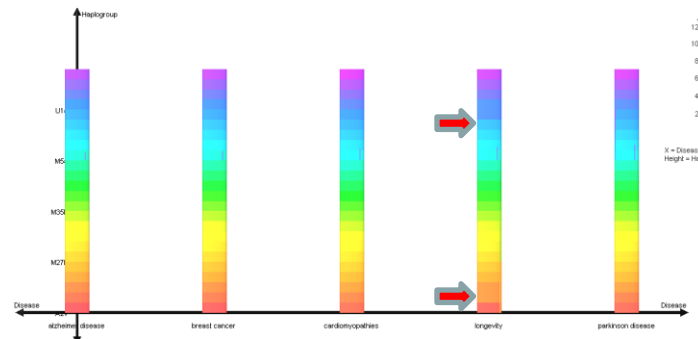
**Figure 9: Mutational signatures observed in the 100 mitochondrial genomes of Parsis.** Graph depicts the quantification of both transitions and transversions on both H&L strands of the 100 mitochondrial genomes of Parsis.

**Figure 10: Observation of Longevity variants across all sub-haplogroups and predisposition of U and M haplogroups to diseases**

A

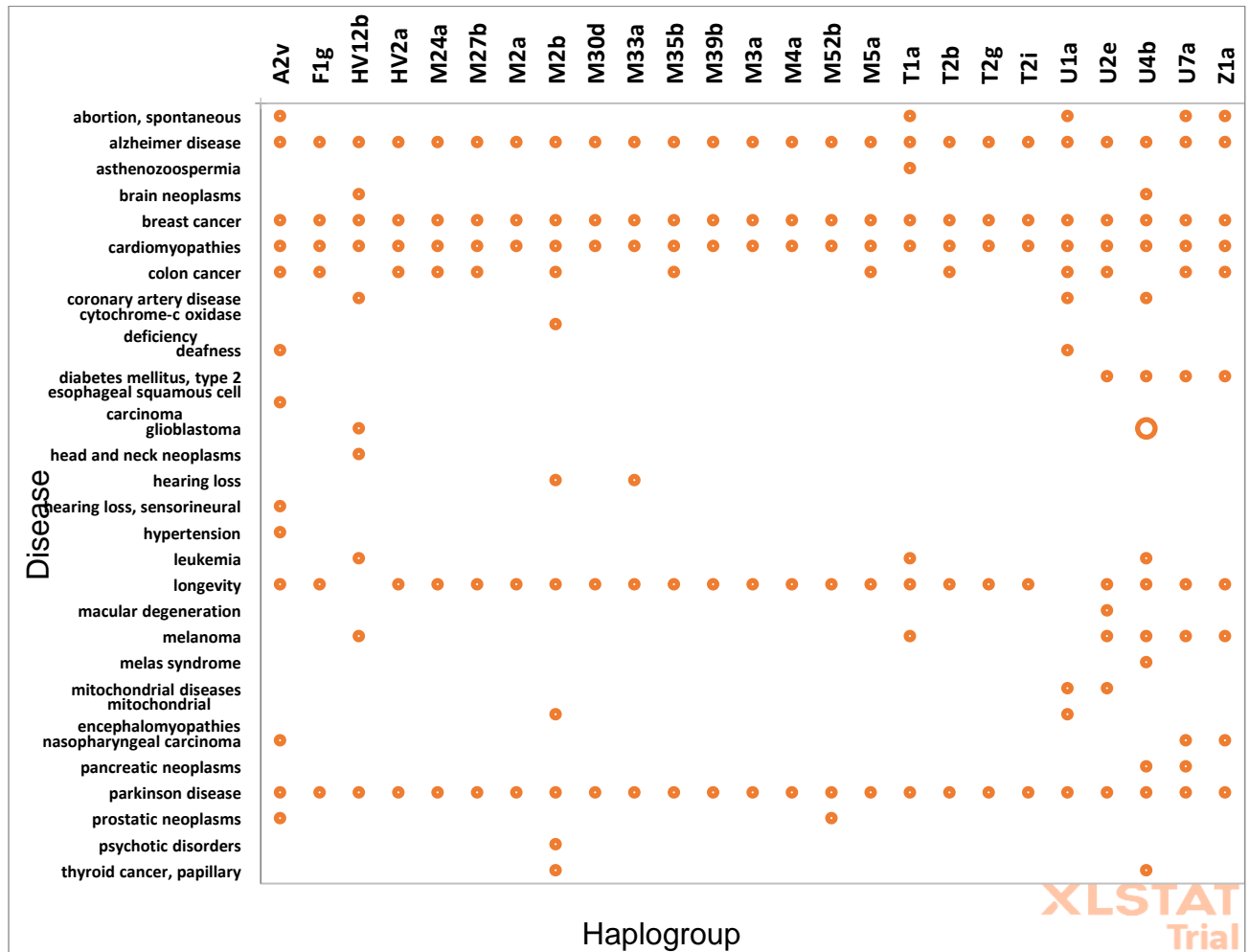


B



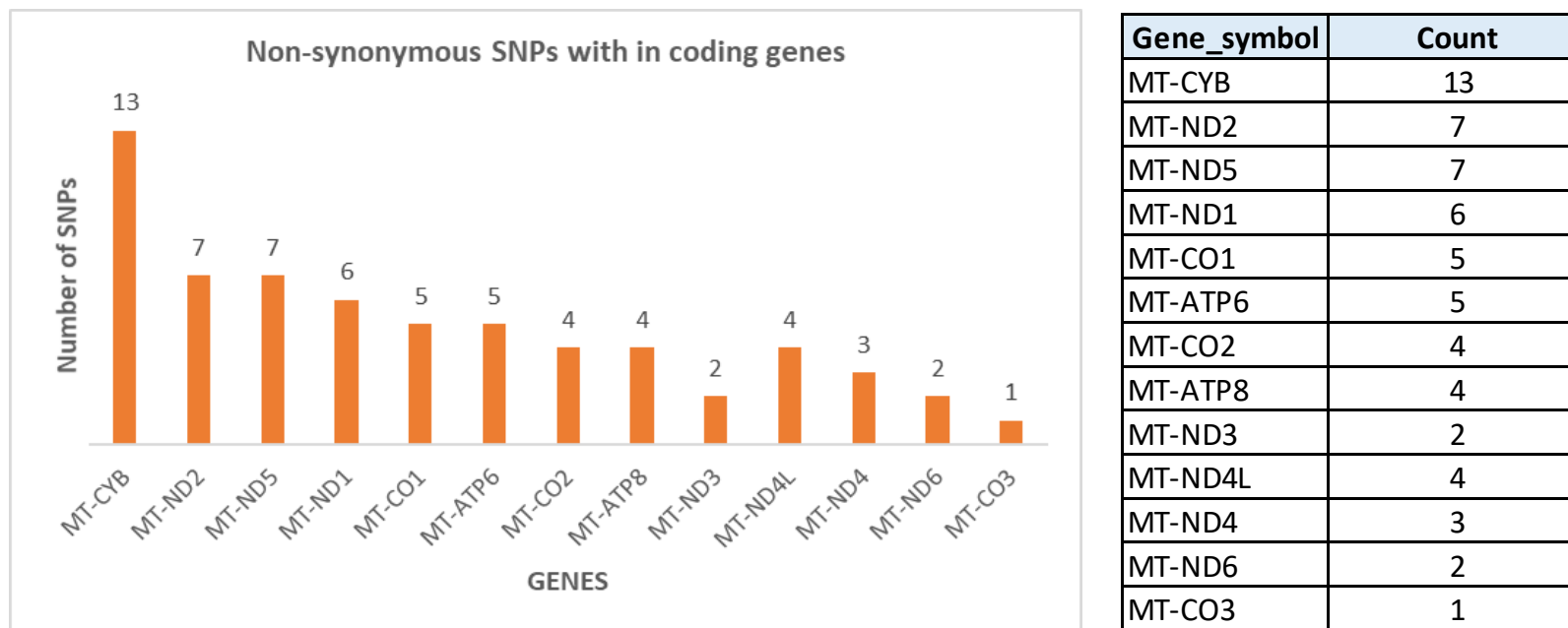
**Figure 10: Haplogroup specific distribution of diseases.** (A) Distribution of 188 diseases across 25 sub-haplogroups of the 100 Parsi subjects analyzed in this study (B) Histogram depicting longevity and disease prevalence across U1a, M52b, M35b, M27b

**Figure 11: PCA analysis shows absence of Longevity variants in U1a and F1g sub-haplogroups**



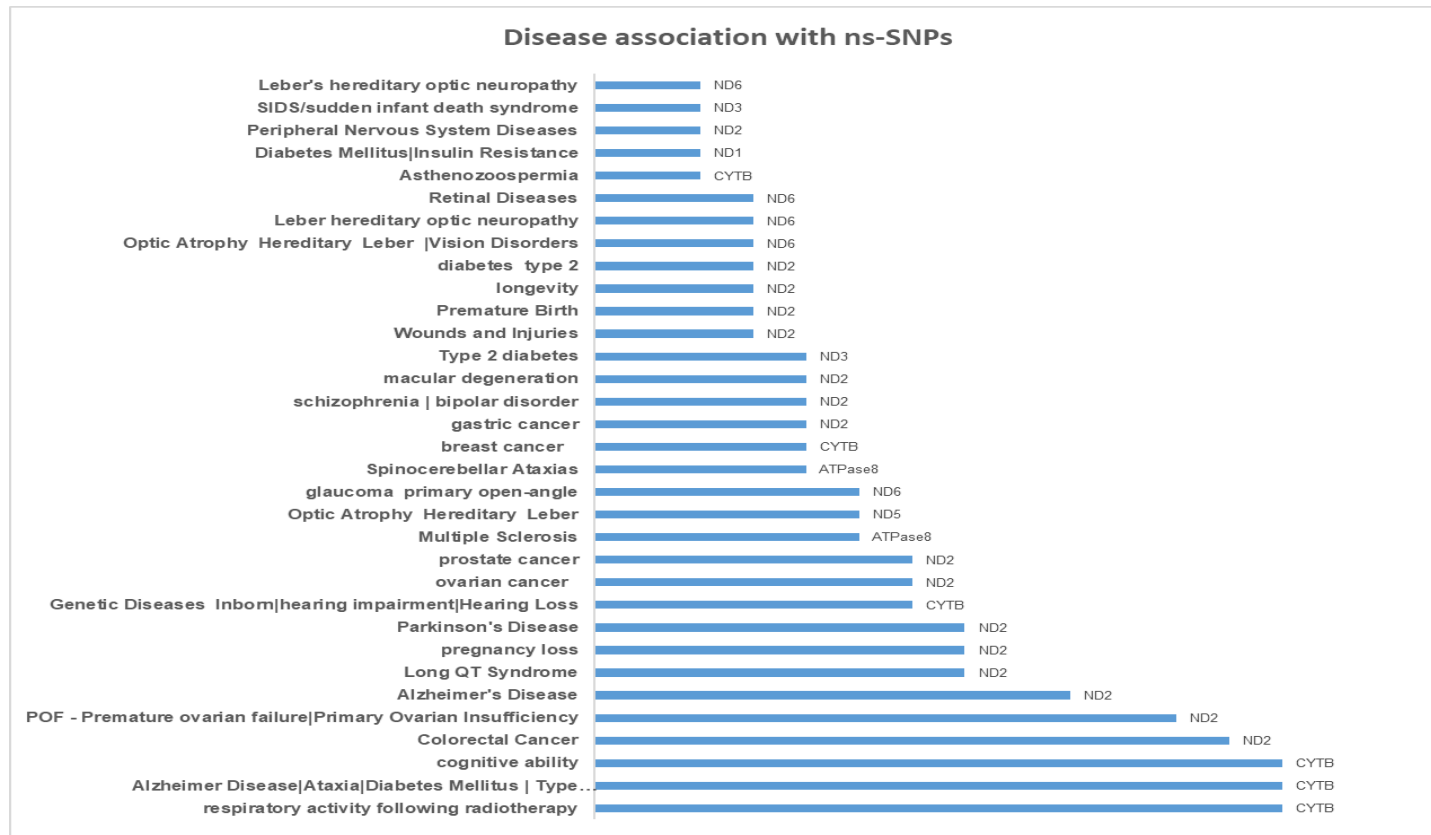
**Figure 11:** Principal Component Analysis of disease associations with sub-haplogroups in the Parsi-Zoroastrian group under study

**Figure 12: CYTB gene has the highest occurrence of non-synonymous variants in this study**



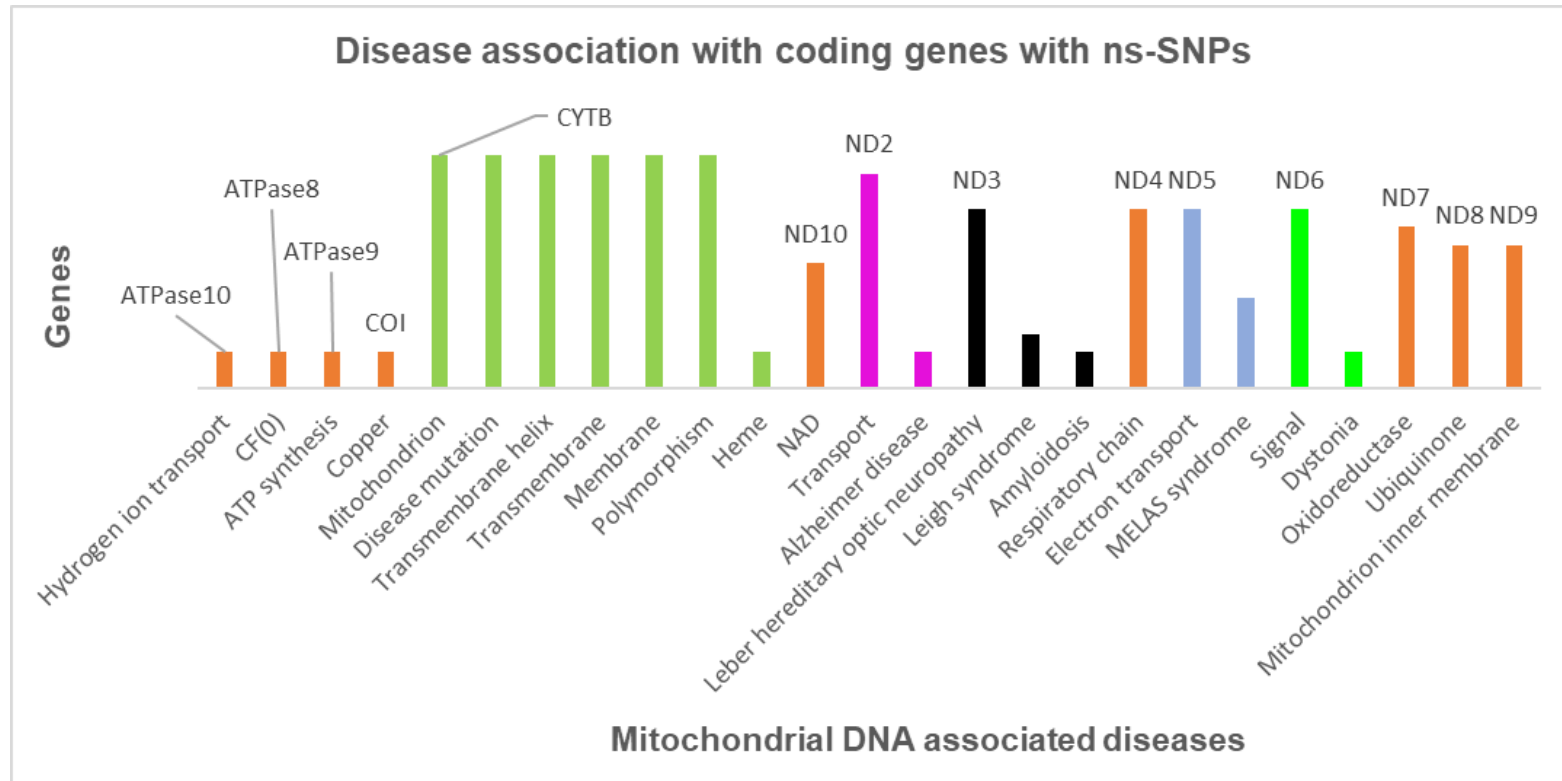
**Figure 12:** Analysis of the non-synonymous variants within 420 variants in the 100 Parsi mitochondrial genome sequences. The histogram and the table show the location of the non-synonymous variants in the coding gene loci in the mitochondrial genome analysed with MitImpact database

**Figure 13: Non-synonymous variants among 420 variants and their disease associations**



**Figure 13:** Analysis of the non-synonymous variants within 420 variants in the 100 Parsi mitochondrial genome sequences for and their disease associations.

**Figure 14: Non-synonymous variants among 420 variants and their associations with mitochondrial function**



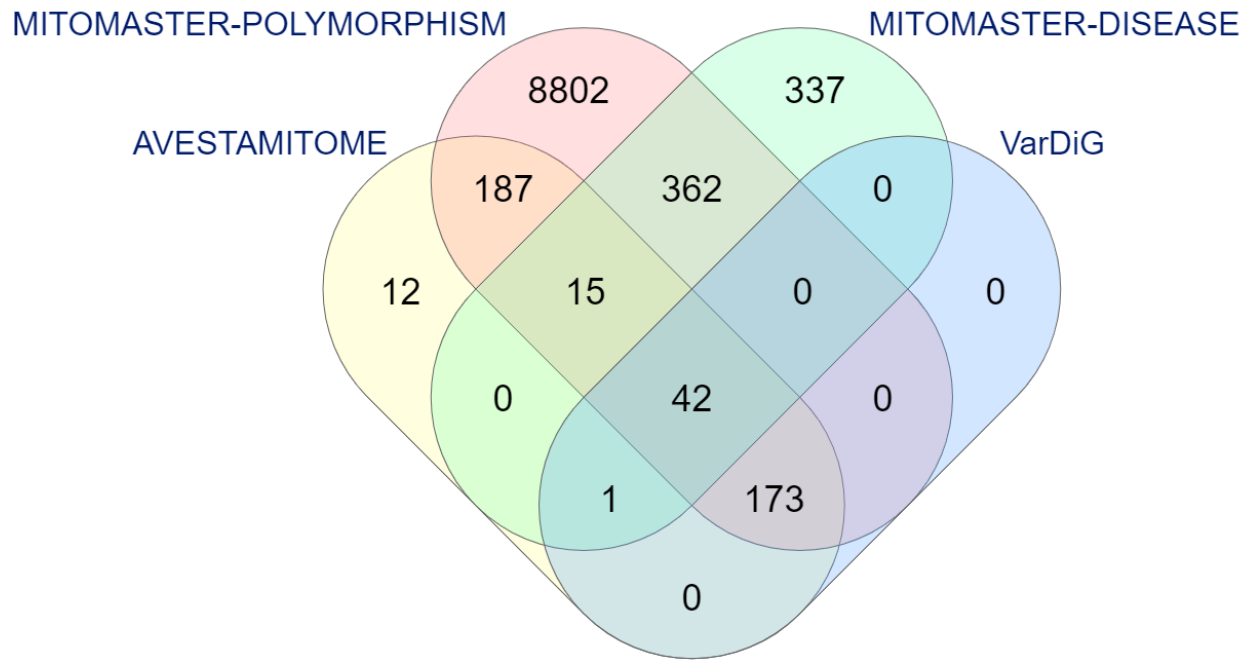
**Figure 14 : Distribution of non-synonymous Variants across coding genes.** Analysis was performed on the 420 Variants linked to the 100 Parsi mitochondrial genomes.

## Figure 15: Gene ontology associated with non-synonymous variants among 420 variants



**Figure 15:** Analysis of non-synonymous mutations and their functional classification, engagement in different pathways respectively using DAVID and UNIPROT annotation tools.

**Figure 16: 12 unique variants found in the current study**



**Figure 16: Comparative analysis of the 420 variants** in the AVESTAMITOME™ Zoroastrian-Parsi community dataset with common and disease associated polymorphisms in MITOMASTER database and VarDiG®-R



# Main Tables

**Table 1: Annotation of 28 variants in the AGENOME-ZPMS-HV2a-1**

Reference_position	72	73	152	195	263	309.1	309.2	310	750	1438	2706	4769	5075	6104	6179	7028	7193
Reference_base	T	A	T	T	A	.	.	T	A	A	A	A	T	C	G	C	T
AGENOME-ZPMS-HV2a-1	C	G	C	C	G	C	T	C	G	G	G	G	C	T	A	T	C
Mitochondrial genome loci	HVR-II								12S-rRNA:RNR1		16S-rRNA:RNR2	ND2		COI			
Amino Acid change	nc	nc	nc	nc	nc	nc	nc	nc	rRNA	rRNA	rRNA	M100M	I202I	F67F	M92M	A375A	F430F
Conservation index									98%	87%	84%	24%	44%	100%	100%	100%	100%
Protein Position												100	202	67	92	375	430
Variant Type												syn	syn	syn	syn	syn	syn
Type of base change	trans	trans	trans	trans	trans	ins	ins	trans	trans	trans	trans	trans	trans	trans	trans	trans	trans

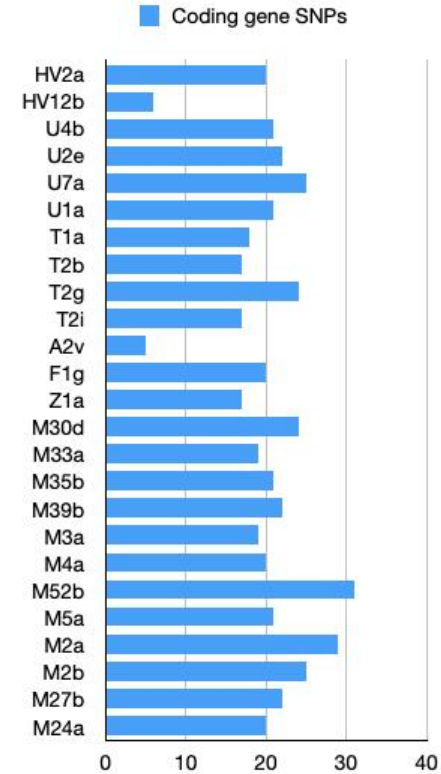
Reference_position	8860	9336	10410	11016	11935	12061	15326	15792	16153	16217	16309	Haplogroup
Reference_base	A	A	T	G	T	C	A	T	G	T	A	
ZPMS-HV-1	G	G	C	A	C	T	G	C	A	C	G	
Mitochondrial genome loci	ATPase6	COIII	tRNA [R]	ND-4			CYTB		HVR-I			
Amino Acid change	T112A	M44V	tRNA	S86N	T392T	N434N	T194A	I349T	nc	nc	nc	
Conservation index	71%	16%	22%	7%	89%	69%	18%	58%				
Protein Position	112	44		86	392	434	194	349				
Variant Type	n-syn	n-syn		n-syn	syn	syn	n-syn	n-syn				
Type of base change	trans	trans	trans	trans	trans	trans	trans	trans	trans	trans	trans	

**Table 1:** Annotation of the de novo Parsi mitochondrial genome AGENOME-ZPMS-HV2a-1. B) The table indicates the Variants (n=28) found in the AGENOME-ZPMS-HV2a-1 in relation to the revised Cambridge Reference Sequence (rCRS, Reference bases

**Table 2: Distribution of 420 variants for each sub-haplogroup for protein coding regions, D-loop of 100 Parsi mitogenomes**

**Association of coding region, D-loop with sub-haplogroup**

Sub-haplogroup	Coding gene SNPs	Gene with max SNPs	D-loop
HV2a	20	6 COI	1
HV12b	6	2 CYTB	0
U4b	21	6 COI	4
U2e	22	4 CYTB, 4 ND2, 4 ND5	2
U7a	25	6 ND5	2
U1a	21	6 ND5	2
T1a	18	5 CYTB	1
T2b	17	3 CYTB	0
T2g	24	6 CYTB	1
T2i	17	5 CYTB	1
A2v	5	2 ND2	0
F1g	20	7 CYTB	0
Z1a	17	3 ND5	1
M30d	24	8 CYTB	1
M33a	19	5 CYTB	1
M35b	21	5 CYTB	1
M39b	22	5 CYTB	0
M3a	19	5 CYTB	1
M4a	20	5 CYTB	1
M52b	31	9 CYTB	2
M5a	21	6 CYTB	1
M2a	29	7 CYTB	2
M2b	25	6 CYTB	2
M27b	22	6 CYTB	1
M24a	20	5 CYTB	1



**Table 2:** Distribution of Variants across coding genes, D-loop across all the 25 sub-haplogroup

**Table 3: Phylogenetic clustering of complete mitogenomes of Parsis with 352 Iranian and 100 relic tribes of Indian origin**

**Table: Clustering of Parsis with population of Persian and Indian descent**

Major haplogroup	Sub-haplogroups	People of Persian origin (PO)	People of Indian & Relic tribal origin (IO)	Max BS value to nearest PO	Max BS value to nearest IO
HV	HV2a	Persian	N.A	0.7270	0
	HV12b	Persian, Qashqai, Mazandarani	N.A	0.6550	0
U	U7a	Persian, Kurd, Tajik	N.A	0.8980	0
	U2e	Persian, Qashqai, Azeri	N.A	1.000	0
	U4b	Persian, Khorasani, Qashqai	N.A	0.5100	0
	U1a	Persian, Armenian	N.A	0.6850	0
T	T1a	Persian	N.A	0.7320	0
	T2g	Persian	N.A	0.4880	0
	T2i	Persian	N.A	0.4480	0
	T2b	Persian	N.A	0.4320	0
M	M5a	Persian	Munda, Mahali	0.9860	0.6270
	M39b	Unique cluster			
	M33a	Azeri	Jenu Kuruba	0.2250	0.0960
	M52b	Indian Shia Muslim	Mathakur, Dirang Monpa	0.7950	0.1170
	M24a	Persian, Qashqai	Pauri Bhaiya, Nihal	0.8560	0.0200
	M3a	Persian	N.A	0.9380	0
	M30d	Unique cluster	1 M30d with Brahmin lyengar, Bhovi	0	0.4020
	M2a	N.A	Lambadi, Hill Kolam, Katkari, Dongri Bhil	0	0.6110
	M4a	Persian	N.A	0.8560	0
	M2b	N.R	Korku, Hill Kolam	0	0.9400
	M35b	Persian	N.A	0.3860	0
	M27b	Indian Shia Muslim	N.A	0.4220	0
A	A2v	Persian	N.A	0.4690	0
F	F1g	Kurd, Turkmen	N.A	0.9970	0
Z	Z1a	Qashqai, Persian	N.A	0.2470	0

**Table 3** : Results of the Phylogenetic clustering of the 100 Parsis mitochondrial genomes with 352 mitochondrial genomes of Iranian origin and 100 mitochondrial genomes of relic tribes of Indian origin through Neighbour Joining method. BS indicates Boot-Strap values between each sample. \*N.A. indicates *No Association*, indicating a lack of representation of samples in the specific sub-haplogroup

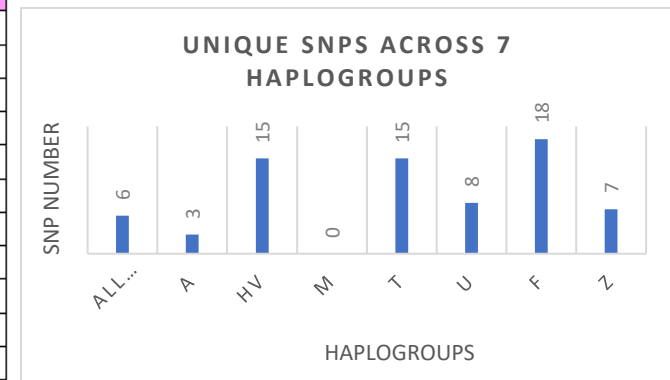
**Table 4: Variants associated with haplogroup specific Zoroastrian Parsi Mitochondrial Reference Genome (n=7) and Zoroastrian Parsi Mitochondrial Consensus Genome (n=1) mitochondrial genome sequences**

Consensus Sequence	Number of Variants	Variants
AGENOME-ZPMCG-V1.0	31	T65TT, A73G, A263G, C309CCCT, T310C, T489C, G513GCA, A567ACCCCC, A750G, A1438G, A2706G, A3158AT, A4769G, C7028T, A8701G, A8860G, T9540C, A10398G, C10400T, T10873C, G11719A, C12705T, C14766T, T14783C, G15043A, G15301A, A15326G, C16169CC, A16182AC, C16223T, T16519C
AGENOME-ZPMRG-A2v-V1.0	11	A263G, C309CCT, T310C, A750G, A1438G, A4769G, A8860G, C11881T, A15326G, C16168T, C16239T
AGENOME-ZPMRG-HV-V1.0	26	T72C, A73G, T152C, T195C, A263G, C309CCCT, T310C, A750G, A1438G, A2706G, A4769G, T5075C, C6104T, G6179A, C7028T, T7193C, A8860G, A9336G, T10410C, G11016A, T11935C, C12061T, A15326G, T15792C, T16217C, A16309G
AGENOME-ZPMRG-M-V1.0	29	T65TT, A73G, A263G, C309CCCT, T310C, T489C, A567ACCCC, A750G, A1438G, A2706G, A4769G, C7028T, A8701G, A8860G, T9540C, A10398G, C10400T, T10873C, G11719A, C12705T, C14766T, T14783C, G15043A, G15301A, A15326G, C16169CC, A16182AC, C16223T, T16519C
AGENOME-ZPMRG-U-V1.0	25	A73G, A263G, C309CCCT, T310C, G499A, G513GCA, A567ACCCCC, A750G, A1438G, A1811G, A2706G, A3158AT, A4769G, C7028T, A8860G, C11332T, A11467G, G11719A, A12308G, G12372A, C14620T, C14766T, A15326G, T16189TT, T16519C
AGENOME-ZPMRG-T-V1.0	28	A73G, A263G, C309CCT, T310C, G709A, A750G, A1438G, G1888A, A2706G, T4216C, A4769G, A4917G, C7028T, G8697A, A8860G, T10463C, A11251G, G11719A, G13368A, C14766T, G14905A, A15326G, C15452A, A15607G, G15928A, T16126C, C16294T, T16519C
AGENOME-ZPMRG-F1g-V1.0	32	A73G, A248d, A263G, C315CC, CA514d, A750G, A1438G, C2389T, A2706G, T3398C, C3970T, T3999C, A4769G, T6392C, G6962A, C7028T, A8589G, A8860G, G10310A, T10609C, G11719A, G12406A, C12882T, G13928C, C14766T, A15326G, T15916C, A16183C, T16189C, C16193CC, T16304C, T16519C
AGENOME-ZPMRG-Z-V1.0	33	A73G, C151T, T152C, A263G, C315CC, T489C, A750G, A1438G, A2072d, A2706G, A4769G, C7028T, A8701G, A8860G, T9540C, A10149T, A10398G, C10400T, C10556T, T10873C, G11719A, G12007A, C12705T, C14766T, T14783C, G15043A, G15301A, A15326G, G15346A, T15784C, C16223T, T16311C, T16519C

**Table 4:** List of unique variants associated with the Haplogroup specific Zoroastrian Parsi Mitochondrial Reference Genomes (ZPMRG) for A2v, HV, M, U, T, F1g, Z and overall unique variants in the Zoroastrian Parsi Mitochondrial Consensus Genome (ZPMCG)

**Table 5: Variants associated with Zoroastrian Parsi Mitochondrial Reference Genome (ZPMRG) and unique variants of each ZPMRG compared to Zoroastrian Parsi Mitochondrial Consensus Genome (ZPMCG)**

AGENOME-ZPMRG-A2v-V1.0	AGENOME-ZPMRG-HV-V1.0	AGENOME-ZPMRG-M-V1.0	AGENOME-ZPMRG-T-V1.0	AGENOME-ZPMRG-U-V1.0	AGENOME-ZPMRG-F-V1.0	AGENOME-ZPMRG-F-V1.0
C11881T	A16G		C6G	A21G	A248d	C151T
C16168T	T72C		G709A	G499A	CA514d	A2072d
C16239T	T195C		G1888A	A1811G	C2389T	C10556T
	T5075C		T4216C	C11332T	T3398C	G12007A
	C6104T		A4917G	A11467G	C3970T	G15346A
	G6179A		G8697A	A12308G	T3999C	T15784C
	T7193C		T10463C	G12372A	T6392C	T16311C
	A9336G		A11251G	C14620T	G6962A	
	T10410C		G13368A		A8589G	
	G11016A		G14905A		G10310A	
	T11935C		C15452A		T10609C	
	C12061T		A15607G		G12406A	
	T15792C		G15928A		C12882T	
	T16217C		T16126C		G13928C	
	A16309G		C16294T		T15916C	
					A16183C	
					C16193CC	
					T16304C	



**Table 5:** (A) Unique Variants found in the haplogroup specific Reference Genomes (ZPMRG) compared to the Zoroastrian-Parsi Consensus Genome (AGENOME-ZPMCG-V1). The histogram (right) lists the exact number of variants in each ZPMRG compared to ZPMCG

**Table 6: mt-t-RNA variants in our study and their disease association**

mt-tRNA	Variation	Probability_of_pathogenicity	Classification	Frequency %	Haplogroup	Disease association
Phe	T593C	0.16	Neutral	0.06	M52b	Non-syndromic hearing loss (Reported)
Val	G1644A	0.67	Pathogenic	0.01	U4b	LS/HCM/MELAS (Reported)
Val	T1654C	0.12	Neutral	0.01	M3a	
Met	T4454C	0.13	Neutral	0.02	M5a	Possible contributor to mito dysfunction / Hypertension (Reported)
Asp	G7521A	0.46	Likely neutral	0.01	U4b	
Asp	T7561C	0.33	Neutral	0.01	U7a	
Asp	T7581C	0.42	Likely neutral	0.01	U1a	
Arg	T10410C	0.17	Neutral	0.14	Hv2a	
Arg	T10463C	0.31	Neutral	0.04	T1a,T2g,T2i	
His	A12172G	0.53	Likely pathogenic	0.01	U4b	
His	C12191G	0.11	Neutral	0.01	M27b	
Leu(CUN)	A12279G	0.37	Likely neutral	0.06	M52b	
Leu(CUN)	A12308G	0.41	Likely neutral	0.21	U4b,U7a	Stroke, CM, CPEO, Breast/Renal/Prostate cancer risk, Altered brain pH(Reported)
Glu	A14696G	0.26	Neutral	0.01	A2v	Progressive Encephalopathy (Reported)
Thr	A15907G	0.23	Neutral	0.03	U2e	
Thr	T15908C	0.5	Likely pathogenic	0.01	M33a	Deaf Helper mutation (Reported)
Thr	T15916C	0.33	Likely neutral	0.01	F1g	

**Table 6:** Analysis of the occurrence of the 420 variants in the tRNA and their disease associations annotated with the PON-mt-tRNA database. A frequency score  $\geq 0.5$  – pathogenic,  $=0.5$  – likely pathogenic,  $<0.5$  – neutral

