

1 **Genome-wide analysis of mitochondrial DNA copy number reveals multiple loci**  
2 **implicated in nucleotide metabolism, platelet activation, and megakaryocyte proliferation**

3 Longchamps RJ<sup>1\*</sup>, Yang SY<sup>1\*</sup>, Castellani CA<sup>1,2</sup>, Shi W<sup>1</sup>, Lane J<sup>3</sup>, Grove ML<sup>4</sup>, Bartz  
4 TM<sup>5</sup>, Sarnowski C<sup>6</sup>, Burrows K<sup>7,8</sup>, Guyatt AL<sup>9</sup>, Gaunt TR<sup>7,8</sup>, Kacprowski T<sup>10</sup>, Yang  
5 J<sup>11</sup>, De Jager PL<sup>12,13</sup>, Yu L<sup>11</sup>, CHARGE Aging and Longevity Group, Bergman A<sup>14</sup>,  
6 Xia R<sup>15</sup>, Fornage M<sup>15,16</sup>, Feitosa MF<sup>17</sup>, Wojczynski MK<sup>17</sup>, Kraja AT<sup>17</sup>, Province MA<sup>17</sup>, Amin  
7 N<sup>18</sup>, Rivadeneira F<sup>19</sup>, Tiemeier H<sup>18,20</sup>, Uitterlinden AG<sup>18,19</sup>, Broer L<sup>19</sup>, Van Meurs JBJ<sup>18,19</sup>, Van  
8 Duijn CM<sup>18</sup>, Raffield LM<sup>21</sup>, Lange L<sup>22</sup>, Rich SS<sup>23</sup>, Lemaitre RN<sup>24</sup>, Goodarzi  
9 MO<sup>25</sup>, Sitlani CM<sup>24</sup>, Mak ACY<sup>26</sup>, Bennett DA<sup>11</sup>, Rodriguez S<sup>7,8</sup>, Murabito JM<sup>27</sup>, Lunetta  
10 KL<sup>6</sup>, Sotoodehnia N<sup>28</sup>, Atzmon G<sup>29</sup>, Kenny Y<sup>30</sup>, Barzilai N<sup>31</sup>, Brody JA<sup>32</sup>, Psaty BM<sup>33</sup>, Taylor  
11 KD<sup>34</sup>, Rotter JI<sup>34</sup>, Boerwinkle E<sup>4,35</sup>, Pankratz N<sup>3</sup>, Arking DE<sup>1</sup>

12

13 **\* Indicates authors contributed equally to this work.**

14 1 McKusick-Nathans Institute, Department of Genetic Medicine, Johns

15 Hopkins University School of Medicine, Baltimore, MD

16 2 Department of Pathology and Laboratory Medicine, Western University, London, ON, Canada

17 3 Department of Laboratory Medicine and Pathology, University of Minnesota Medical School,

18 Minneapolis, MN

19 4 Human Genetics Center, Department of Epidemiology, Human Genetics, and Environmental

20 Sciences, School of Public Health, The University of Texas Health Science Center at Houston,

21 Houston, TX

22 5 Cardiovascular Health Research Unit, Departments of Medicine and Biostatistics, University of

23 Washington, Seattle, WA

24 6 Department of Biostatistics, Boston University School of Public Health, Boston, MA

25 7 MRC Integrative Epidemiology Unit at the University of Bristol, University of Bristol, Oakfield

26 House, Oakfield Grove, Bristol, UK

- 27 8 Population Health Sciences, Bristol Medical School, University of Bristol, Oakfield House,  
28 Oakfield Grove, Bristol, UK
- 29 9 Department of Health Sciences, University of Leicester, University Road, Leicester, UK
- 30 10 Department of Functional Genomics, Interfaculty Institute for Genetics and Functional  
31 Genomics, University of Greifswald, Greifswald, Germany; Division Data Science in  
32 Biomedicine, Peter L. Reichertz Institute for Medical Informatics, TU Braunschweig and  
33 Hannover Medical School, Brunswick, Germany
- 34 11 Rush Alzheimer's Disease Center & Department of Neurological Sciences, Rush University  
35 Medical Center, Chicago, IL
- 36 12 Center for Translational and Systems Neuroimmunology, Department of Neurology,  
37 Columbia University Medical Center, New York, NY
- 38 13 Program in Medical and Population Genetics, Broad Institute, Cambridge, MA
- 39 14 Department of Systems and Computational Biology, Albert Einstein College of Medicine,  
40 Bronx, New York, USA
- 41 15 Institute of Molecular Medicine, The University of Texas health Science Center at Houston,  
42 Houston TX
- 43 16 Human Genetics Center, The University of Texas Health Science Center at Houston,  
44 Houston
- 45 17 Division of Statistical Genomics, Department of Genetics, Washington University School of  
46 Medicine
- 47 18 Department of Epidemiology, Erasmus Medical Center, Rotterdam, The Netherlands
- 48 19 Department of Internal Medicine, Erasmus Medical Center, Rotterdam, The Netherlands
- 49 20 Department of Social and Behavioral Science, Harvard T.H. School of Public Health, Boston,  
50 USA
- 51 21 Department of Genetics, University of North Carolina at Chapel Hill, Chapel Hill, NC

- 52 22 Department of Medicine, University of Colorado Denver, Anschutz Medical Campus, Aurora,  
53 CO
- 54 23 Center for Public Health Genomics, University of Virginia, Charlottesville, Virginia, USA
- 55 24 Cardiovascular Health Research Unit, Department of Medicine, University of Washington,  
56 Seattle, WA
- 57 25 Division of Endocrinology, Diabetes and Metabolism, Cedars-Sinai Medical Center, Los  
58 Angeles, CA
- 59 26 Cardiovascular Research Institute and Institute for Human Genetics, University of California,  
60 San Francisco, California
- 61 27 Boston University School of Medicine, Boston University, Boston, MA
- 62 28 Cardiovascular Health Research Unit, Division of Cardiology, University of Washington,  
63 Seattle, WA
- 64 29 Department of Natural science, University of Haifa, Haifa, Israel; Departments of Medicine  
65 and Genetics, Albert Einstein College of Medicine, Bronx, NY, 10461, USA.
- 66 30 Department of Epidemiology and Population Health, Albert Einstein College of Medicine,  
67 Bronx, NY, 10461, USA.
- 68 31 Departments of Medicine and Genetics, Albert Einstein College of Medicine, Bronx, NY,  
69 10461, USA.
- 70 32 Cardiovascular Health Research Unit, Department of Medicine, University of Washington,  
71 Seattle, WA
- 72 33 Cardiovascular Health Research Unit, Departments of Epidemiology, Medicine and Health  
73 Services, University of Washington, Seattle, WA
- 74 34 The Institute for Translational Genomics and Population Sciences, Department of Pediatrics,  
75 The Lundquist Institute for Biomedical Innovation at Harbor-UCLA Medical Center, Torrance,  
76 CA
- 77 35 Baylor College of Medicine, Human Genome Sequencing Center, Houston, TX

78

## 79 **Abstract**

80           Blood-derived mitochondrial DNA copy number (mtDNA-CN) is a minimally invasive  
81 proxy measure of mitochondrial function that exhibits both inter-individual and intercellular  
82 variation. While mtDNA-CN has been previously associated with various aging-related diseases,  
83 little is known about the genetic factors that may modulate this phenotype. We performed a  
84 genome-wide association study (GWAS) in 465,809 individuals of White (European) ancestry  
85 from the Cohorts for Heart and Aging Research in Genomic Epidemiology (CHARGE)  
86 consortium and the UK Biobank (UKB). We identified 129 SNPs with statistically significant,  
87 independent effects associated with mtDNA-CN across 96 loci. A combination of fine-mapping,  
88 variant annotation, co-localization, and gene set enrichment analyses were used to prioritize  
89 genes within each of the 129 independent sites. Putative causal genes were enriched for known  
90 mitochondrial DNA depletion syndromes ( $p = 3.09 \times 10^{-15}$ ) and the gene ontology (GO) terms for  
91 mtDNA metabolism ( $p = 1.43 \times 10^{-8}$ ) and mtDNA replication ( $p = 1.2 \times 10^{-7}$ ). A clustering  
92 approach leveraged pleiotropy between mtDNA-CN associated SNPs and 42 mtDNA-CN  
93 associated phenotypes to identify functional domains, revealing five distinct groups, including  
94 platelet activation, megakaryocyte proliferation, and mtDNA metabolism. In conclusion, in a  
95 GWAS of mtDNA-CN conducted in >450,000 individuals, we identified SNPs within loci that  
96 implicate novel pathways that provide a framework for defining the underlying mechanisms  
97 involved in genetic control of mtDNA-CN.

98

99

## 100 **Introduction**

101 Mitochondria are the cellular organelles primarily responsible for producing the chemical energy  
102 required for metabolism, as well as signaling the apoptotic process, maintaining homeostasis,  
103 and synthesizing several macromolecules such as lipids, heme and iron-sulfur clusters<sup>1,2</sup>.

104 Mitochondria possess their own genome (mtDNA); a circular, intron-free, double-stranded,  
105 haploid, ~ 16.6 kb maternally inherited molecule encoding 37 genes vital for proper  
106 mitochondrial function. Due to the integral role of mitochondria in cellular metabolism,  
107 mitochondrial dysfunction is known to play a critical role in the underlying etiology of several  
108 aging-related diseases<sup>3-5</sup>.

109  
110 Unlike the nuclear genome, a large amount of variation exists in the number of copies of mtDNA  
111 present within cells, tissues, and individuals. The relative copy number of mtDNA (mtDNA-CN)  
112 has been shown to be positively correlated with oxidative stress<sup>6</sup>, energy reserves, and  
113 mitochondrial membrane potential<sup>7</sup>. As a minimally invasive proxy measure of mitochondrial  
114 dysfunction<sup>8</sup>, decreased blood-derived mtDNA-CN has been previously associated with aging-  
115 related disease states including frailty<sup>9</sup>, cardiovascular disease<sup>10-12</sup>, chronic kidney disease<sup>13</sup>,  
116 neurodegeneration<sup>14,15</sup>, and cancer<sup>16</sup>.

117  
118 Although mtDNA-CN measured from whole blood presents itself as an easily accessible and  
119 minimally invasive biomarker, cell type composition has been shown to be an important  
120 confounder, complicating analyses<sup>17,18</sup>. For example, while platelets generally have fewer  
121 mtDNA molecules than leukocytes, the lack of a platelet nuclear genome drastically skews  
122 mtDNA-CN estimates. As a result, not only is controlling for cell composition extremely vital for  
123 accurate mtDNA-CN estimation, but interpreting the results in relation to the impact of cell  
124 composition becomes a necessity<sup>18-20</sup>.

125  
126 Although the comprehensive mechanism through which mtDNA-CN is modulated is largely  
127 unknown<sup>21,22</sup>, twin studies have estimated broad-sense heritability ~0.65, consistent with  
128 moderate genetic control<sup>23</sup>. Several nuclear genes have been shown to directly modulate  
129 mtDNA-CN, specifically those within the mtDNA replication machinery such as the mitochondrial

130 polymerase, *POLG* and *POLG2*<sup>24,25</sup>, as well as the mitochondrial DNA helicase, *TWNK*, and the  
131 mitochondrial single-stranded binding protein, *mtSSB*<sup>26</sup>. Furthermore, nuclear genes which  
132 maintain proper mitochondrial nucleotide supply including *DGUOK* and *TK2* have also been  
133 shown to regulate mtDNA-CN<sup>27-29</sup>. To further elucidate the genetic control over mtDNA-CN,  
134 several genome-wide association studies (GWAS) of mtDNA-CN have been published<sup>30-33</sup>,  
135 including a study that was published while the current manuscript was in preparation, analyzing  
136 ~300,000 participants from the UK Biobank (UKB), and identifying 50 independent loci<sup>33</sup>.

137

138 In the present study, we report mtDNA-CN GWAS results from 465,809 individuals across the  
139 Cohorts for Heart and Aging Research in Genomic Epidemiology (CHARGE) consortium<sup>34</sup> and  
140 the UK Biobank (UKB)<sup>35</sup>. Using multiple gene prioritization and functional annotation methods,  
141 we assign genes to loci that reach genome-wide significance. Finally, we perform gene  
142 expression analyses and gene-set enrichment through PHEWAS-based SNP-clustering to  
143 identify functional domains related to mtDNA-CN.

144

## 145 **Subjects and Methods**

146

### 147 **Study Populations**

148 470,579 individuals participated in this GWAS, 465,809 of whom self-identified as White.

149 Participants were derived from 7 population-based cohorts representing the Cohorts for Heart  
150 and Aging Research in Genetic Epidemiology (CHARGE) consortium (Avon Longitudinal Study  
151 of Parents and Children [ALSPAC], Atherosclerosis Risk in Communities [ARIC], Cardiovascular  
152 Health Study [CHS], Multi-Ethnic Study of Atherosclerosis [MESA], Religious Orders Study and  
153 Memory and Aging Project [ROSMAP], Study of Health in Pomerania [SHIP]) and from the UK  
154 Biobank (UKB) (Supplemental Table 1). Detailed descriptions of each participating cohort, their  
155 quality control practices, study level analyses, and ethic statements are available in the

156 Supplemental Methods. All study participants provided written informed consent and all centers  
157 obtained approval from their institutional review boards.

158

## 159 **Methods for Mitochondrial DNA Copy Number Estimation (CHARGE cohorts)**

### 160 ***qPCR***

161 mtDNA-CN was determined using a quantitative PCR assay as previously described<sup>32,36</sup>. Briefly,  
162 the cycle threshold (Ct) value of a nuclear-specific and mitochondrial-specific probe were  
163 measured in triplicate for each sample. In CHS, a multiplex assay using the mitochondrial *ND1*  
164 probe and nuclear *RPPH1* probe was used, whereas ALSPAC used a mitochondrial probe  
165 targeting the D-Loop and a nuclear probe targeting *B2M*. In CHS, we observed plate effects, as  
166 well as a linear increase in  $\Delta$ Ct due to the pipetting order of each replicate. These effects were  
167 corrected in the analysis using linear mixed model regression, with pipetting order included as a  
168 fixed effect and plate as a random effect to create a raw measure of mtDNA-CN. In ALSPAC,  
169 run-to-run variability was controlled using 3 calibrator samples added to every plate, to allow for  
170 adjustment by a per-plate calibration factor<sup>32</sup>.

171

172

### 173 ***Microarray***

174 Microarray probe intensities were used to estimate mtDNA-CN using the Genvisis software  
175 package<sup>37</sup> as previously described<sup>10,36</sup>. Briefly, Genvisis uses the median mitochondrial probe  
176 intensity across all homozygous mitochondrial SNPs as an initial estimate of mtDNA-CN.  
177 Technical artifacts such as DNA input quality, DNA input quantity, and hybridization efficiency  
178 were captured through either surrogate variable (SV) or principal component (PC) analyses.  
179 SVs or PCs were adjusted for through stepwise linear regression by adding successive  
180 components until each successive surrogate variable or principal component no longer  
181 significantly improved the model.

182

### 183 ***Whole Genome Sequencing (ARIC)***

184 Whole genome sequencing read counts were used to estimate mtDNA-CN as previously  
185 described<sup>36</sup>. Briefly, the total number of reads in a sample were web scraped from the NCBI  
186 sequence read archive. Mitochondrial reads were downloaded directly from dbGaP through  
187 Samtools (1.3.1). There was no overlap between ARIC microarray and ARIC whole-genome  
188 sequencing samples. A ratio of mitochondrial reads to total aligned reads was used as a raw  
189 measure of mtDNA-CN.

190

### 191 ***Adjusting for Covariates***

192 Each method described above represents a raw measure of mtDNA-CN, adjusted for technical  
193 artifacts; however, several potential confounding variables (e.g., age, sex, blood cell  
194 composition) have been identified previously<sup>18</sup>. Raw mtDNA-CN values were adjusted for white  
195 blood cell count in ARIC, SHIP and CHS (which also adjusted for platelet count), depending on  
196 available data. Standardized residuals (mean = 0, standard deviation = 1) of mtDNA-CN were  
197 used after adjusting for covariates (Supplemental Table 1.

198

### 199 ***Estimation of Mitochondrial DNA Copy Number (UKB)***

200 Due to the availability of more detailed cell count data, as well as a different underlying  
201 biochemistry for the Affymetrix Axiom array compared to the genotyping arrays used in the  
202 CHARGE cohorts, mtDNA-CN in the UKB was estimated differently (Supplemental Methods).  
203 Briefly, mtDNA-CN estimates derived from whole exome sequencing data, available on ~50,000  
204 individuals, were generated first using customized Perl scripts to aggregate the number of  
205 mapped sequencing reads and correct for covariates through both linear and spline regression  
206 models. Concurrently, mitochondrial probe intensities from the Affymetrix Axiom arrays,  
207 available on the full ~500,000 UKB cohort, were adjusted for technical artifacts through principal



208 components generated from nuclear probe intensities. Probe intensities were then regressed  
209 onto the whole exome sequencing mtDNA-CN metric, and beta estimates from that regression  
210 were used to estimate mtDNA-CN in the full UKB cohort. Finally, we used a 10-fold cross  
211 validation method to select the cell counts to include in the final model (Supplemental Table 2).  
212 The final UKB mtDNA-CN phenotype is the standardized residuals (mean = 0, standard  
213 deviation = 1) from a linear model adjusting for covariates (age, sex, cell counts) as described in  
214 the Supplemental Methods.

215

### 216 **Genome-Wide Association Study**

217 For each individual cohort, regression analysis was performed with residualized mtDNA-CN as  
218 the dependent variable adjusting for age, sex, and cohort-specific covariates (e.g., principal  
219 components, DNA collection site, family structure, cell composition). Cohorts with multiple  
220 mtDNA-CN estimation platforms were stratified into separate analyses. Ancestry-stratified meta-  
221 analyses were performed using Metasoft software using the Han and Eskin random effects  
222 model to control for unobserved heterogeneity due to differences in mtDNA-CN estimation  
223 method<sup>38</sup>. Effect size estimates for SNPs were calculated using a random effect meta-analysis  
224 from cohort summary statistics, as the Han and Eskin model relaxes the assumption under the  
225 null hypothesis without modifying the effect size estimates that occur under the alternative  
226 hypothesis<sup>38</sup>. In total, three complementary analyses were performed in self-identified White  
227 individuals, (1) a meta-analysis using all available studies, (2) a meta-analysis of studies with  
228 available data for cell count adjustments, and (3) an analysis of UKB-only data. As the vast  
229 majority of samples are derived from the UKB study, and given the difficulty in interpreting effect  
230 size estimates from a random effects model, further downstream analyses were all performed  
231 using effect size estimates from UKB-only data.

232

### 233 **Identification of Independent GWAS Loci**

234 To identify the initial genome-wide significant (lead) SNPs in each locus, the most significant  
235 SNP that passed genome-wide significance ( $p < 5 \times 10^{-8}$ ) within a 1 Mb window was selected.  
236 To avoid Type I error, SNPs were only retained for further analyses if there were either (a) at  
237 least two genome-wide significant SNPs in the 1 Mb window or (b) if the lead SNP was directly  
238 genotyped. Conditional analyses were performed in UKB, where the lead SNPs from the original  
239 GWAS were used as additional covariates in order to identify additional independent  
240 associations.

241

## 242 **Fine-mapping**

243 The susieR package was used to identify all potential causal variants for each independent  
244 locus associated with mtDNA CN<sup>39</sup>. UKB imputed genotype data for unrelated White subjects  
245 were used and variants were extracted using a 500 kb window around the lead SNP for each  
246 locus with minor allele frequency (MAF) > 0.001. 95% credible sets (CS) of SNPs, containing a  
247 potential causal variant within a locus, were generated. The minimum absolute correlation within  
248 each CS is 0.5 and the scaled prior variance is 0.01. When the CS did not include the lead SNP  
249 identified from the GWAS, some of the parameters were slightly relaxed [minimum absolute  
250 correlation is 0.2, estimate prior variance is TRUE]. The SNP with the highest posterior inclusion  
251 probability (PIP) within each CS was also identified (Supplemental Table 3). With a few  
252 exceptions, final lead SNPs were selected by prioritizing initially identified SNPs unless the SNP  
253 with the highest PIP had a PIP greater than 0.2 and was 1.75 times larger than the SNP with the  
254 second highest PIP.

255

## 256 **Functional Annotation and Gene Prioritization**

### 257 ***Functional Annotation***

258 ANNOVAR was used for functional annotation of variants identified in the fine-mapping step<sup>40</sup>.  
259 First, variants were converted to an ANNOVAR-ready format using the dbSNP version 150

260 database<sup>41</sup>. Then, variants were annotated with ANNOVAR using the RefSeq Gene database<sup>42</sup>.  
261 The annotation for each variant includes the associated gene and region (e.g., exonic, intronic,  
262 intergenic). For intergenic variants, ANNOVAR provides flanking genes and the distance to  
263 each gene. For exonic variants, annotations also include likely functional consequences (e.g.,  
264 synonymous/nonsynonymous, insertion/deletion), the gene affected by the variant, and the  
265 amino acid sequence change (Supplemental Table 4).

266

### 267 ***Co-localization Analyses***

268 Co-localization analyses were performed using the approximate Bayes factor method in the R  
269 package *coloc*<sup>43</sup>. Briefly, *coloc* utilizes eQTL data and GWAS summary statistics to evaluate  
270 the probability that gene expression and GWAS data share a single causal SNP (colocalize).  
271 *Coloc* returns multiple posterior probabilities; H0 (no causal variant), H1 (causal variant for gene  
272 expression only), H2 (causal variant for mtDNA-CN only), H3 (two distinct causal variants), and  
273 H4 (shared causal variant for gene expression and mtDNA-CN). In the event of high H4, we can  
274 assume that the gene is potentially causal for the GWAS phenotype of interest. eQTL summary  
275 statistics were obtained from the eQTLGen database<sup>44</sup>. Genes with significant associations with  
276 lead SNPs were tested for co-localization using variants within a 500 kb window of the sentinel  
277 SNP. Occasionally, some of the eQTLGen p-values for certain SNPs were identical due to R's  
278 (ver 4.0.3) limitation in handling small numbers. To account for this, if the absolute value for a  
279 SNP's z-score association with a gene was greater than 37.02, z-scores were rescaled so that  
280 the largest z-score would result in a p-value of  $5 \times 10^{-300}$ . Additionally, a few clearly co-localized  
281 genes did not result in high H4 PPs due to the strong effect for each phenotype of a single SNP  
282 (Supplemental Figure 1), possibly due to differences in linkage disequilibrium (LD) between the  
283 populations. To account for this, we summed mtDNA-CN GWAS p-values and eQTLGen p-  
284 values for each SNP and removed the SNP with the lowest combined p-value. Co-localization

285 analyses were then repeated without the lowest SNP. Genes with H4 greater than 50% were  
286 classified as genes with significant evidence of co-localization.

### 287 **DEPICT**

288 Gene prioritization was performed with Depict, an integrative tool that incorporates gene co-  
289 regulation and GWAS data to identify the most likely causal gene at a given locus<sup>45</sup>. Across  
290 GWAS SNPs which overlapped with the DEPICT database, we identified SNPs representing  
291 119 independent loci with LD pruning defined as  $p < 5 \times 10^{-8}$ ,  $r^2 < 0.05$  and  $> 500$  kb from other  
292 locus boundaries. Only genes with a nominal p-value of less than 0.05 were considered for  
293 downstream prioritization.

### 294 **Gene Assignment**

295 To prioritize genes for each identified locus, we utilized functional annotations, eQTL co-  
296 localization analyses, and DEPICT gene prioritization results (Supplemental Figure 2). First,  
297 genes with missense variants within *SusieR* fine-mapped credible sets were assigned to loci. If  
298 loci co-localized with a gene's expression with a posterior probability (PP) of greater than 0.50  
299 and there were no other co-localized genes with a PP within 5%, the gene with the highest  
300 posterior probability was assigned. If there was still no assigned gene, the most significant  
301 DEPICT gene was assigned. If there was no co-localization or DEPICT evidence, the nearest  
302 gene was assigned.

303

### 304 **Gene Set Enrichment Analyses**

305 Using the finalized gene list from the prioritization pipeline, GO and KEGG pathway enrichment  
306 analyses were performed using the "goana" and "kegga" functions from the R package *limma*<sup>46</sup>,  
307 treating all known genes as the background universe<sup>47</sup>. Only one gene per locus was used for  
308 "goana" and "kegga" gene set enrichment analysis, prioritizing genes assigned to primary  
309 independent hits. If there were multiple assigned genes, one gene was randomly selected to  
310 avoid biasing results through loci with multiple genes. To identify an appropriate p-value cutoff,

311 96 genes were randomly selected from the genome and run through the same enrichment  
312 analysis. This permutation was repeated 1000 times to generate a null distribution of the  
313 smallest p-values from each permutation. For cluster 5 gene set enrichment analyses,  
314 permutation testing used 47 random genes. To ensure robustness of results, gene set  
315 enrichment analysis was repeated 50 times with random selection of genes at loci with multiple  
316 assigned genes. GO and KEGG terms that passed permutation cutoffs at least 40/50 times  
317 were retained.

318

### 319 **Gene-based Association Test**

320 We used metaXcan, which employs gene expression prediction models to evaluate associations  
321 between phenotypes and gene expression<sup>48</sup>. We obtained pre-calculated expression prediction  
322 models and SNP covariance matrices, computed using whole blood from European ancestry  
323 individuals in version 7 of the Genotype-Tissue expression (GTEx) database<sup>49</sup>. Using prediction  
324 performance p-values of less than 0.05, a total of 6,285 genes were predicted. Of these genes,  
325 74 passed Bonferroni correction of  $p < 7.95 \times 10^{-6}$ . Gene set enrichment analyses were  
326 performed on Bonferroni-significant genes as previously described. REVIGO<sup>50</sup> was used on the  
327 “medium” setting (allowed similarity = 0.7) to visualize significantly enriched GO terms.

328

329 We used a one-sided Fisher’s exact test to test for enrichment of genes that have been  
330 previously identified as causal for mtDNA depletion syndromes<sup>51–53</sup>.

331

### 332 **PHEWAS-based SNP Clustering**

#### 333 ***mtDNA-CN Phenome-wide Association Study (PHEWAS)***

334 We used the PHEnome Scan ANalysis Tool (PHESANT)<sup>54</sup> to identify mtDNA-CN associated  
335 quantitative traits in the UKB. Briefly, we tested for the association of mtDNA-CN with 869  
336 quantitative traits (Supplemental Table 5), limiting analyses to 365,781 White, unrelated

337 individuals (used.in.pca.calculation=1), and excluding individuals with extreme cell count  
338 measurements (see Supplementary Methods). Analyses were adjusted for age, sex, and  
339 assessment center.

340

### 341 ***SNP Phenotype Associations***

342 SNP genotypes were regressed on mtDNA-associated quantitative traits using linear  
343 regression, adjusted for sex, age with a natural spline (df=2), assessment center, genotyping  
344 array, and 40 principal components (provided as part of the UKB data download).

345

### 346 ***SNP Clustering***

347 To identify distinct clusters of mtDNA-CN GWS SNPs based on phenotypic associations, beta  
348 estimates from the SNP phenotype associations were first divided by the beta estimate of the  
349 mtDNA-CN SNP-mtDNA-CN association, so that all SNP-phenotype associations are relative to  
350 the mtDNA-CN increasing allele and scaled to the effect of the SNP on mtDNA-CN. The  
351 adjusted beta estimates were subjected to a dimensionality reduction method, Uniform Manifold  
352 and Approximation Projection (UMAP), as implemented in the R package *umap*<sup>55</sup>  
353 (random\_state=123, n\_neighbors=10, n\_components=2, n\_epochs=5000). SNPs were  
354 assigned to clusters using Density Based Clustering of Applications with Noise (DBSCAN) as  
355 implemented in the R package *dbscan*<sup>56</sup> (minPts=3). Clusters represent groups of SNPs with  
356 similar phenotypic associations.

357

358 To identify reproducible sub-clusters within cluster 5, we ran *umap* 100 times, varying the initial  
359 random state (n\_neighbors=4, n\_epochs=5000, min\_dist=0.05, n\_components=3). For each  
360 *umap* run, sub-clusters were assigned using *dbscan* as described above. Final SNP sub-cluster  
361 assignments were determined by identifying SNPs that all mapped to the same cluster >50% of  
362 the time (i.e., all pairwise correlations for each SNP within the cluster against all other SNPs in

363 the cluster was >50%). The final umap run for plotting in 2 dimensions used the following  
364 parameters: random\_state=123, n\_neighbors=4, n\_epochs=5000, min\_dist=0.05.

365

### 366 ***Phenotype Enrichment and Permutation Testing***

367 To test for enrichment of specific phenotypes within clusters, we compared the median mtDNA-  
368 CN scaled phenotype beta estimates within the cluster to the median beta estimates for all  
369 SNPs not in the cluster, with significance determined using 20,000 permutations in which cluster  
370 assignment was permuted. For multi-test correction of the overall cluster and sub-cluster  
371 analyses, we performed 300 permutations of the initial cluster assignment (separately for the  
372 cluster and sub-cluster analyses), followed by the comparison of median beta estimates as  
373 described above. We retained only the most significant result from across all phenotypes and  
374 clusters from each of the 300 permutations, and then selected the 15<sup>th</sup> most significant value as  
375 the study-wide threshold for multi-test corrected significance of  $p < 0.05$ .

376

377 All statistical analyses were performed using R version 4.0.3

378

## 379 **Results**

### 380 **Sample Characteristics**

381 The current study included 465,809 individuals of White (European) ancestry (53.9% female)  
382 with an average age of 56.6 yrs (sd = 8.2 yrs) (Supplemental Table 1). Follow-up validation  
383 analyses were performed in 4,770 Blacks (60.2% female) with an average age of 61.2 yrs (sd =  
384 7.4 yrs). The majority of the data originated from the UKB (93%). The bulk of the DNA used for  
385 mtDNA-CN estimation was derived from buffy coat (95.5%) while the rest was derived from  
386 peripheral leukocytes (2.2%) or whole blood (2.3%). mtDNA-CN estimated from Affymetrix  
387 genotyping arrays consisted of 97.9% of the data while the remainder was derived from qPCR-  
388 (1.8%) and WGS (0.3%).

389

## 390 **GWAS**

391 Previous work has demonstrated that the method used to measure mtDNA-CN can impact the  
392 strength of association<sup>36</sup>. To account for potential differences across studies due to the  
393 different mtDNA-CN measurements used, as well the inclusion of blood cell counts in only a  
394 subset of the cohorts, we took two approaches. First, we used a random effects model to  
395 perform meta-analyses, allowing for different genetic effect size estimates across cohorts.  
396 Second, we performed three complementary analyses in individuals who self-identified as  
397 White: 1) meta-analysis of all available studies (n = 465,809); 2) meta-analysis of studies  
398 with available data for cell count adjustment (n = 456,151); and 3) GWAS of UKB only (n =  
399 440,266). 77 loci were significant in all three meta-analyses, and we identified 93 independent  
400 loci, of which 92 were genome-wide significant in the UKB data alone (Supplemental Figure 3,  
401 Figure 1). Given that > 90% of the samples come from the UKB study, and the challenge of  
402 interpreting effect size estimates from a random effects model, downstream analyses all use  
403 effect size estimates from the UKB only (Supplemental Table 6).

404 The most significant SNP associated with mtDNA-CN was a missense mutation in  
405 *LONP1* ( $p = 3.00 \times 10^{-141}$ ), a gene that encodes a mitochondrial protease that can directly bind  
406 mtDNA, and has been shown to regulate *TFAM*, a transcription factor involved in mtDNA  
407 replication and repair (for review see Gibellini *et al.*)<sup>57</sup>.

408

## 409 **Fine-mapping and Secondary Hits**

410 To identify additional independent SNPs within novel loci whose effects were masked by the  
411 original significant SNP, as well as identify additional loci, we took two approaches. First, a  
412 conditional analysis adjusting for the top 93 SNPs from the initial (primary) GWAS run revealed  
413 3 novel loci and 19 additional independent significant SNPs within existing loci. We also  
414 performed fine-mapping with SuSiE<sup>39</sup> and discovered an additional 14 independent SNPs within



415 existing loci. The majority of loci had only one 95% credible set of SNPs; further, twenty of the  
416 credible sets contained only one SNP. However, many of the credible sets contained greater  
417 than 50 SNPs after fine-mapping, and 12 of the 122 credible sets had a missense SNP as the  
418 SNP with the highest PIP in the set. Using these two methods, we identified in total 129  
419 independent SNPs across 96 loci (Supplemental Figure 4).

420

### 421 **Associations in African Ancestry (AA) Populations**

422 Examining the 129 SNPs from the Whites-only analysis, 99 were available in the Blacks-  
423 only meta-analysis ( $n = 4770$ ). After multiple testing correction, one of these SNPs was  
424 significant (rs73349121,  $p = 0.0001$ ), 9 were nominally significant ( $p < 0.05$ , with 5 expected),  
425 and 58/99 had a direction of effect that was consistent with the White-only analyses (one-sided  
426  $p = 0.04$ , Figure 2). Despite being under-powered, these results in the Black-only analyses  
427 provide evidence for similar genetic effects in a different ancestry.

428

### 429 **Gene Prioritization and Enrichment of mtDNA Depletion Syndrome Genes**

430 We integrated results from three different gene prioritization and functional annotation methods  
431 (ANNOVAR<sup>40</sup>, COLOC<sup>43</sup>, and DEPICT<sup>45</sup>) so that loci with nonsynonymous variants in gene  
432 exons were prioritized first, with eQTL co-localization results considered second (Supplemental  
433 Table 7), and those from DEPICT (Supplemental Table 8) were considered last (Supplemental  
434 Figure 2). For 20 loci, multiple genes were assigned as analyses could not identify a single  
435 priority gene (Supplemental Table 9). We noted the identification of a number of mtDNA  
436 depletion syndrome genes in the priority list and tested for enrichment of these known causal  
437 genes using a one-sided Fisher's exact test. For this analysis, all genes for loci assigned to  
438 multiple genes were used, and genes for all primary and secondary loci were considered. Our  
439 gene prioritization approach identified 7 of 16 mtDNA depletion genes (Supplemental Table 10),  
440 consistent with a highly significant enrichment (one-sided  $p = 3.09 \times 10^{-15}$ ).

441

## 442 **Gene Set Enrichment Analyses**

443 To avoid bias from a single locus with multiple functionally related genes contributing to a false-  
444 positive signal, only one gene per unique locus was used, prioritizing genes assigned to primary  
445 loci. One gene was randomly selected for loci with multiple assigned genes. To test for  
446 robustness of gene set enrichment results, random selection was repeated 50 times, and only  
447 gene sets that were significantly enriched for at least 40 iterations were retained. In all, a total of  
448 96 genes were utilized for GO term and KEGG pathway enrichment analyses. Using a  
449 Bonferroni-corrected p-value cutoff, 15 gene sets were significantly enriched for all 50 iterations,  
450 including mitochondrial DNA metabolic process, mitochondrial DNA replication, coagulation,  
451 hemostasis, amyloid-beta clearance, and mitochondrial genome maintenance (Supplemental  
452 Table 11). No KEGG terms were significant across multiple iterations.

453

## 454 **MetaXcan Gene Expression Analysis**

455 As a complementary approach to single-SNP analyses, we explored the associations between  
456 mtDNA-CN and predicted gene expression using MetaXcan<sup>48</sup> MetaXcan incorporates multiple  
457 SNPs within a locus along with a reference eQTL dataset to generate predicted gene  
458 expression levels. As our study estimated mtDNA-CN derived from blood, we used whole blood  
459 gene expression eQTLs from the Gene-Tissue Expression (GTEx) consortium<sup>58</sup> to predict gene  
460 expression in the UKB dataset. We identified 6,285 genes that had a predicted performance p-  
461 value of less than 0.05 (i.e., they had sufficient data to generate robust gene expression levels)  
462 and were tested for association with mtDNA-CN. Of these genes, 74 were significantly  
463 associated with mtDNA-CN ( $p < 7.95 \times 10^{-6}$ , Figure 3), including 8 that were not identified  
464 through single-SNP analyses. Many of the significant genes have known mitochondrial  
465 functions, notably the mtDNA transcription factor *TFAM* ( $p = 1.09 \times 10^{-29}$ ) and mitochondrial  
466 exonuclease *MGME1* ( $p = 5.87 \times 10^{-23}$ ), genes known as causal for mtDNA depletion

467 syndromes<sup>51,52</sup>. Additionally, *LONP1*, *MRPL43*, and *BAK1*, are all genes with known  
468 mitochondrial functions<sup>59-61</sup>. Bonferroni significant MetaXcan genes were used for gene  
469 enrichment analysis, finding enrichment for “nucleobase metabolic process” ( $p = 1.47 \times 10^{-4}$ )  
470 and “mitochondrial fusion” ( $p = 1.86 \times 10^{-4}$ , Supplemental Figure 5).

471

## 472 **PHEWAS-based SNP Clustering and Gene Set Enrichment**

473 mtDNA-CN is associated with numerous quantitative and qualitative phenotypes, many of which  
474 are relevant to aging-related disease<sup>3-5,9,10,13-16</sup>. We hypothesized that this pleiotropy may  
475 reflect different underlying functional domains captured by mtDNA-CN that may be reflected in  
476 GWAS-identified SNPs and their likely causal genes. To test this hypothesis, we used the UKB  
477 data to identify quantitative traits associated with mtDNA-CN and selected 42 highly significant,  
478 non-redundant traits to test for association with the mtDNA-CN GWAS SNPs (Supplemental  
479 Table 5, in PHEWAS = 1). We clustered SNPs using the trait effect size (beta) divided by the  
480 mtDNA-CN effect size estimate, so that all effects are standardized to the effect of the mtDNA-  
481 CN increasing allele for each locus. We identified 5 clusters of SNPs (Figure 4A), with clusters  
482 1, 2, and 3 containing SNPs in which the mtDNA-CN increasing allele is associated with  
483 decreased platelet count (PLT) (Figure 4B), increased mean platelet volume (MPV) (Figure 4C),  
484 and platelet distribution width (PDW) (Figure 4D), consistent with a role in platelet activation<sup>62</sup>.  
485 Cluster 4 is most strongly enriched for SNPs in which the mtDNA-CN increasing allele is  
486 associated with increased PLT, plateletcrit (PLTCRIT, a measure of total platelet mass), serum  
487 calcium (Figure 4E), serum phosphate, as well as decreased mean corpuscular volume (MCV)  
488 and mean spherical cellular volume (Figure 4F) (Supplemental Table 12). The cluster 4  
489 phenotypes, and supported by the genes identified for this cluster, implicate megakaryocyte  
490 proliferation and proplatelet formation (*MYB*, *AK3*, *JAK2*)<sup>63</sup>, and apoptosis and autophagy  
491 (*BAK1*, *BCL2*, *TYMP*)<sup>64</sup>. Cluster 5 did not yield any specific trait enrichment (all significant  
492 results reflected the strong enrichment observed in clusters 1-4); however, gene set enrichment

493 for this cluster identified multiple mtDNA-related gene ontology terms, including mitochondrial  
494 genome maintenance, regulation of phospholipid efflux, and amyloid-beta clearance  
495 (Supplemental Table 13). Cluster 5 contains a highly diverse set of SNP-trait associations. A  
496 secondary clustering only with SNPs mapped to this cluster identified 9 sub-clusters  
497 (Supplemental Figure 6A), 3 having clear evidence for enrichment of specific phenotypes.  
498 Cluster 6 is enriched for decreased MCV (Supplemental Figure 6B). Cluster 7 is enriched for  
499 decreased MPV (Supplemental Figure 6C), PDW, serum phosphate, serum calcium, and  
500 PLTCRT. Cluster 10 is enriched for decreased apoA, apoB, total cholesterol (Supplemental  
501 Figure 6D), LDL, and HDL, and increased vitamin D, pulse rate, and direct bilirubin. SNPs that  
502 were not assigned to a cluster were enriched for decreased fasting blood glucose and maximum  
503 workload from fitness testing (Supplemental Table 14).

504

## 505 **Discussion**

506 We conducted a GWAS for mtDNA-CN using 465,809 individuals from the CHARGE consortium  
507 and the UKB. In addition to replicating two previously reported loci, we discovered 94 novel loci  
508 and report multiple independent hits for 26 loci. Examining our GWS SNPs in a Black  
509 population, we observed a concordant signal, suggesting that the genetic etiology of mtDNA-CN  
510 may be broadly similar across populations. Using several functional follow-up methods, genes  
511 were assigned for each identified independent hit and significant enrichment was observed for  
512 genes involved in mitochondrial DNA metabolism, homeostasis, cell activation, and amyloid-  
513 beta clearance. In total, we assigned 124 unique genes to independent GWAS signals  
514 associated with mtDNA-CN. We also identified 8 additional genes whose predicted gene  
515 expression is associated with mtDNA-CN that could not be mapped back to GWS loci. Finally,  
516 using a clustering approach based on SNP associations with various mtDNA-CN associated  
517 phenotypes, we were able to functionally categorize SNPs, providing insight into biological  
518 pathways that impact mtDNA-CN. We note that during the preparation of this manuscript, a

519 GWAS including 295,150 unrelated individuals from the UK Biobank was published, which  
520 reported 50 genome-wide significant regions<sup>33</sup>. While many of our loci overlap with their  
521 findings, our study reports twice as many loci largely due to the increased power of our study.

522 We were able to identify a substantial proportion of the genes involved in mtDNA  
523 depletion syndromes (7/16,  $p = 3.09 \times 10^{-15}$  for enrichment), including *TWINK*, *TFAM*, *DGUOK*,  
524 *MGME1*, *RRM2B*, *TYMP*, and *POLG*. mtDNA depletion syndromes can be broken down into 5  
525 subtypes based on their constellation of phenotypes<sup>65</sup>, and with the exception of  
526 cardiomyopathic subtypes (associated with mutations in *AGK* and *SLC25A4*), we were able to  
527 identify at least 1 gene from the other 4 subtypes, suggesting that our mtDNA-CN measurement  
528 in blood-derived DNA can identify genes widely relevant to non-blood phenotypes. This finding  
529 is consistent with a large body of work showing that mtDNA-CN measured in blood is  
530 associated with numerous aging-related phenotypes for which the primary tissue of interest is  
531 not blood (e.g. chronic kidney disease<sup>13</sup>, heart failure<sup>11</sup>, and diabetes<sup>66</sup>). Also consistent with this  
532 finding is recent work demonstrating that mtDNA-CN measured in blood is associated with  
533 mtRNA expression across numerous non-blood tissues, suggesting a link between  
534 mitochondrial function measured in blood and other tissues<sup>67</sup>.

535 In addition to identifying the mtDNA depletion syndrome genes directly linked to  
536 mitochondrial DNA metabolic processes, DNA replication, and genome maintenance, we also  
537 identify genes which play a role in mitochondrial function. The top GWAS hit is a missense  
538 mutation in *LONP1*, which encodes a mitochondrial protease that has been shown to cause  
539 mitochondrial cytopathy and reduced respiratory chain activity<sup>68,69</sup>. Interestingly, this missense  
540 mutation was recently found to be associated with mitochondrial tRNA methylation levels<sup>70</sup>.  
541 Additional genes known to impact mitochondrial function include *MFN1*, which encodes a  
542 mediator of mitochondrial fusion<sup>71,72</sup>, *STMP1*, which plays a role in mitochondrial respiration<sup>73</sup>,  
543 and *MRPS35*, which encodes a ribosomal protein involved in protein synthesis in the  
544 mitochondrion<sup>74,75</sup>.

545 Using a combination of gene-based tests and gene prioritization using functional  
546 annotation, pathway analyses reveal enrichment for numerous mitochondrial related pathways,  
547 as well as those involved in regulation of cell activation ( $p < 3.65 \times 10^{-5}$ ), homeostatic processes  
548 ( $p < 1.82 \times 10^{-5}$ ), and regulation of immune system processes ( $p < 2.75 \times 10^{-5}$ ) (Supplemental  
549 Table 11). These results provide additional evidence for the broad role played by mitochondria  
550 in numerous aspects of cellular function. Of particular interest, the GO term for amyloid beta is  
551 significantly enriched, reinforcing a link between mtDNA-CN and neurodegenerative disease<sup>76–</sup>  
552 <sup>78</sup>. Previous work from our lab using the UKB has shown that increased mtDNA-CN is  
553 associated with lower rates of prevalent neurodegenerative disease, as well as predictive for  
554 decreased risk of incident neurodegenerative disease<sup>67</sup>. mtDNA-CN is also known to be  
555 decreased in the frontal cortex of Alzheimer's disease (AD) patients<sup>79</sup>. Interestingly, the four  
556 GWAS-identified genes driving the enrichment for amyloid-beta clearance are all related to  
557 regulation of lipid levels, and lipid homeostasis within the brain is known to play an important  
558 role in Alzheimer's disease<sup>80</sup>. *APOE*, one of the most well-known risk genes for Alzheimer's  
559 disease, is a cholesterol carrier involved in lipid transport, and the ApoE- $\epsilon$ 4 isoform involved in  
560 AD pathogenesis is associated with mitochondrial dysfunction and oxidative distress in the  
561 human brain<sup>81</sup>; *CD36* is a platelet glycoprotein which mediates the response to amyloid-beta  
562 accumulation<sup>82</sup>; *LDLR* is a low-density lipoprotein receptor associated with AD<sup>83</sup>; and *ABCA7* is  
563 a phospholipid transporter<sup>84</sup>. *ABCA7* loss of function variants are enriched in both AD and  
564 Parkinson's disease (PD) patients<sup>85</sup>, suggesting a broad role across neurodegenerative  
565 diseases.

566 Given the integral role of mitochondria in cellular function, not just with ATP  
567 formation/energy production, but signaling through ROS, and its key role in apoptosis, there is a  
568 strong reason to *a priori* assume that genetic variants associated with mtDNA-CN are likely to  
569 be highly pleiotropic. Indeed, mtDNA-CN itself is associated with numerous phenotypes  
570 (Supplemental Table 5). Through our PHEWAS-based clustering approach using 42 mtDNA-CN

571 associated phenotypes, we uncovered phenotypic associations between five distinct clusters of  
572 GWS mtDNA-CN associated SNPs. Cluster 1-3 were characterized by increased MPV, PDW,  
573 and decreased PLT (note that measured MPV and PLT are generally inversely correlated to  
574 maintain hemostasis), which are the hallmarks of platelet activation<sup>62</sup>. The link between platelets  
575 and mtDNA-CN has typically revolved around platelet count, as platelets have functional  
576 mitochondria, but do not have a nucleus. Given that the mtDNA-CN measurement is the ratio  
577 between mtDNA and nuclear DNA, increased platelets, all else being equal, would directly  
578 equate with increased mtDNA-CN. We note that the mtDNA-CN metric used in this GWAS was  
579 adjusted for platelet count, likely increasing the ability to detect variants that impact mtDNA-CN  
580 through increased platelet activation. Examining the genes within this cluster suggests role for  
581 actin formation/regulation (*CARMIL1*, *TPM4*, *PACSIN2*)<sup>86-88</sup> and vesicular transport/endocytic  
582 trafficking (*DNM3*, *EHD3*)<sup>89,90</sup> in platelet activation.

583 Cluster 4 is most strongly enriched for SNPs in which the mtDNA-CN increasing allele is  
584 associated with increased PLT/PLTCRIT and serum calcium/phosphate. Examining the genes  
585 assigned to the cluster, we implicate megakaryocyte proliferation and proplatelet formation  
586 (*MYB*, *AK3*, *JAK2*)<sup>63</sup>, and apoptosis and autophagy (*BAK1*, *BCL2*, *TYMP*)<sup>64</sup>. Megakaryocytes  
587 are used to form proplatelets, and the process includes an important role for both intra- and  
588 extracellular calcium levels<sup>91</sup>. A role for apoptosis, and specifically *BCL2*, in proplatelet  
589 formation and platelet release has been suggested<sup>92,93</sup>, however work in mice has suggested  
590 that apoptosis does not play a direct role in these processes<sup>94</sup>. Nevertheless, apoptosis is  
591 important for platelet lifespan<sup>95</sup>.

592 Cluster 5 was particularly challenging to interpret, given that no particular phenotype was  
593 enriched relative to the non-cluster 5 SNPs. We note that this cluster appeared to be enriched  
594 for the mtDNA depletion syndrome genes, containing 6/7 genes identified in the GWAS, and  
595 significantly enriched for GO Terms related mitochondrial DNA. Examining sub-clusters within  
596 cluster 5 identified MCV as an important phenotype (while only sub-cluster 6 was formally

597 associated, looking at Figure 3B suggests more widespread clustering based on MCV). MCV is  
598 a measure of the average volume of a red blood corpuscle, and the link between red blood  
599 volume and mtDNA-CN is unlikely to be direct, given that red blood cells contain neither nuclei  
600 nor mitochondria. Sub-cluster 7 is surprisingly enriched for variants associated with decreased  
601 MPV/PDW and serum phosphate/calcium, inverse of what we observe for cluster 4. Finally, sub-  
602 cluster 10 is strongly associated with lipids, and includes a number of genes known to be  
603 associated with lipid levels (*LIPC*, *CETP*, *LDLR*, *APOE*). While lipids play a role in both energy  
604 metabolism (largely through fatty acids) and cellular membrane formation, a link to mtDNA-CN  
605 and/or mitochondrial function is not well-established. One potentially interesting result is  
606 provided by Olkowicz and colleagues, who demonstrated that ApoE<sup>-/-</sup>/LDLR<sup>-/-</sup> mice had  
607 increased cardiac mitochondrial oxidative metabolism, with proteomic analysis suggesting  
608 increased mitochondrial abundance in mouse hearts<sup>96</sup>. However, we note that our results show  
609 an association between decreased lipids and increased mtDNA-CN, rather than the reverse,  
610 shown in the Olkowicz study.

611 Although the question of causality is of great interest, this study was unable to determine  
612 the directionality of effect between mtDNA-CN and phenotypes of interest, as we are  
613 underpowered for Mendelian randomization (MR), with less than 1% of the variance in mtDNA-  
614 CN explained by GWS loci when predicted into ARIC. As an additional limitation, we note that  
615 despite the large sample size and numerous loci identified, we are likely missing a great deal of  
616 the true signal, as previous studies have estimated mtDNA-CN heritability to be 65%<sup>23</sup>. Finally,  
617 while we have adjusted our mtDNA-CN metric for a variety of confounders, it is important to  
618 note that mtDNA-CN can be influenced by a variety of environmental factors including  
619 smoking<sup>97</sup> and drugs, which have not been adjusted for in these analyses.

620 In summary, we performed the largest-to-date GWAS for mtDNA-CN, including almost  
621 500,000 individuals. We identified distinct groups of SNPs associated with mtDNA-CN that are  
622 related to platelet activation, megakaryocyte formation and apoptotic processes, and showed



623 clear enrichment for genes involved in mtDNA depletion and nucleotide regulation. Given the  
624 role of mtDNA-CN in aging-related disease, this work begins to unravel the many varied  
625 underlying mechanisms for genetic control of mtDNA-CN.

626

### 627 **Supplemental Data**

628 Supplemental Data include six supplemental figures, fourteen supplemental tables, and  
629 supplemental methods.

630

### 631 **Declaration of Interests**

632 Psaty serves on the Steering Committee of the Yale Open Data Access Project funded by  
633 Johnson & Johnson. All other authors declare no competing interests.

634

### 635 **Acknowledgements**

636 This work was supported by National Heart, Lung and Blood Institute, National Institutes of  
637 Health (NIH) grants R01HL13573 and R01HL144569 (RJL, SYY, CAC, DEA), NIH grant P01-  
638 AG027734 (GA, YK, NB, AB) and the National Center for Advancing Translational Sciences,  
639 NIH, through Grant KL2TR002490 (LMR). The content is solely the responsibility of the authors  
640 and does not necessarily represent the official views of the NIH. LMR was also funded by T32  
641 HL129982.

642

643 This research was conducted using data from the Genotype-Tissue Expression (GTEx) project  
644 (dbGaP accession: phs000424.v8.p2). The GTEx project was supported by the Common  
645 Fund of the Office of the Director of the National Institutes of Health, and by NCI, NHGRI,  
646 NHLBI, NIDA, NIMH, and NINDS.

647

648 This research was also conducted using the UK Biobank Resource under Application Number  
649 17731.

650 The Atherosclerosis Risk in Communities study has been funded in whole or in part with Federal  
651 funds from the National Heart, Lung, and Blood Institute, National Institutes of Health,  
652 Department of Health and Human Services (contract numbers HSN268201700001I,  
653 HHSN268201700002I, HHSN268201700003I, HHSN268201700004I and  
654 HHSN268201700005I), R01HL087641, R01HL059367 and R01HL086694; National Human  
655 Genome Research Institute contract U01HG004402; and National Institutes of Health contract  
656 HHSN268200625226C. Funding support for “Building on GWAS for NHLBI diseases: the U.S.  
657 CHARGE consortium” was provided by the NIH through the American Recovery and  
658 Reinvestment Act of 2009 (ARRA) (5RC2HL102419). Sequencing was carried out at the Baylor  
659 College of Medicine Human Genome Sequencing Center and supported by the National Human  
660 Genome Research Institute grants U54 HG003273 and UM1 HG008898. The authors thank the  
661 staff and participants of the ARIC study for their important contributions. Infrastructure was  
662 partly supported by Grant Number UL1RR025005, a component of the National Institutes of  
663 Health and NIH Roadmap for Medical Research.

664 This CHS research was supported by NHLBI contracts HHSN268201200036C,  
665 HHSN268200800007C, HHSN268201800001C, N01HC55222, N01HC85079, N01HC85080,  
666 N01HC85081, N01HC85082, N01HC85083, N01HC85086, 75N92021D00006; and NHLBI  
667 grants U01HL080295, R01HL087652, R01HL105756, R01HL103612, R01HL120393, and  
668 U01HL130114 with additional contribution from the National Institute of Neurological Disorders  
669 and Stroke (NINDS). Additional support was provided through R01AG023629 from the National  
670 Institute on Aging (NIA). A full list of principal CHS investigators and institutions can be found  
671 at CHS-NHLBI.org. The provision of genotyping data was supported in part by the National  
672 Center for Advancing Translational Sciences, CTSI grant UL1TR001881, and the National

673 Institute of Diabetes and Digestive and Kidney Disease Diabetes Research Center (DRC) grant  
674 DK063491 to the Southern California Diabetes Endocrinology Research Center.

675 The FHS phenotype-genotype analyses were supported by the National Institute of Aging  
676 (U34AG051418). This research was conducted in part using data and resources from the  
677 Framingham Heart Study of the National Heart Lung and Blood Institute of the National  
678 Institutes of Health and Boston University School of Medicine. This work was partially supported  
679 by the National Heart, Lung and Blood Institute's Framingham Heart Study (Contract No. N01-  
680 HC-25195, HHSN268201500001) and its contract with Affymetrix, Inc for genotyping services  
681 (Contract No. N02-HL-6-4278). Genotyping, quality control and calling of the Illumina  
682 HumanExome BeadChip in the Framingham Heart Study was supported by funding from the  
683 National Heart, Lung and Blood Institute Division of Intramural Research (Daniel Levy and  
684 Christopher J. O'Donnell, Principle Investigators). The authors thank the participants for their  
685 dedication to the study. The authors are pleased to acknowledge that the computational work  
686 reported on in this paper was performed on the Shared Computing Cluster which is  
687 administered by Boston University's Research Computing Services. URL:  
688 [www.bu.edu/tech/support/research/](http://www.bu.edu/tech/support/research/).

689

690 MESA and the MESA SHARe projects are conducted and supported by the National Heart,  
691 Lung, and Blood Institute (NHLBI) in collaboration with MESA investigators. Support for MESA  
692 is provided by contracts 75N92020D00001, HHSN268201500003I, N01-HC-95159,  
693 75N92020D00005, N01-HC-95160, 75N92020D00002, N01-HC-95161, 75N92020D00003,  
694 N01-HC-95162, 75N92020D00006, N01-HC-95163, 75N92020D00004, N01-HC-95164,  
695 75N92020D00007, N01-HC-95165, N01-HC-95166, N01-HC-95167, N01-HC-95168, N01-HC-  
696 95169, UL1-TR-000040, UL1-TR-001079, and UL1-TR-001420. Funding for SHARe  
697 genotyping was provided by NHLBI Contract N02-HL-64278. Genotyping was performed at

698 Affymetrix (Santa Clara, California, USA) and the Broad Institute of Harvard and MIT (Boston,  
699 Massachusetts, USA) using the Affymetrix Genome-Wide Human SNP Array 6.0. The provision  
700 of genotyping data was supported in part by the National Center for Advancing Translational  
701 Sciences, CTSI grant UL1TR001881, and the National Institute of Diabetes and Digestive and  
702 Kidney Disease Diabetes Research Center (DRC) grant DK063491 to the Southern California  
703 Diabetes Endocrinology Research Center.

704

705 ROS/MAP is supported by the Translational Genomics Research Institute and National Institute  
706 on Aging (NIA) through grants U01AG46152, U01AG61256, P30AG10161, R01AG17917,  
707 RF1AG15819, R01AG30146.

708

709 SHIP and SHIP-TREND are part of the Community Medicine Research net of the University of  
710 Greifswald, Germany, which is funded by the Federal Ministry of Education and Research (grants  
711 no. 01ZZ9603, 01ZZ0103, and 01ZZ0403), the Ministry of Cultural Affairs as well as the Social  
712 Ministry of the Federal State of Mecklenburg-West Pomerania, and the network 'Greifswald  
713 Approach to Individualized Medicine (GANI\_MED)' funded by the Federal Ministry of Education  
714 and Research (grant 03IS2061A).

715 This research was funded in whole, or in part, by the Wellcome Trust [Grant ref: 217065/Z/19/Z].  
716 For the purpose of Open Access, the author has applied a CC BY public copyright license to any  
717 Author Accepted Manuscript version arising from this submission.

## 718 **Web Resources**

719 REVIGO was accessed at <http://revigo.irb.hr/>.

720

## 721 **Data and Code Availability**

722 All data used in this manuscript is available through either the UKBiobank and CHARGE  
723 consortiums. Code and scripts are available in a zipped file at  
724 <https://www.arkinglab.org/resources/>.

725

## 726 **References**

- 727 1. Wallace, D.C. (1992). Diseases of the Mitochondrial Dna. *Annual Review of Biochemistry* *61*,  
728 1175–1212.
- 729 2. Vakifahmetoglu-Norberg, H., Ouchida, A.T., and Norberg, E. (2017). The role of mitochondria  
730 in metabolism and cell death. *Biochem. Biophys. Res. Commun.* *482*, 426–431.
- 731 3. Herst, P.M., Rowe, M.R., Carson, G.M., and Berridge, M.V. (2017). Functional Mitochondria  
732 in Health and Disease. *Front Endocrinol (Lausanne)* *8*,
- 733 4. Dai, D.-F., Rabinovitch, P.S., and Ungvari, Z. (2012). Mitochondria and cardiovascular aging.  
734 *Circ. Res.* *110*, 1109–1124.
- 735 5. Cui, H., Kong, Y., and Zhang, H. (2012). Oxidative stress, mitochondrial dysfunction, and  
736 aging. *J Signal Transduct* *2012*, 646354.
- 737 6. Liu, C.-S., Tsai, C.-S., Kuo, C.-L., Chen, H.-W., Lii, C.-K., Ma, Y.-S., and Wei, Y.-H. (2003).  
738 Oxidative stress-related alteration of the copy number of mitochondrial DNA in human  
739 leukocytes. *Free Radic Res* *37*, 1307–1317.
- 740 7. Guha, M., and Avadhani, N.G. (2013). Mitochondrial retrograde signaling at the crossroads of  
741 tumor bioenergetics, genetics and epigenetics. *Mitochondrion* *13*, 577–591.
- 742 8. Malik, A.N., and Czajka, A. (2013). Is mitochondrial DNA content a potential biomarker of  
743 mitochondrial dysfunction? *Mitochondrion* *13*, 481–492.
- 744 9. Ashar, F.N., Moes, A., Moore, A.Z., Grove, M.L., Chaves, P.H.M., Coresh, J., Newman, A.B.,  
745 Matteini, A.M., Bandeen-Roche, K., Boerwinkle, E., et al. (2015). Association of mitochondrial  
746 DNA levels with frailty and all-cause mortality. *J. Mol. Med.* *93*, 177–186.
- 747 10. Ashar, F.N., Zhang, Y., Longchamps, R.J., Lane, J., Moes, A., Grove, M.L., Mychaleckyj,  
748 J.C., Taylor, K.D., Coresh, J., Rotter, J.I., et al. (2017). Association of Mitochondrial DNA Copy  
749 Number With Cardiovascular Disease. *JAMA Cardiol* *2*, 1247–1255.
- 750 11. Hong, Y.S., Longchamps, R.J., Zhao, D., Castellani, C.A., Loehr, L.R., Chang, P.P.,  
751 Matsushita, K., Grove, M.L., Boerwinkle, E., Arking, D.E., et al. (2020). Mitochondrial DNA Copy  
752 Number and Incident Heart Failure: The Atherosclerosis Risk in Communities (ARIC) Study.  
753 *Circulation* *141*, 1823–1825.
- 754 12. Zhao, D., Bartz, T.M., Sotoodehnia, N., Post, W.S., Heckbert, S.R., Alonso, A.,  
755 Longchamps, R.J., Castellani, C.A., Hong, Y.S., Rotter, J.I., et al. (2020). Mitochondrial DNA  
756 copy number and incident atrial fibrillation. *BMC Medicine* *18*, 246.

- 757 13. Tin, A., Grams, M.E., Ashar, F.N., Lane, J.A., Rosenberg, A.Z., Grove, M.L., Boerwinkle, E.,  
758 Selvin, E., Coresh, J., Pankratz, N., et al. (2016). Association between Mitochondrial DNA Copy  
759 Number in Peripheral Blood and Incident CKD in the Atherosclerosis Risk in Communities  
760 Study. *J Am Soc Nephrol* 27, 2467–2473.
- 761 14. Pyle, A., Anugraha, H., Kurzawa-Akanbi, M., Yarnall, A., Burn, D., and Hudson, G. (2016).  
762 Reduced mitochondrial DNA copy number is a biomarker of Parkinson’s disease. *Neurobiol.*  
763 *Aging* 38, 216.e7-216.e10.
- 764 15. Wei, W., Keogh, M.J., Wilson, I., Coxhead, J., Ryan, S., Rollinson, S., Griffin, H., Kurzawa-  
765 Akanbi, M., Santibanez-Koref, M., Talbot, K., et al. (2017). Mitochondrial DNA point mutations  
766 and relative copy number in 1363 disease and control human brains. *Acta Neuropathologica*  
767 *Communications* 5, 13.
- 768 16. Reznik, E., Miller, M.L., Şenbabaoğlu, Y., Riaz, N., Sarungbam, J., Tickoo, S.K., Al-  
769 Ahmadi, H.A., Lee, W., Seshan, V.E., Hakimi, A.A., et al. (2016). Mitochondrial DNA copy  
770 number variation across human cancers. *Elife* 5.
- 771 17. Hurtado-Roca, Y., Ledesma, M., Gonzalez-Lazaro, M., Moreno-Loshuertos, R., Fernandez-  
772 Silva, P., Enriquez, J.A., and Laclaustra, M. (2016). Adjusting MtDNA Quantification in Whole  
773 Blood for Peripheral Blood Platelet and Leukocyte Counts. *PLOS ONE* 11, e0163770.
- 774 18. Knez, J., Winckelmans, E., Plusquin, M., Thijs, L., Cauwenberghs, N., Gu, Y., Staessen,  
775 J.A., Nawrot, T.S., and Kuznetsova, T. (2016). Correlates of Peripheral Blood Mitochondrial  
776 DNA Content in a General Population. *Am J Epidemiol* 183, 138–146.
- 777 19. Kumar, P., Efstathopoulos, P., Millischer, V., Olsson, E., Wei, Y.B., Brüstle, O., Schalling,  
778 M., Villaescusa, J.C., Ösby, U., and Lavebratt, C. (2018). Mitochondrial DNA copy number is  
779 associated with psychosis severity and anti-psychotic treatment. *Sci Rep* 8, 12743.
- 780 20. Urata, M., Koga-Wada, Y., Kayamori, Y., and Kang, D. (2008). Platelet contamination  
781 causes large variation as well as overestimation of mitochondrial DNA content of peripheral  
782 blood mononuclear cells. *Ann Clin Biochem* 45, 513–514.
- 783 21. Clay Montier, L.L., Deng, J.J., and Bai, Y. (2009). Number matters: control of mammalian  
784 mitochondrial DNA copy number. *J Genet Genomics* 36, 125–131.
- 785 22. Tang, Y., Schon, E.A., Wilichowski, E., Vazquez-Memije, M.E., Davidson, E., and King, M.P.  
786 (2000). Rearrangements of Human Mitochondrial DNA (mtDNA): New Insights into the  
787 Regulation of mtDNA Copy Number and Gene Expression. *Mol Biol Cell* 11, 1471–1485.
- 788 23. Xing, J., Chen, M., Wood, C.G., Lin, J., Spitz, M.R., Ma, J., Amos, C.I., Shields, P.G.,  
789 Benowitz, N.L., Gu, J., et al. (2008). Mitochondrial DNA content: its genetic heritability and  
790 association with renal cell carcinoma. *J. Natl. Cancer Inst.* 100, 1104–1112.
- 791 24. Carling, P.J., Cree, L.M., and Chinnery, P.F. (2011). The implications of mitochondrial DNA  
792 copy number regulation during embryogenesis. *Mitochondrion* 11, 686–692.
- 793 25. Harvey, A., Gibson, T., Lonergan, T., and Brenner, C. (2011). Dynamic regulation of  
794 mitochondrial function in preimplantation embryos and embryonic stem cells. *Mitochondrion* 11,  
795 829–838.

- 796 26. Copeland, W.C. (2014). Defects of Mitochondrial DNA Replication. *J Child Neurol* 29, 1216–  
797 1224.
- 798 27. Mandel, H., Szargel, R., Labay, V., Elpeleg, O., Saada, A., Shalata, A., Anbinder, Y.,  
799 Berkowitz, D., Hartman, C., Barak, M., et al. (2001). The deoxyguanosine kinase gene is  
800 mutated in individuals with depleted hepatocerebral mitochondrial DNA. *Nat. Genet.* 29, 337–  
801 341.
- 802 28. Wang, L., Limongelli, A., Vila, M.R., Carrara, F., Zeviani, M., and Eriksson, S. (2005).  
803 Molecular insight into mitochondrial DNA depletion syndrome in two patients with novel  
804 mutations in the deoxyguanosine kinase and thymidine kinase 2 genes. *Mol. Genet. Metab.* 84,  
805 75–82.
- 806 29. Rusecka, J., Kaliszewska, M., Bartnik, E., and Tońska, K. (2018). Nuclear genes involved in  
807 mitochondrial diseases caused by instability of mitochondrial DNA. *J Appl Genet* 59, 43–57.
- 808 30. Cai, N., Li, Y., Chang, S., Liang, J., Lin, C., Zhang, X., Liang, L., Hu, J., Chan, W., Kendler,  
809 K.S., et al. (2015). Genetic Control over mtDNA and Its Relationship to Major Depressive  
810 Disorder. *Curr Biol* 25, 3170–3177.
- 811 31. Workalemahu, T., Enquobahrie, D.A., Tadesse, M.G., Hevner, K., Gelaye, B., Sanchez, S.,  
812 and Williams, M.A. (2017). Genetic Variations Related to Maternal Whole Blood Mitochondrial  
813 DNA Copy Number: A Genome-Wide and Candidate Gene Study. *J Matern Fetal Neonatal Med*  
814 30, 2433–2439.
- 815 32. Guyatt, A.L., Brennan, R.R., Burrows, K., Guthrie, P.A.I., Ascione, R., Ring, S.M., Gaunt,  
816 T.R., Pyle, A., Cordell, H.J., Lawlor, D.A., et al. (2019). A genome-wide association study of  
817 mitochondrial DNA copy number in two population-based cohorts. *Human Genomics* 13, 6.
- 818 33. Hägg, S., Jylhävä, J., Wang, Y., Czene, K., and Grassmann, F. (2020). Deciphering the  
819 genetic and epidemiological landscape of mitochondrial DNA abundance. *Hum Genet.*
- 820 34. Longchamps, R.J., Castellani, C.A., Yang, S.Y., Newcomb, C.E., Sumpter, J.A., Lane, J.,  
821 Grove, M.L., Guallar, E., Pankratz, N., Taylor, K.D., et al. (2020). Evaluation of mitochondrial  
822 DNA copy number estimation techniques. *PLOS ONE* 15, e0228166.
- 823 35. MitoPipeline: Generating Mitochondrial copy number estimates from SNP array data in  
824 Genvisis.
- 825 36. Han, B., and Eskin, E. (2011). Random-Effects Model Aimed at Discovering Associations in  
826 Meta-Analysis of Genome-wide Association Studies. *The American Journal of Human Genetics*  
827 88, 586–598.
- 828 37. Wang, G., Sarkar, A., Carbonetto, P., and Stephens, M. (2020). A simple new approach to  
829 variable selection in regression, with application to genetic fine mapping. *Journal of the Royal*  
830 *Statistical Society: Series B (Statistical Methodology)* 82, 1273–1300.
- 831 38. Wang, K., Li, M., and Hakonarson, H. (2010). ANNOVAR: functional annotation of genetic  
832 variants from high-throughput sequencing data. *Nucleic Acids Research* 38, e164–e164.

- 833 39. Sherry, S.T., Ward, M., and Sirotkin, K. (1999). dbSNP-database for single nucleotide  
834 polymorphisms and other classes of minor genetic variation. *Genome Res* 9, 677–679.
- 835 40. O’Leary, N.A., Wright, M.W., Brister, J.R., Ciufu, S., Haddad, D., McVeigh, R., Rajput, B.,  
836 Robbertse, B., Smith-White, B., Ako-Adjei, D., et al. (2016). Reference sequence (RefSeq)  
837 database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic*  
838 *Acids Res* 44, D733-745.
- 839 41. Giambartolomei, C., Vukcevic, D., Schadt, E.E., Franke, L., Hingorani, A.D., Wallace, C.,  
840 and Plagnol, V. (2014). Bayesian Test for Colocalisation between Pairs of Genetic Association  
841 Studies Using Summary Statistics. *PLOS Genetics* 10, e1004383.
- 842 42. Vösa, U., Claringbould, A., Westra, H.-J., Bonder, M.J., Deelen, P., Zeng, B., Kirsten, H.,  
843 Saha, A., Kreuzhuber, R., Kasela, S., et al. (2018). Unraveling the polygenic architecture of  
844 complex traits using blood eQTL metaanalysis. *BioRxiv* 447367.
- 845 43. Pers, T.H., Karjalainen, J.M., Chan, Y., Westra, H.-J., Wood, A.R., Yang, J., Lui, J.C.,  
846 Vedantam, S., Gustafsson, S., Esko, T., et al. (2015). Biological interpretation of genome-wide  
847 association studies using predicted gene functions. *Nature Communications* 6, 5890.
- 848 44. Smyth, G., Hu, Y., Ritchie, M., Silver, J., Wettenhall, J., McCarthy, D., Wu, D., Shi, W.,  
849 Phipson, B., Lun, A., et al. (2021). limma: Linear Models for Microarray Data (Bioconductor  
850 version: Release (3.12)).
- 851 45. Young, M.D., Wakefield, M.J., Smyth, G.K., and Oshlack, A. (2010). Gene ontology analysis  
852 for RNA-seq: accounting for selection bias. *Genome Biology* 11, R14.
- 853 46. Barbeira, A.N., Dickinson, S.P., Bonazzola, R., Zheng, J., Wheeler, H.E., Torres, J.M.,  
854 Torstenson, E.S., Shah, K.P., Garcia, T., Edwards, T.L., et al. (2018). Exploring the phenotypic  
855 consequences of tissue specific gene expression variation inferred from GWAS summary  
856 statistics. *Nat Commun* 9, 1825.
- 857 47. Barbeira, A.N., Pividori, M., Zheng, J., Wheeler, H.E., Nicolae, D.L., and Im, H.K. (2019).  
858 Integrating predicted transcriptome from multiple tissues improves association detection. *PLOS*  
859 *Genetics* 15, e1007889.
- 860 48. Supek, F., Bošnjak, M., Škunca, N., and Šmuc, T. (2011). REVIGO summarizes and  
861 visualizes long lists of gene ontology terms. *PLoS One* 6, e21800.
- 862 49. Stiles, A.R., Simon, M.T., Stover, A., Eftekharian, S., Khanlou, N., Wang, H.L., Magaki, S.,  
863 Lee, H., Partynski, K., Dorrani, N., et al. (2016). Mutations in TFAM, encoding mitochondrial  
864 transcription factor A, cause neonatal liver failure associated with mtDNA depletion. *Mol Genet*  
865 *Metab* 119, 91–99.
- 866 50. Kornblum, C., Nicholls, T.J., Haack, T.B., Schöler, S., Peeva, V., Danhauser, K., Hallmann,  
867 K., Zsurka, G., Rorbach, J., Iuso, A., et al. (2013). Loss-of-function mutations in MGME1 impair  
868 mtDNA replication and cause multisystemic mitochondrial disease. *Nat Genet* 45, 214–219.
- 869 51. El-Hattab, A.W., and Scaglia, F. (2013). Mitochondrial DNA depletion syndromes: review  
870 and updates of genetic basis, manifestations, and therapeutic options. *Neurotherapeutics* 10,  
871 186–198.



- 872 52. Millard, L.A.C., Davies, N.M., Gaunt, T.R., Davey Smith, G., and Tilling, K. (2018). Software  
873 Application Profile: PHEASANT: a tool for performing automated phenome scans in UK Biobank.  
874 *Int J Epidemiol* *47*, 29–35.
- 875 53. Konopka, T. (2020). umap: Uniform Manifold Approximation and Projection.
- 876 54. Hahsler, M., Piekenbrock, M., Arya, S., and Mount, D. (2019). dbscan: Density Based  
877 Clustering of Applications with Noise (DBSCAN) and Related Algorithms.
- 878 55. Gibellini, L., De Gaetano, A., Mandrioli, M., Van Tongeren, E., Bortolotti, C.A., Cossarizza,  
879 A., and Pinti, M. (2020). The biology of Lonp1: More than a mitochondrial protease. *Int Rev Cell*  
880 *Mol Biol* *354*, 1–61.
- 881 56. GTEx Consortium (2013). The Genotype-Tissue Expression (GTEx) project. *Nat Genet* *45*,  
882 580–585.
- 883 57. Liu, T., Lu, B., Lee, I., Ondrovicová, G., Kutejová, E., and Suzuki, C.K. (2004). DNA and  
884 RNA binding by the mitochondrial lon protease is regulated by nucleotide and protein substrate.  
885 *J Biol Chem* *279*, 13902–13910.
- 886 58. Sharma, M.R., Koc, E.C., Datta, P.P., Booth, T.M., Spremulli, L.L., and Agrawal, R.K.  
887 (2003). Structure of the mammalian mitochondrial ribosome reveals an expanded functional role  
888 for its component proteins. *Cell* *115*, 97–108.
- 889 59. Shimizu, S., Narita, M., and Tsujimoto, Y. (1999). Bcl-2 family proteins regulate the release  
890 of apoptogenic cytochrome c by the mitochondrial channel VDAC. *Nature* *399*, 483–487.
- 891 60. Vagdatli, E., Gounari, E., Lazaridou, E., Katsibourlia, E., Tsikopoulou, F., and Labrianou, I.  
892 (2010). Platelet distribution width: a simple, practical and specific marker of activation of  
893 coagulation. *Hippokratia* *14*, 28–32.
- 894 61. PathCards :: Factors involved in megakaryocyte development and platelet production  
895 Pathway and related pathways.
- 896 62. PathCards :: Apoptosis and Autophagy Pathway and related pathways.
- 897 63. Basel, D. (2020). Mitochondrial DNA Depletion Syndromes. *Clin Perinatol* *47*, 123–141.
- 898 64. DeBarmore, B., Longchamps, R.J., Zhang, Y., Kalyani, R.R., Guallar, E., Arking, D.E.,  
899 Selvin, E., and Young, J.H. (2020). Mitochondrial DNA copy number and diabetes: the  
900 Atherosclerosis Risk in Communities (ARIC) study. *BMJ Open Diabetes Res Care* *8*,
- 901 65. Yang, S.Y., Castellani, C.A., Longchamps, R.J., Pillalamarri, V.K., O'Rourke, B., Guallar, E.,  
902 and Arking, D.E. (2020). Blood-derived mitochondrial DNA copy number is associated with gene  
903 expression across multiple tissues and is predictive for incident neurodegenerative disease.  
904 *BioRxiv* 2020.07.17.209023.
- 905 66. Hannah-Shmouni, F., MacNeil, L., Brady, L., Nilsson, M.I., and Tarnopolsky, M. (2019).  
906 Expanding the Clinical Spectrum of LONP1-Related Mitochondrial Cytopathy. *Front Neurol* *10*,  
907 981.

- 908 67. Grainha, T.R.R., Jorge, P.A. da S., Pérez-Pérez, M., Pérez Rodríguez, G., Pereira,  
909 M.O.B.O., and Lourenço, A.M.G. (2018). Exploring anti-quorum sensing and anti-virulence  
910 based strategies to fight *Candida albicans* infections: an in silico approach. *FEMS Yeast Res*  
911 *18*,.
- 912 68. Ali, A.T., Idaghdour, Y., and Hodgkinson, A. (2020). Analysis of mitochondrial m1A/G RNA  
913 modification reveals links to nuclear genetic variants and associated disease processes.  
914 *Communications Biology* *3*, 1–11.
- 915 69. Schrepfer, E., and Scorrano, L. (2016). Mitofusins, from Mitochondria to Metabolism. *Mol*  
916 *Cell* *61*, 683–694.
- 917 70. Ishihara, N., Eura, Y., and Mihara, K. (2004). Mitofusin 1 and 2 play distinct roles in  
918 mitochondrial fusion reactions via GTPase activity. *J Cell Sci* *117*, 6535–6546.
- 919 71. Zhang, D., Xi, Y., Coccimiglio, M.L., Mennigen, J.A., Jonz, M.G., Ekker, M., and Trudeau,  
920 V.L. (2012). Functional prediction and physiological characterization of a novel short trans-  
921 membrane protein 1 as a subunit of mitochondrial respiratory complexes. *Physiol Genomics* *44*,  
922 1133–1140.
- 923 72. Cavdar Koc, E., Burkhart, W., Blackburn, K., Moseley, A., and Spremulli, L.L. (2001). The  
924 small subunit of the mammalian mitochondrial ribosome. Identification of the full complement of  
925 ribosomal proteins present. *J Biol Chem* *276*, 19363–19374.
- 926 73. Márquez-Jurado, S., Díaz-Colunga, J., das Neves, R.P., Martínez-Lorente, A., Almazán, F.,  
927 Guantes, R., and Iborra, F.J. (2018). Mitochondrial levels determine variability in cell death by  
928 modulating apoptotic gene expression. *Nat Commun* *9*, 389.
- 929 74. Dölle, C., Flønes, I., Nido, G.S., Miletic, H., Osuagwu, N., Kristoffersen, S., Lilleng, P.K.,  
930 Larsen, J.P., Tysnes, O.-B., Haugarvoll, K., et al. (2016). Defective mitochondrial DNA  
931 homeostasis in the substantia nigra in Parkinson disease. *Nat Commun* *7*, 13548.
- 932 75. Chen, C., Turnbull, D.M., and Reeve, A.K. (2019). Mitochondrial Dysfunction in Parkinson's  
933 Disease-Cause or Consequence? *Biology (Basel)* *8*,.
- 934 76. Pinto, M., and Moraes, C.T. (2014). Mitochondrial genome changes and neurodegenerative  
935 diseases. *Biochim Biophys Acta* *1842*, 1198–1207.
- 936 77. Rodríguez-Santiago, B., Casademont, J., and Nunes, V. (2001). Is mitochondrial DNA  
937 depletion involved in Alzheimer's disease? *Eur J Hum Genet* *9*, 279–285.
- 938 78. Chew, H., Solomon, V.A., and Fonteh, A.N. (2020). Involvement of Lipids in Alzheimer's  
939 Disease Pathology and Potential Therapies. *Front Physiol* *11*, 598.
- 940 79. Yin, J., Reiman, E.M., Beach, T.G., Serrano, G.E., Sabbagh, M.N., Nielsen, M., Caselli,  
941 R.J., and Shi, J. (2020). Effect of ApoE isoforms on mitochondria in Alzheimer disease.  
942 *Neurology* *94*, e2404–e2411.
- 943 80. El Khoury, J.B., Moore, K.J., Means, T.K., Leung, J., Terada, K., Toft, M., Freeman, M.W.,  
944 and Luster, A.D. (2003). CD36 mediates the innate host response to beta-amyloid. *J Exp Med*  
945 *197*, 1657–1666.

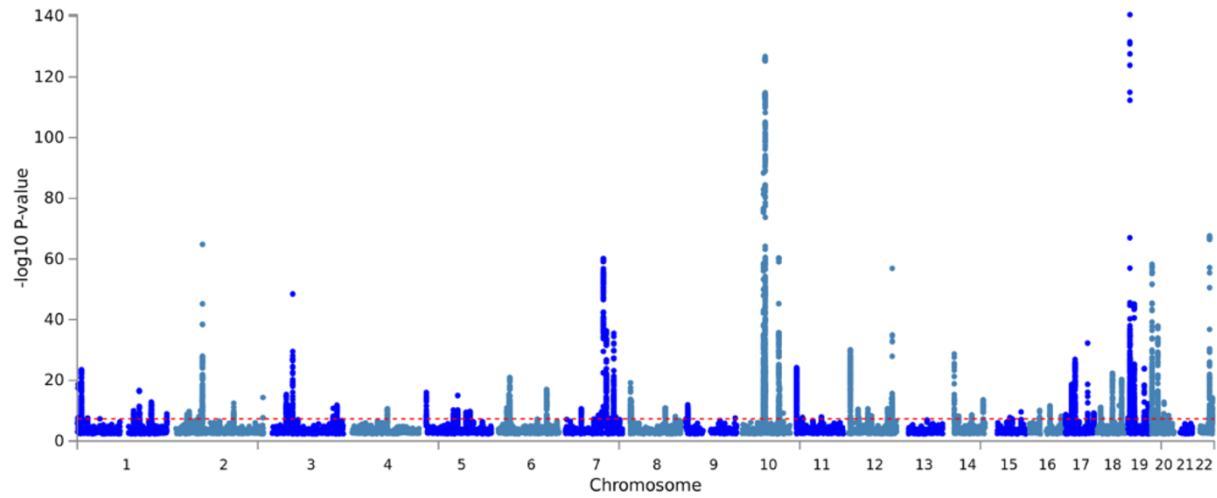
- 946 81. Lämsä, R., Helisalmi, S., Herukka, S.-K., Tapiola, T., Pirttilä, T., Vepsäläinen, S., Hiltunen,  
947 M., and Soininen, H. (2008). Genetic study evaluating LDLR polymorphisms and Alzheimer's  
948 disease. *Neurobiol Aging* 29, 848–855.
- 949 82. Tomioka, M., Toda, Y., Mañucat, N.B., Akatsu, H., Fukumoto, M., Kono, N., Arai, H., Kioka,  
950 N., and Ueda, K. (2017). Lysophosphatidylcholine export by human ABCA7. *Biochim Biophys*  
951 *Acta Mol Cell Biol Lipids* 1862, 658–665.
- 952 83. Nuytemans, K., Maldonado, L., Ali, A., John-Williams, K., Beecham, G.W., Martin, E., Scott,  
953 W.K., and Vance, J.M. (2016). Overlap between Parkinson disease and Alzheimer disease in  
954 ABCA7 functional variants. *Neurol Genet* 2, e44.
- 955 84. Yang, C., Pring, M., Wear, M.A., Huang, M., Cooper, J.A., Svitkina, T.M., and Zigmond, S.H.  
956 (2005). Mammalian CARMIL inhibits actin filament capping by capping protein. *Dev Cell* 9, 209–  
957 221.
- 958 85. Crabos, M., Yamakado, T., Heizmann, C.W., Cerletti, N., Bühler, F.R., and Erne, P. (1991).  
959 The calcium binding protein tropomyosin in human platelets and cardiac tissue: elevation in  
960 hypertensive cardiac hypertrophy. *Eur J Clin Invest* 21, 472–478.
- 961 86. Kostan, J., Salzer, U., Orlova, A., Törö, I., Hodnik, V., Senju, Y., Zou, J., Schreiner, C.,  
962 Steiner, J., Meriläinen, J., et al. (2014). Direct interaction of actin filaments with F-BAR protein  
963 pacsin2. *EMBO Rep* 15, 1154–1162.
- 964 87. Sever, S. (2002). Dynamin and endocytosis. *Curr Opin Cell Biol* 14, 463–467.
- 965 88. Cai, B., Giridharan, S.S.P., Zhang, J., Saxena, S., Bahl, K., Schmidt, J.A., Sorgen, P.L.,  
966 Guo, W., Naslavsky, N., and Caplan, S. (2013). Differential roles of C-terminal Eps15 homology  
967 domain proteins as vesiculators and tubulators of recycling endosomes. *J Biol Chem* 288,  
968 30172–30180.
- 969 89. Di Buduo, C.A., Moccia, F., Battiston, M., De Marco, L., Mazzucato, M., Moratti, R., Tanzi,  
970 F., and Balduini, A. (2014). The importance of calcium in the regulation of megakaryocyte  
971 function. *Haematologica* 99, 769–778.
- 972 90. De Botton, S., Sabri, S., Daugas, E., Zermati, Y., Guidotti, J.E., Hermine, O., Kroemer, G.,  
973 Vainchenker, W., and Debili, N. (2002). Platelet formation is the consequence of caspase  
974 activation within megakaryocytes. *Blood* 100, 1310–1317.
- 975 91. Josefsson, E.C., James, C., Henley, K.J., Debrincat, M.A., Rogers, K.L., Dowling, M.R.,  
976 White, M.J., Kruse, E.A., Lane, R.M., Ellis, S., et al. (2011). Megakaryocytes possess a  
977 functional intrinsic apoptosis pathway that must be restrained to survive and produce platelets. *J*  
978 *Exp Med* 208, 2017–2031.
- 979 92. Josefsson, E.C., Burnett, D.L., Lebois, M., Debrincat, M.A., White, M.J., Henley, K.J., Lane,  
980 R.M., Moujalled, D., Preston, S.P., O'Reilly, L.A., et al. (2014). Platelet production proceeds  
981 independently of the intrinsic and extrinsic apoptosis pathways. *Nat Commun* 5, 3455.
- 982 93. McArthur, K., Chappaz, S., and Kile, B.T. (2018). Apoptosis in megakaryocytes and  
983 platelets: the life and death of a lineage. *Blood* 131, 605–610.

- 984 94. Olkowicz, M., Tomczyk, M., Debski, J., Tyrankiewicz, U., Przyborowski, K., Borkowski, T.,  
985 Zabielska-Kaczorowska, M., Szupryczynska, N., Kochan, Z., Smeda, M., et al. (2021).  
986 Enhanced cardiac hypoxic injury in atherogenic dyslipidaemia results from alterations in the  
987 energy metabolism pattern. *Metabolism* 114, 154400.
- 988 95. Vyas, C.M., Ogata, S., Reynolds, C.F., Mischoulon, D., Chang, G., Cook, N.R., Manson,  
989 J.E., Crous-Bou, M., De Vivo, I., and Okereke, O.I. (2020). Lifestyle and behavioral factors and  
990 mitochondrial DNA copy number in a diverse cohort of mid-life and older adults. *PLoS One* 15,  
991 e0237235.
- 992

993 **Figures**

994 **Figure 1. Manhattan plot of GWS loci from UKB-only analyses.**

995

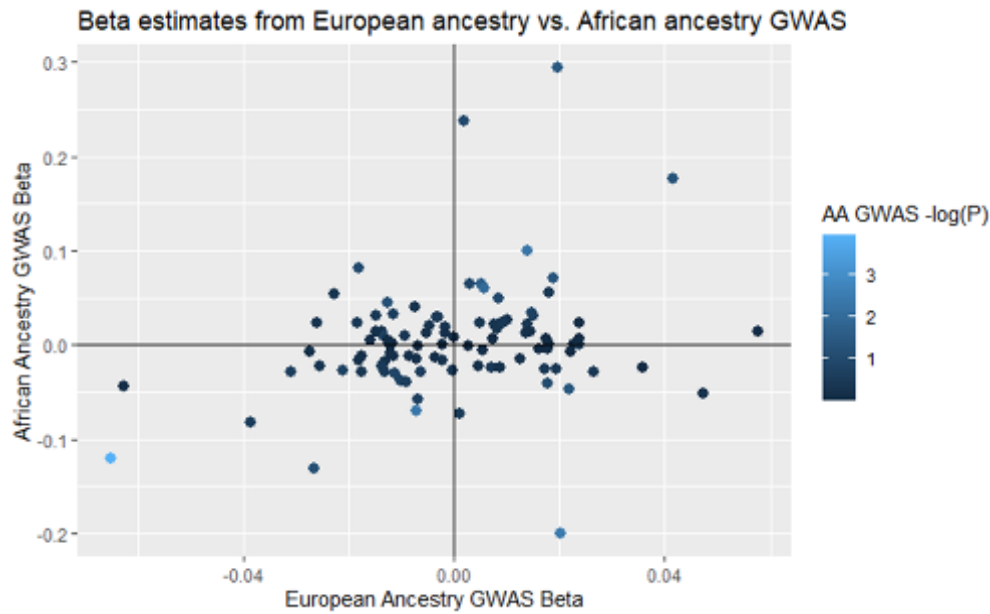


996

997 Manhattan plot showing genome-wide significant loci for the UK Biobank-only analyses.

998

999 **Figure 2. Scatterplot displaying effect size estimates between Whites and Blacks GWAS**  
1000 **results for the 129 loci identified in the Whites analyses.**



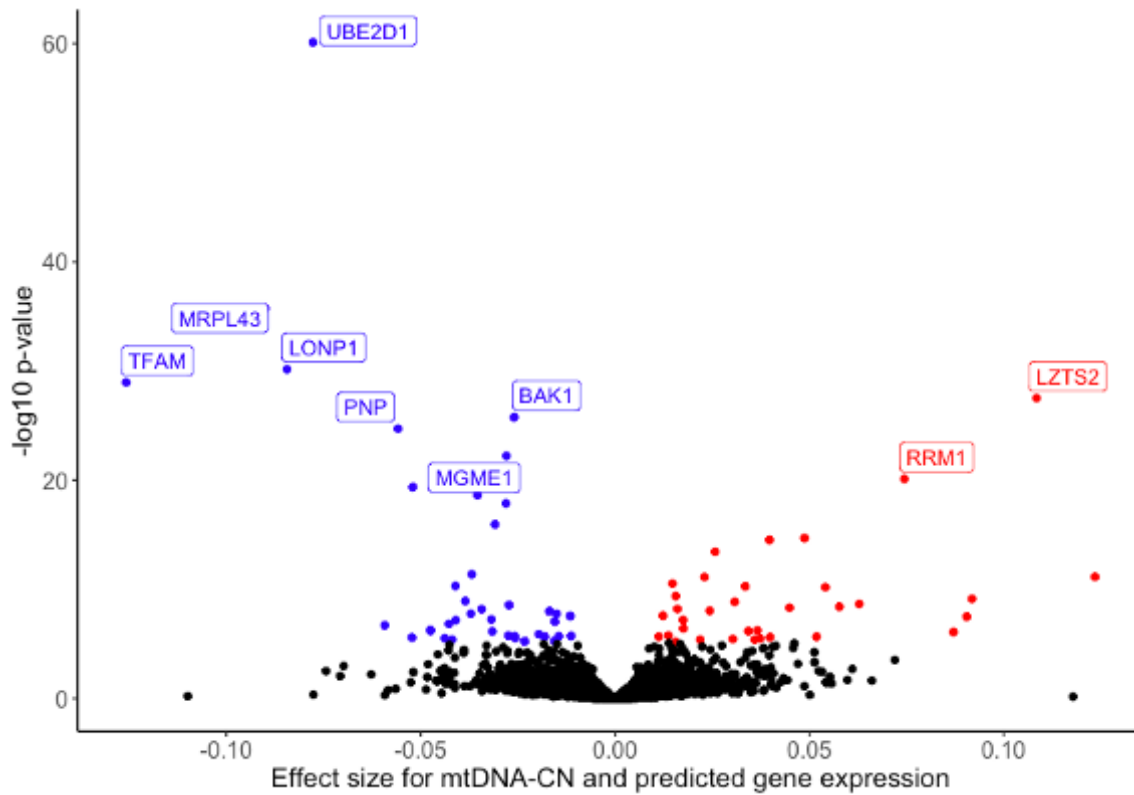
1001

1002 Scatterplot showing comparison between effect size estimates for White (European ancestry)  
1003 and Black (African Ancestry) individuals. Color represents significance of effect for each locus  
1004 in Blacks GWAS analyses.

1005

1006

1007 **Figure 3. Volcano plot of genes whose predicted gene expression is significantly**  
1008 **associated with mtDNA-CN.**



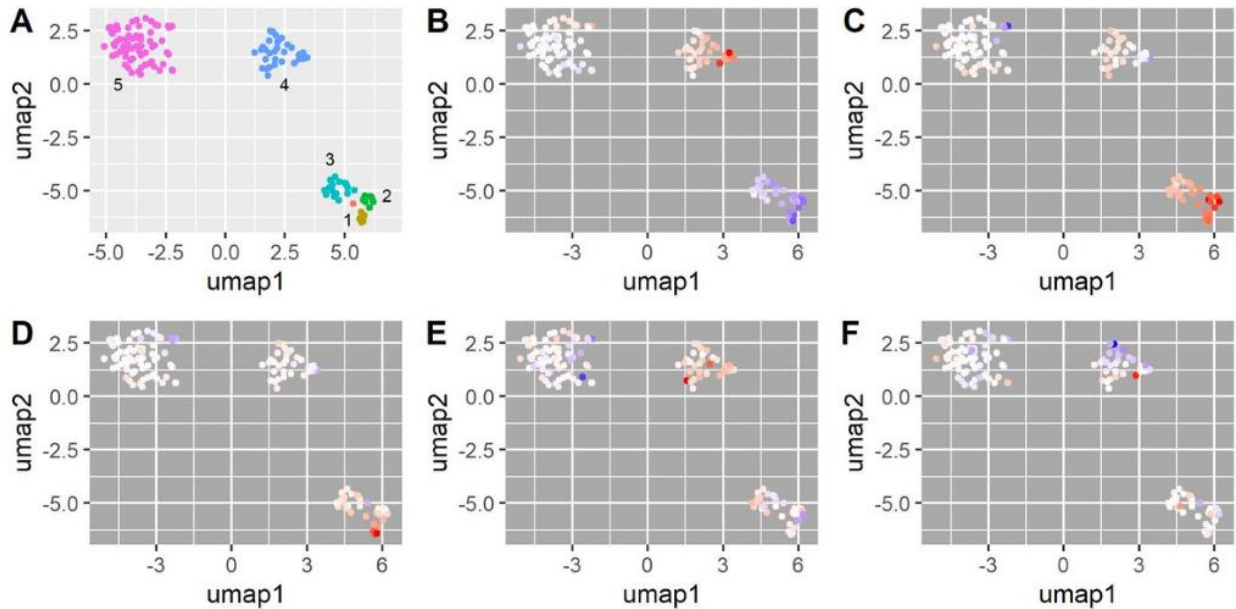
1009

1010 Volcano plot showing genes whose predicted gene expression is significantly associated with  
1011 mtDNA-CN. Red indicates positive associations, blue indicates negative associations. Three  
1012 genes (ARRDC1, EHMT1, PNPLA7) had extreme effect size estimates greater than 0.3 but  
1013 were non-significant and removed from the plot for readability.

1014

1015  
1016

**Figure 4. PHEWAS-based clustering of mtDNA-CN associated SNPs.**



1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025

**UMAP clusters created from PHEWAS associations for mtDNA-CN associated SNPs. (A) Five clusters were identified as labeled in the panel; orange indicates no cluster. (B-F) SNPs are colored based on their effect estimate size, standardized to the effect on mtDNA-CN (red = positive, blue = negative estimates), for (B) platelet count, (C) mean platelet volume, (D) platelet distribution width, (E) serum calcium levels, (F) mean spherical cellular volume.**