

1 **Mitochondrial DNA variants modulate N-formylmethionine, proteostasis and risk of**
2 **late-onset human diseases**

3
4 **Author list:**

5
6 Na Cai^{1,2,§,†}, Aurora Gomez-Duran^{3,4,£,†}, Ekaterina Yonova-Doing^{5,^}, Kousik Kundu¹,
7 Annette I. Burgess^{6,§}, Zoe J. Golder^{3,4,§}, Claudia Calabrese^{3,4}, Marc J. Bonder^{2,#}, Marta
8 Camacho³, Rachael A. Lawson⁷, Lixin Li⁶, Caroline H Williams-Gray³, ICICLE-PD Study
9 Group[‡], Emanuele Di Angelantonio^{5,6,9,10}, David J. Roberts^{9,11,12}, Nick A. Watkins¹³, Willem
10 H. Ouwehand^{1,6,13,14}, Adam S. Butterworth^{5,6,9,10}, Isobel D. Stewart¹⁵, Maik Pietzner¹⁵, Nick J.
11 Wareham¹⁵, Claudia Langenberg¹⁵, John Danesh^{1,5,6,9,10}, Klaudia Walter¹, Peter M. Rothwell⁶,
12 Joanna M. M. Howson^{5,^,*}, Oliver Stegle^{2,16,#,*}, Patrick F. Chinnery^{3,4,*}, Nicole Soranzo^{1,14,*}

13
14 **Affiliations:**

- 15
1. Human Genetics Department, Wellcome Sanger Institute (WT), Hinxton, UK
 2. European Bioinformatics Institute (EMBL-EBI), Hinxton, UK
 3. Department of Clinical Neurosciences, School of Clinical Medicine, University of Cambridge, Cambridge Biomedical Campus, Cambridge, UK
 4. Medical Research Council Mitochondrial Biology Unit, University of Cambridge, Cambridge Biomedical Campus, Cambridge, UK
 5. British Heart Foundation Cardiovascular Epidemiology Unit, Department of Primary Public Health and Primary care, University of Cambridge, Cambridge, UK
 6. Stroke Prevention Research Unit, Nuffield Department of Clinical Neurosciences, John Radcliffe Hospital, Oxford, UK
 7. Translational and Clinical Research Institute, Newcastle University, Newcastle, UK
 8. British Heart Foundation Centre of Research Excellence, University of Cambridge, Cambridge, UK
 9. National Institute for Health Research Blood and Transplant Research Unit in Donor Health and Genomics, University of Cambridge, Cambridge, UK
 10. Health Data Research UK Cambridge, Wellcome Genome Campus and University of Cambridge, Cambridge, UK
 11. NHS Blood and Transplant-Oxford Centre, Level 2, John Radcliffe Hospital, Oxford, UK
 12. Radcliffe Department of Medicine, University of Oxford, Oxford, UK
 13. NHS Blood and Transplant, Cambridge Biomedical Campus, Long Road, Cambridge, UK
 14. Department of Haematology, University of Cambridge, Cambridge, UK
 15. MRC Epidemiology Unit, University of Cambridge, Cambridge, UK
 16. European Molecular Biology Laboratory, Meyerhofstraße 1, Heidelberg, Germany

16
17 § Current address: Helmholtz Pioneer Campus, Helmholtz Zentrum München, Ingolstädter
18 Landstraße 1, Neuherberg, Germany

19 £ Current address: Centro de Investigaciones Biológicas “Margarita Salas”, Consejo Superior
20 de Investigaciones Científicas (CIB-CSIC), Madrid, Spain.

21 ^ Current address: Novo Nordisk Research Centre Oxford, Innovation Building, Old Road
22 Campus, Oxford, UK

23 # Current address: Deutsches Krebsforschungszentrum, Im Neuenheimer Feld 280,
24 Heidelberg, Germany

25

26 † These authors contributed equally

27 * Jointly supervised the work

28 ‡ A list of authors and their affiliations appears at the end of the paper

29

30 Correspondence should be addressed to O.S. (o.stegle@dkfz-heidelberg.de), P.F.C.

31 (pfc25@cam.ac.uk) and N. S. (ns6@sanger.ac.uk)

32

33 **Abstract:**

34

35 **Mitochondrial DNA (mtDNA) variants influence the risk of late-onset human diseases,**
36 **but the reasons are poorly understood. Undertaking an hypothesis-free analysis of 5,689**
37 **blood-derived biomarkers with mtDNA variants in 16,220 healthy donors, here we show**
38 **that variants defining mtDNA haplogroups Uk and H4 modulate the level of circulating**
39 **N-formylmethionine (fMet), which initiates mitochondrial protein translation. In**
40 **human cybrid lines, fMet modulated both mitochondrial and cytosolic proteins on**
41 **multiple levels - through transcription, post-translational modification, and proteolysis**
42 **by an N-degron pathway - abolishing known differences between mtDNA haplogroups.**
43 **In a further 11,966 individuals, fMet levels contributed to all-cause mortality and the**
44 **disease risk of several common cardiovascular disorders. Together these findings**
45 **indicate that fMet plays a key role in common age-related disease through pleiotropic**
46 **effects on cell proteostasis.**

47

48 **Introduction**

49

50 The 16.5kb human mitochondrial genome (mtDNA) encodes 13 proteins of the electron
51 transport chain (ETC) and the tRNA and rRNA machinery necessary for their transcription
52 and translation *in situ*¹. Genetic diversity in the maternally inherited² mtDNA with high
53 mutation rates³ has been effectively used for studies of human evolution and phylogenies⁴⁻⁷.
54 Because of the ubiquitous and essential roles of mtDNA encoded proteins in cellular
55 metabolism⁸, mtDNA sequence variations have been examined extensively for their effects
56 on cellular metabolism⁹⁻¹² and human health and diseases.

57

58 A compendium of rare mutations in mtDNA genes encoding ETC subunits have been
59 identified to cause severe multisystemic diseases^{13,14}, commonly due to the primary
60 biochemical consequences of the mutations on oxidative phosphorylation (OXPHOS) and the
61 synthesis of adenosine triphosphate (ATP)¹⁵. mtDNA variations and somatic mutations have
62 also been shown to be important in the pathology of cancers¹⁶⁻¹⁸ and inducing the Warburg
63 effect¹⁹. Common mtDNA variants with less severe phenotypes have been shown to affect
64 the risk of complex late-onset human diseases²⁰, including neurodegenerative diseases like
65 Alzheimer's disease²¹ and Parkinson's disease²², cardiovascular diseases like ischemic
66 stroke²³, myocardial infarction²⁴ and coronary artery disease²⁵, and metabolic diseases like
67 type 2 diabetes^{26,27}. In some cases, mtDNA haplogroups have pleiotropic effects on multiple
68 diseases, while in others they have opposite effects between diseases^{28,29}.

69

70 The relevance of mtDNA variations to health and disease is apparent³⁰, but
71 understanding the molecular underpinnings of genetic associations on the mtDNA is not
72 straightforward, especially when they do not directly implicate OXPHOS³¹, implicating
73 hitherto unknown mechanisms. Studies have identified mitophagy due to accumulation of
74 reactive oxygen species (ROS)^{32,33} and impairment of intra-mitochondrial protein
75 synthesis^{34,35} as potential mechanisms behind mtDNA associations with diseases, in addition
76 to those that directly impact ETC complex functions and OXPHOS efficiency^{32,36}. Most of
77 the mtDNA associations examined, however, are rare variations with large effects on rare and
78 severe diseases. Functional analysis of mtDNA variations altering risks to common diseases
79 are relatively underexplored, and more difficult due to their smaller effect sizes.

80
81 To bridge this gap and discover new ways that mtDNA variants contribute to
82 physiology and complex diseases, we took a phenome-driven and unbiased approach to
83 survey the impact of mtDNA polymorphisms on a wide set of 5,689 molecular and metabolic
84 traits beyond ATP synthesis (**Figure 1A**). Using mtDNA variations identified from whole-
85 genome and whole-exome sequencing in 16,220 healthy individuals, we found novel
86 associations between mtDNA variations in Haplogroup Uk with the metabolite N-
87 formylmethionine (fMet). We followed up on this finding with analysis of mtDNA effects on
88 gene expression in 44 tissues from the GTEx Consortium³⁷, and dissected the molecular
89 consequences of the associations using human cytoplasmic hybrid (cybrids) cell lines,
90 including the biogenesis of mitochondrial complexes, efficiency of OXPHOS, and both
91 cytoplasmic and intra-mitochondrial protein synthesis and degradation. Finally, we examined
92 the relevance of our findings in health and disease by verifying the effects of fMet and
93 mtDNA Haplogroup Uk in disease and longitudinal cohorts. Our findings explain how
94 common mtDNA variants can impact cellular proteostasis, providing a hitherto unknown
95 genetic checkpoint and an easily measured circulating biomarker for late-onset diseases.

96 97 **Results**

98 99 *mtDNA genotyping in 16,220 individuals*

100
101 To identify mtDNA polymorphic sites for association analysis, we obtained short-read
102 whole-genome sequencing (WGS, mean coverage = 26.8x, SD = 3.1x, in N=12,111
103 participants) and whole-exome sequencing (WES, mean coverage = 48.0x, SD = 8.6x, in
104 N=4,470 participants) data in a total of 16,220 unrelated European descent participants in the
105 INTERVAL study³⁸ (**Supplementary Figure 1**). We recovered a mean coverage of 2022.6x
106 (SD = 566.5x) and 30.6x (SD = 13.5x) on the mtDNA from WGS and WES respectively, and
107 identified 5,247 homoplasmic variants from WGS and WES (of which 5,161 are single
108 nucleotide polymorphisms [SNPs]) using GATK HaplotypeCaller v4³⁹ (**Supplementary**
109 **Figure 1**).

110
111 We took rigorous steps to avoid spurious variant calls, leveraging reference mtDNA
112 genotyping results from Affymetrix array data to recalibrate genotype quality filters, and high
113 coverage on mtDNA from the WGS for identifying loci with high frequency of heteroplasmic
114 variants that may be misidentified as homoplasmic (**Supplementary Figure 2**). For
115 subsequent analyses we considered the 396 high-confidence SNPs that could be accurately
116 identified from both WES and WGS; of these, 184 were common with minor population
117 allele frequency (MAF) greater than 0.01 (**Supplementary Figure 2, Supplementary Table**
118 **1**). We use these 184 SNPs, consistently in all following analyses on the INTERVAL cohort.
119 Principal component analysis using common mtDNA SNPs demonstrated the expected

120 clustering of individuals by their mitochondrial haplogroup (**Figure 1B**), distinct from
121 population structure identified using SNPs from the nDNA (**Figure 1 C,D**).

122

123 ***Common mtDNA variants are associated with blood N-formylmethionine***

124 We considered 5,689 distinct quantitative phenotypes representing a broad spectrum of
125 biological processes and pathways for association with mtDNA variations. These included a
126 total of 36 haematological traits, 1,344 small-molecule metabolites and 4,309 proteins
127 measured in whole blood or plasma using nine high-throughput phenotyping platforms
128 (**Figure 1A, Supplementary Table 2**). Many of these measures have been previously
129 reported to be affected by genetic variants in the nDNA^{40,41}. To specifically investigate the
130 effect of mtDNA variants on these biomarkers, we tested the association of each of the 184
131 common mtDNA SNPs with each of the phenotypes using linear mixed models (LMM,
132 **Online Methods**)⁴².

133

134 One metabolite, N-formylmethionine (fMet), was significantly associated with
135 mtDNA variants (5% false discovery rate [FDR], adjusted for SNPs and phenotypes;
136 equivalent to $P=1.65 \times 10^{-7}$), driven by eight mtDNA SNPs (**Figure 2A**). These represent
137 three independent associations (**Figure 2B,C**), two with increased fMet levels (top SNP
138 mt.1811A>G in *MT-RNR2*, MAF= 0.129, $P=3.03 \times 10^{-10}$, Beta[SE]= 0.12[0.02], and
139 mt.1189T>C in *MT-RNR1*, MAF= 0.073, $P=7.57 \times 10^{-8}$, Beta[SE]= 0.13[0.03]), and one with
140 decreased fMet (mt.3992C>T, T229M in *MT-ND1*, MAF= 0.02, $P=7.86 \times 10^{-6}$, Beta[SE]=
141 -0.19[0.04]). Twenty-eight more mtDNA SNPs reached suggestive levels of association with
142 one or more of the proteins or metabolites ($P \leq 5 \times 10^{-5}$, **Supplementary Table 3**). To test the
143 robustness of the association with fMet, we first sought to assess the contribution of potential
144 biological pathways or bacterial exposure that may increase exogenous fMet. Using the
145 Covariates in Multi-phenotype Studies (CMS)⁴³, we selected informative covariates from 896
146 metabolites and 36 blood cell measures (which can act as proxies for immune reactions
147 against infections) and we assessed the association between mtDNA SNPs and fMet
148 accounting for their contributions. Accounting for these additional covariates increased the
149 statistical significance of the associations between mtDNA SNPs and fMet levels (e.g. at
150 mt.1811A>G, CMS $P=4.63 \times 10^{-18}$, Beta[SE]= 0.06[0.01]) and identified seven additional
151 SNPs, yielding a total of 15 fMet-associated variants (**Figure 2B**). We verified that these
152 findings are not likely due to mtDNA copy number differences between haplogroups (**Online**
153 **Methods, Supplementary Figure 3**), or artefacts introduced during our processing of
154 metabolite data (**Online Methods, Supplementary Figure 4**).

155

156 We found that 14 of the 15 mtDNA SNPs associated with fMet levels (**Figure 2B**)
157 defined a single phylogenetic lineage, super-haplogroup Uk (containing branches Uk1 and
158 Uk2, **Figure 2C,D**). This includes two out of the three independent SNPs found to be
159 associated with fMet, mt.1811A>G, which resides on super-haplogroup Uk, and
160 mt.1189T>C, which resides on haplogroup Uk1. The remaining SNP mt.3992C>T resides on
161 the H4 branch of Haplogroup H4'9⁴⁴. In a LMM association analysis between all 896
162 metabolites and haplogroups (instead of mtDNA SNPs), we found that only Haplogroup Uk
163 was significantly associated with fMet ($P=4.41 \times 10^{-8}$, Beta[SE]= 0.25[0.05], at $P < 0.0029 =$
164 $0.05/17$ haplogroups tested). We then reassessed all 15 fMet-associated mtDNA SNPs
165 conditioning on Haplogroup Uk, and confirmed that associations of SNPs mt.1811A>G
166 (conditional $P=1.37 \times 10^{-3}$, Beta[SE]= 0.10[0.03] and mt.3992C>T (conditional $P=2.57 \times 10^{-5}$,
167 Beta[SE]= -0.18[0.04]) were independent from haplogroup Uk (**Supplementary Table 4**).

168

169

170 Finally, we sought replication in plasma samples from 11,538 participants from the
171 EPIC-Norfolk study⁴⁵. Of the 15 variants associated with fMet levels in INTERVAL, eight
172 were available in EPIC-Norfolk, with at least one SNP representing each of the three
173 independent signals (). All eight SNPs replicated the fMet association ($P < 0.0063 = 0.05/8$,
174 Bonferroni adjusted for eight SNPs tested, **Figure 2B**).

176 *Relative contribution of mtDNA and nDNA to fMet levels*

177
178 A first question is whether the three mtDNA variants (mt.1811A>G on Uk,
179 mt.1189T>C on Uk1 and mt.3992C>T on H4) exerted effects on fMet independently of the
180 nuclear genome. We considered 5,577,007 array-imputed nuclear DNA (nDNA) SNPs with
181 $MAF \geq 0.05$, and used a LMM model to test for association with fMet (**Online Methods**).
182 This analysis yielded an association at intronic variants (top SNP rs550045, chr9:130477160,
183 $MAF = 0.488$, $P = 1.14 \times 10^{-26}$, $Beta[SE] = 0.19[0.02]$, **Figure 3A**) in *PTRHI*, encoding a
184 human homolog of the yeast peptidyl-tRNA hydrolase 1 gene, with unverified hydrolase
185 function in humans. Conditional analyses showed that all mtDNA effects on fMet were
186 independent of the nDNA association at *PTRHI* (LMM $P < 1.86 \times 10^{-3}$; **Supplementary Table**
187 **4**). Further, no significant interaction effects were found between the top nDNA SNP
188 rs550045 and the top mtDNA SNP mt.1811A>G (interaction $P = 0.95$, **Figure 3B**). Overall,
189 mtDNA SNPs collectively contribute to 5.85% (SE= 1.20%) of variance in fMet levels, as
190 compared to 14.00% (SE= 3.32%) explained by nDNA SNPs and 59.03% (SE= 1.75%) by
191 unknown factors captured through all plasma metabolites and blood cell counts (**Figure 3C**,
192 **Supplementary Discussion**).

194 *fMet modulates mitochondrial function throughout coordination of mitochondrial* 195 *transcription-translation*

196
197 fMet is the initiation amino acid intra-mitochondrial protein translation^{46,47}. Thus, we
198 hypothesized that variations in fMet levels would be accompanied by changes in mtDNA
199 gene expression and proteostasis. We first sought to assess whether fMet-associated
200 mitochondrial variants impact on intra-mitochondrial gene expression, and assessed mtDNA
201 SNP effects on mtDNA-encoded transcript levels in 41 primary tissues from a total of 456
202 unrelated donors of European descent in the GTEx Consortium v7³⁷ (**Figure 1D,E**
203 **Supplementary Table 5**, **Supplementary Table 6**, **Supplementary Figure 5**). Five
204 transcripts were associated with mtDNA SNPs (*MT-ND1*, *MT-ND3*, *MT-ND4*, *MT-CO3*, *MT-*
205 *CYB*) at 5% study-wide FDR, in 29 out of the 41 tissues tested (**Figure 3D**, **Supplementary**
206 **Tables 7,8**). Strikingly, the 14 fMet-associated mtDNA SNPs on Haplogroup Uk accounted
207 for 93.75% of top eQTLs on mtDNA encoded genes at 5% study-wide FDR (30 out of 32,
208 **Supplementary Table 7**), with associations with *MT-ND3* accounting for 87.5% (28 out of
209 32, **Figure 3E**) of the top eQTLs, and associations with genes encoding subunits in Complex
210 I accounting for 93.3% (28 out of 30).

211
212 We next assessed the effects of fMet levels *in vitro* using trans-mitochondrial
213 cytoplasmic hybrids (cybrids). Cybrids were generated by fusing a single donor cell line
214 depleted of mtDNA⁴⁸ with cytoplasts from donors with mtDNA of haplogroups Uk and H
215 respectively¹⁰ (**Figure 4A**, **Online Methods**). We used four transmitochondrial cybrid cell
216 lines generated from four different healthy donors for each mtDNA haplogroup. As all
217 cybrid lines contain the same nuclear DNA so that any functional differences between them
218 is due to differences in the mtDNA sequence. (**Supplementary Table 9**). We found higher
219 fMet levels in Uk cybrids than H cybrids ($P = 0.02$, two-tailed t-test), in line with our

220 population-level findings (**Figure 4B**). As shown in **Figure 4C** and **Extended Figure 1A**,
221 fMet is synthesized by products of the one-carbon metabolism in mitochondria.

222

223 To determine the source of the different fMet levels between the different mtDNA
224 haplogroups we first looked for evidence of increased synthesis. This included an analysis of
225 one carbon metabolism proteins which supplement the metabolites necessary for fMet
226 synthesis (**Extended Data Figure 1A-B**); methionyl-tRNA formyltransferase (MTFMT)
227 (**Extended Data Figure 1A-D**); and mitochondrial ribosomal proteins (for example,
228 MRPL19). We found no differences between haplogroup Uk and H cybrids (**Extended Data**
229 **Figure 1C**). Next, we looked for differences in the degradation of fMet by studying levels of
230 the peptide deformylase (PDF) (**Extended Data Figure 1C**) and its deformylation bioproduct
231 formate (**Extended Data Figure 1D**). Again, we saw no differences between haplogroup Uk
232 and H cybrids. Next, we explored whether the increased fMet in haplogroup Uk was caused
233 by the accumulation of fMet due to decreased protein synthesis, by studying serum samples
234 from patients with a rare genetic mitochondrial disorder caused by a mutation in the mtDNA-
235 encoded gene for tRNA Leucine (m.3243A>G). Despite the known profound defect of intra-
236 mitochondrial translation^{49,50}, we did not observe any difference in fMet levels between the
237 patient serums and controls (**Extended Data Figure 1E**), indicating that blocking protein
238 synthesis itself does not lead to an increase in fMet levels. This raises the possibility of other
239 sources not known at present, including the release of fMet from unstable supercomplexes⁵¹.

240

241 Our findings raised the possibility that previously reported differences in intra-
242 mitochondrial protein synthesis between the haplogroups H and Uk¹⁰ could be accounted for
243 by the observed differences in fMet. In previous work, increasing fMet in fibroblasts by over-
244 expressing MTFMT also decreased intra-mitochondrial protein synthesis⁵². In keeping with
245 this hypothesis, treating the cybrid lines with exogenous fMet at a similar concentration seen
246 in plasma significantly increased the intracellular fMet levels ~1.19-fold (**Extended Data**
247 **Figure 2A**). This was similar to the fold-differences observed in INTERVAL between the
248 individuals carrying mtDNA SNPs alleles from haplogroup Uk as compared to those carrying
249 mtDNA SNP alleles from haplogroup H (fold change = 1.25, SE=0.05, P= 1.7x10⁻⁷,
250 **Extended Data Figure 2B**). This increase significantly decreased intra-mitochondrial protein
251 synthesis in the cells of both haplogroups (P = 0.04, 2-way-ANOVA; **Figure 4D**).

252

253 This was accompanied by increased levels of *MT-CO3* transcripts (**Figure 4E**)
254 (consistent with eQTL analysis in GTEx as shown in **Extended Data Figure 3**), which is the
255 characteristic compensatory response to reduced intra-mitochondrial translation⁴⁹. We
256 therefore asked if altering fMet levels through downregulation of mitochondrial methionyl-
257 tRNA formyltransferase (MTFMT) using siRNA (**Extended Data Figure 4A**) would result
258 in a similar compensatory response. Downregulation of MTFMT reduces fMet and the
259 synthesis of fMet-dependent proteins in complex I and IV^{46,47}, resulted in lower levels of the
260 fMet-dependent protein MT-CO1 (**Extended Data Figure 4A**) and a parallel increase in
261 transcript levels for the fMet-independent *MT-CO3* (**Extended Data Figure 4B**). Taken
262 together, these findings indicate that different fMet levels associated with common mtDNA
263 polymorphisms modulate intra-mitochondrial protein synthesis under homeostatic conditions
264 (**Figure 4D**).

265

266 Next, we studied the downstream consequences of fMet on mitochondrial oxidative
267 phosphorylation. In keeping with the observed decrease in mtDNA-encoded protein MT-CO1
268 in the cell lines from the haplogroup H (P= 0.016, 2-way-ANOVA; **Extended Data Figure**
269 **4D**), fMet supplementation reduced the abundance of the N-formylation dependent⁵³

270 complexes I and IV ($P = 0.02$, 2-way-ANOVA; **Figure 4F, Extended Data Figure 4E&F**).
271 This reduction was accompanied by lowered enzyme activities⁵⁴ ($P \leq 0.0005$; **Figure 4G**),
272 abolishing the differences between the haplogroups. We, however, did not observe this
273 decrease in the mitochondrial complexes III and V, which are less dependent on fMet⁵¹
274 (**Extended Data Figure 4G-H**). As expected, the effect of fMet supplementation on ETC
275 complexes I and IV abundance was associated with a decrease in oxygen consumption in
276 both haplogroup H and Uk cybrids ($P \leq 0.03$, 2-way-ANOVA; **Figure 4H**), and consequent
277 increase in both glycolytic ATP ($P = 0.02$; **Figure 4I**) and cytoplasmic ROS ($P \leq 0.0005$;
278 **Figure 4J**) as previously observed in mouse fibroblasts lacking MTFMT⁵¹, the main enzyme
279 required for fMet synthesis. fMet had no significant effect on mitochondrial ATP levels ($P >$
280 0.05 , **Extended Data Figure 4J**), mitochondrial mass, membrane potential or intra-
281 mitochondrial reactive oxygen species (ROS) ($P > 0.05$, **Extended Data Figure 4K**),
282 reflecting the relative sparing of complex III and V⁵¹. Thus, the differences in intra-
283 mitochondrial protein synthesis modulated by fMet have downstream effects on
284 mitochondrial respiratory chain function, and transcription mediated through the N-
285 formylation of specific intra-mitochondrial proteins. Our findings suggest that fMet levels are
286 physiologically balanced in the haplogroup H and Uk cybrids due to additional factors that
287 are poorly understood at present. Adding more fMet disrupted this equilibrium resulting in
288 the down-stream consequences that we observed.

289

290 *fMet regulates cellular stress and cytosolic proteostasis*

291

292 Although N-formylation of methionine is important for intra-mitochondrial protein
293 synthesis⁴⁷, its effects are less well documented in the cytosol, particularly in mammals.
294 Interestingly, we observed that fMet supplementation at levels seen in plasma globally
295 suppressed the *de novo* cytosolic translation ($P = 0.0001$, 3-way-ANOVA; **Figure 5A**). In
296 accordance with this, we observed activation of the eukaryotic translation initiation factor 2A
297 (EIF2A) the downstream kinase of the highly conserved integrated stress response (ISR)
298 pathway that leads to global repression of protein translation⁵⁵ (**Figure 5B, Extended Data**
299 **Figure 5A**). Consistently with the activation of EIF2A we also saw increased mRNA
300 expression of its downstream target activation factor 4 (ATF4)⁵⁶ and the CCAAT-enhancer-
301 binding protein homologous protein (CHOP)⁵⁷, enough to abolish the basal differences
302 between the haplogroups ($P \leq 0.03$, 2-way-ANOVA, **Figure 5C**). These observations were
303 independent of the ATF5-mitochondrial chaperone dependent mitochondrial unfolded
304 response⁵⁸ (**Extended Data Figure 5B**). Exogenous fMet also abolished the difference
305 between haplogroups Uk and H seen under basal conditions, implicating fMet in modulating
306 whole cell differences in protein homeostasis associated with the different mtDNA
307 polymorphisms.

308

309 The EIF2A/ATF4 pathway mediates repression of protein translation through
310 inhibition of the mammalian target of rapamycin (mTORC1) and thus growth⁵⁹. However, we
311 did not find any differences in mTORC1 activation (**Extended Data Figure 5A, C**) or
312 growth (**Extended Data Figure 5D-E**), implying an alternative mechanism for fMet in
313 cytosolic proteostasis. In *Saccharomyces cerevisiae*, the formylation of methionine in the
314 cytoplasm, has been proposed as a new protein degradation mechanism (N-degron) under
315 stress conditions, in a process dependent on the amino acid sensor GNC2 (general control
316 nonderepressible 2)⁶⁰ and presumably its downstream target EIF2A⁶¹. In line with these
317 observations, supplementation with fMet significantly increased the levels of ubiquitinated
318 proteins in both H and Uk cybrids ($P \leq 0.0002$, 2-way-ANOVA; **Figure 5D, Extended Data**
319 **Figure 5F**). These are in keeping with fMet also being an N-degron in humans. Thus, fMet

320 modulates cytosolic protein homeostasis at multiple levels, including translation (**Figure 5A**)
321 and degradation (**Figure 5D, Extended Data Figure 5F**).

322

323 Given our previous findings that fMet modulates mitochondrial and cytosolic protein
324 homeostasis, we determined the potential downstream consequences of the mitochondrial
325 background H and Uk by comparing the transcriptome of individuals with haplogroup Uk
326 against those with haplogroup H across 49 tissues using again data from the GTEx
327 Consortium³⁷. This revealed 1 to 619 differentially expressed genes (total 4,244 genes) in 47
328 out of 49 tissues using a quasi-likelihood F test in edgeR^{62,63} at 5% tissue-wide FDR
329 (**Supplementary Tables 10,11**), with the majority (94.5%) of differentially expressed genes
330 being more highly expressed in haplogroup Uk across all tissues (**Supplementary Table**
331 **10,11**). Pathway (**Figure 5E**) and gene set enrichment analysis (GSEA) (**Supplementary**
332 **Tables 12,13**) of all the differentially expressed genes showed an enrichment for pathways
333 involved in metabolism and immunity in addition to protein homeostasis pathways and
334 ribosomal translation initiation (**Figure 5F**). In particular, we found three significant
335 differential expression signals (at 5% tissue-wide FDR) among nuclear DNA encoded
336 mitochondrial ribosomal genes and genes involved in the processing of mitochondrial
337 rRNAs: *RMRP* (logFC=2.58, FDR=3.93x10⁻⁸) in coronary artery, and *MRPS6* (logFC=0.82,
338 FDR=2.24x10⁻²) and *MRPL14* (logFC=0.68, FDR=3.07x10⁻²) in the prostate
339 (**Supplementary Figure 6**). Taken together, these independent observations validate our
340 earlier *in vitro* findings, and indicate that mtDNA haplogroup-associated differences in fMet
341 have multiple potential downstream consequences for cellular function beyond bioenergetics,
342 oxidative phosphorylation, and mitochondrial ATP synthesis (**Extended Data Figure 4J-K**).

343

344 *fMet levels mediate late-onset disease risk*

345 Haplogroup Uk has been previously associated with reduced risk of developing late-onset
346 diseases including ischemic stroke (IS)²³ and Parkinson's disease (PD)²², but the reasons for
347 this are not known. To explore these potential mechanisms, we measured blood fMet levels
348 in an IS cohort (282 cases, 181 controls, **Supplementary Table 14**) where we previously
349 described a protective effect of haplogroup Uk²³. In order to remove haplogroup effects on IS
350 when testing for fMet associations with the disease, we enriched for individuals of
351 haplogroup Uk in both cases and controls (32% haplogroup Uk in cases, 23% haplogroup Uk
352 in controls). We then asked whether fMet is associated with IS, and if its effects can be due to
353 haplogroup differences. We found a marginal and negative association between fMet and IS
354 (OR= 0.83, SE= 0.08, logistic regression *P*= 0.06). This association is present in non-Uk
355 individuals (OR= 0.77, SE= 0.08, logistic *P*= 0.02), but not in those with haplogroup Uk
356 (OR= 2.24, SE= 1.06, logistic *P*= 0.09) (**Figure 6A**). Our results are consistent with a
357 potential involvement of fMet in IS etiology, only part of which is due to mtDNA
358 haplogroups. We also considered a PD cohort without enrichment for haplogroup Uk (120
359 cases, 43 controls, **Supplementary Table 14**), which however did not reveal any
360 associations, possibly due to low sample sizes.

361

362 Next, we investigated whether fMet levels may be associated with the risk of other
363 ageing-related diseases (cardiometabolic and common cancers, **Supplementary Table 15**),
364 and if those associations were mtDNA haplogroup-dependent. We used Cox-proportional
365 hazards models to test associations between fMet and incident risk of 24 non-communicable
366 diseases and all-cause mortality in 11,966 individuals from the EPIC-Norfolk study covering
367 more than 20-years of follow-up (**Online Methods**). We observed significant (*P*< 0.002;
368 0.05/23 tests) positive associations between fMet levels and incident renal disease, heart
369 failure, coronary artery disease, abdominal aortic aneurysms, peripheral artery disease and

370 chronic obstructive pulmonary disease (COPD) as well as mortality (**Figure 6B**,
371 **Supplementary Table 16**). Hazard ratios ranged between 1.10 and 1.29 per 1 standard
372 deviation (SD) increase in log-transformed fMet levels. As fMet was correlated with age
373 (Pearson correlation of age and fMet $r=0.33$, $P=2.48 \times 10^{-307}$, **Figure 6C**), all Cox models
374 accounted for age. To ensure the age association does not violate the proportional hazards
375 assumption, we investigated Schoenfeld residuals and age-interaction terms, and none
376 showed evidence of violation (**Supplementary Table 16**). Associations did not differ
377 significantly between the Uk (N= 895) and other haplogroups (N= 9,887), however
378 confidence intervals were wide in the smaller Uk haplogroup for some outcomes
379 (**Supplementary Table 16, Figure 6D**). While this highlights the low power of the current
380 study to detect heterogeneity between groups, the potential role of fMet as a marker of
381 ageing-related diseases in different haplogroups warrants further investigations.

382

383 Discussion

384

385 In this study, we profiled more than 5,000 molecular traits in a healthy population-
386 based cohort, and found novel associations between three mtDNA variants in Haplogroups
387 Uk and H4 and the metabolite fMet. Whilst it is possible that the differences in fMet levels
388 and mitochondrial transcription and translation are due to independent effects of haplogroup-
389 specific variants rather than through a common causal pathway, it is not clear how this would
390 occur, and results from our experiments are consistent with the latter. In fact, two of the
391 variants associated with fMet (**Figure 2B**) affect the non-coding D-loop which is involved in
392 regulation mtDNA transcription, and two also involve the rRNA genes directly involved in
393 protein synthesis. Four variants affect the amino acid sequence of critical respiratory chain
394 proteins, potentially influencing their function^{10,64}, assembly or stability⁶⁵, with the non-
395 synonymous variants being associated indirectly through co-inheritance on the same mtDNA
396 haplogroup. This complex scenario highlights the need of future experiments dissecting the
397 effect of each independent variant.

398

399 fMet is the initiation amino acid for intra-mitochondrial translation⁶⁶. Previous studies
400 have shown that fMet is not necessary for initiation of translation or stability of newly
401 synthesized polypeptides. However, a lack of fMet decreases synthesis of mtDNA encoded
402 proteins and their integration into OXPHOS complexes and supercomplexes: MTFMT
403 knockouts or mutants⁴⁷ display inefficient OXPHOS and increased risk of disease^{47,51}. On the
404 other hand, increasing fMet by MTFMT overexpression⁵² and our experiments with fMet
405 supplementation also compromised mitochondrial protein synthesis, OXPHOS complex
406 levels and respiratory supercomplex function⁵¹. This implies that under physiological
407 conditions, fMet is maintained within a narrow window, and increasing or decreasing fMet
408 can have detrimental effects. Our findings demonstrate a role for fMet-associated mtDNA
409 variants regulating levels of intra-mitochondrial gene expression, and modulating intra-
410 mitochondrial protein synthesis and OXPHOS complex formation under homeostatic
411 conditions through fMet. This may have tissue-specific and age-cumulative effects on
412 metabolism and disease risk.

413

414 One of our most intriguing findings is that of a mechanism of cellular proteostasis that
415 is modulated by mtDNA. Within mitochondria, mtDNA haplogroup-dependent fMet levels
416 directly and specifically affect the abundance of mtDNA encoded, N-formylated ETC
417 subunits and complexes. In the cytosol, fMet exerts indirect influence on global intra-
418 mitochondrial protein synthesis through transcription and proteolysis without effects on cell
419 growth (**Extended Data Figure 5D,E**) and/or mitochondrial membrane stability (**Extended**

420 **Figure 4J,K**); higher fMet levels in individuals from haplogroup Uk increase the ubiquitin-
421 targeted N-degron mediated proteolysis (**Figure 5G**), and thus decrease the formation of
422 protein aggregates⁶⁷ and the regulation of apoptosis⁶⁸. This can explain the previously found
423 protective effects of haplogroup Uk on late-onset neurodegenerative disorders⁶⁹. In addition
424 to the elimination of mis-folded proteins, N-end rule pathways also play a role in controlling
425 subunit stoichiometries⁷⁰ in protein complexes such as the ETC, and the elimination of
426 proteins mis-localised from their primary cellular compartment⁷¹. Furthermore, the profound
427 reduction in cytosolic protein synthesis is likely to have multiple downstream effects on cell
428 function. This emphasizes the importance of maintaining fMet levels within a narrow
429 physiological range and the need of future studies dissecting its fluctuations in different
430 tissues and disease models.

431

432 It is therefore plausible that fMet is involved in ‘matching’ protein synthesis with the
433 mitochondrial and cytosolic compartments⁷² in response to cellular bioenergetic needs in a
434 tissue-specific manner. Subtle differences in fMet, partly attributable to mtDNA haplogroup
435 effects, could have a cumulative effect on proteostasis and degradation throughout life, and
436 thereby modify the risk of developing several late-onset diseases. fMet levels were
437 significantly associated with late-onset diseases in the EPIC Norfolk cohort independent of
438 age. Further experimental work is needed to definitively prove a causal role for fMet,
439 however, given that it can be readily measured in serum as a circulating biomarker of cellular
440 proteostasis, fMet is likely to be valuable for monitoring new treatments across a wide range
441 of common human disorders.

442

443 In conclusion, the use of deep molecular phenotyping based on high-throughput
444 metabolomics, transcriptomics, and proteomics is proving effective in identifying molecular
445 hypotheses underpinning genetic associations with health and disease endpoints^{40,41,73,74}. Our
446 findings open up possibilities for further investigations into mtDNA control over metabolism
447 and cellular physiology, and its implications on human health and disease. fMet may not be
448 the only metabolite mtDNA variants regulated; most plasma biomarkers included in our
449 study were only assayed in around 3,000 individuals, limiting statistical power. Continued
450 investigation of mtDNA effects on biomarkers may lead to further elucidation of
451 mitochondria’s role in cellular physiology and function.

452

453 **References**

- 454 1. Anderson, S. *et al.* Sequence and organization of the human mitochondrial genome.
 455 *Nature* **290**, 457–465 (1981).
- 456 2. Giles, R. E., Blanc, H., Cann, H. M. & Wallace, D. C. Maternal inheritance of human
 457 mitochondrial DNA. *Proc. Natl. Acad. Sci. U. S. A.* **77**, 6715–6719 (1980).
- 458 3. Howell, N. Mutational analysis of the human mitochondrial genome branches into the
 459 realm of bacterial genetics. *American journal of human genetics* vol. 59 749–755 (1996).
- 460 4. Lippold, S. *et al.* Human paternal and maternal demographic histories: insights from
 461 high-resolution Y chromosome and mtDNA sequences. *Investig. Genet.* **5**, 13 (2014).
- 462 5. Stone, A. C. & Stoneking, M. mtDNA analysis of a prehistoric Oneota population:
 463 implications for the peopling of the New World. *Am. J. Hum. Genet.* **62**, 1153–1170
 464 (1998).
- 465 6. Vigilant, L., Stoneking, M., Harpending, H., Hawkes, K. & Wilson, A. C. African
 466 populations and the evolution of human mitochondrial DNA. *Science* **253**, 1503–1507
 467 (1991).
- 468 7. Ruiz-Pesini, E., Mishmar, D., Brandon, M., Procaccio, V. & Wallace, D. C. Effects of
 469 purifying and adaptive selection on regional variation in human mtDNA. *Science* **303**,
 470 223–226 (2004).
- 471 8. Lane, N. & Martin, W. The energetics of genome complexity. *Nature* **467**, 929–934
 472 (2010).
- 473 9. Ji, F. *et al.* Mitochondrial DNA variant associated with Leber hereditary optic
 474 neuropathy and high-altitude Tibetans. *Proc. Natl. Acad. Sci. U. S. A.* **109**, 7391–7396
 475 (2012).
- 476 10. Gómez-Durán, A. *et al.* Unmasking the causes of multifactorial disorders: OXPHOS
 477 differences between mitochondrial haplogroups. *Hum. Mol. Genet.* **19**, 3343–3353
 478 (2010).
- 479 11. Suissa, S. *et al.* Ancient mtDNA genetic variants modulate mtDNA transcription and
 480 replication. *PLoS Genet.* **5**, e1000474 (2009).
- 481 12. Tranah, G. J. *et al.* Mitochondrial DNA variation in human metabolic rate and energy
 482 expenditure. *Mitochondrion* **11**, 855–861 (2011).
- 483 13. Gorman, G. S. *et al.* Prevalence of nuclear and mitochondrial DNA mutations related to
 484 adult mitochondrial disease. *Ann. Neurol.* **77**, 753–759 (2015).
- 485 14. Ruiz-Pesini, E. *et al.* An enhanced MITOMAP with a global mtDNA mutational
 486 phylogeny. *Nucleic Acids Res.* **35**, D823–8 (2007).
- 487 15. Stewart, J. B. & Chinnery, P. F. The dynamics of mitochondrial DNA heteroplasmy:
 488 implications for human health and disease. *Nat. Rev. Genet.* **16**, 530–542 (2015).
- 489 16. Zong, W.-X., Rabinowitz, J. D. & White, E. Mitochondria and Cancer. *Mol. Cell* **61**,
 490 667–676 (2016).
- 491 17. Larman, T. C. *et al.* Spectrum of somatic mitochondrial mutations in five cancers. *Proc.*
 492 *Natl. Acad. Sci. U. S. A.* **109**, 14087–14091 (2012).
- 493 18. Yuan, Y. *et al.* Comprehensive molecular characterization of mitochondrial genomes in
 494 human cancers. *Nat. Genet.* **52**, 342–352 (2020).
- 495 19. Vander Heiden, M. G., Cantley, L. C. & Thompson, C. B. Understanding the Warburg
 496 effect: the metabolic requirements of cell proliferation. *Science* **324**, 1029–1033 (2009).
- 497 20. Marom, S., Friger, M. & Mishmar, D. MtDNA meta-analysis reveals both phenotype
 498 specificity and allele heterogeneity: a model for differential association. *Sci. Rep.* **7**,
 499 43449 (2017).
- 500 21. Bi, R. *et al.* Mitochondrial DNA haplogroup B5 confers genetic susceptibility to
 501 Alzheimer’s disease in Han Chinese. *Neurobiol. Aging* **36**, 1604.e7–16 (2015).

- 502 22. Hudson, G. *et al.* Two-stage association study and meta-analysis of mitochondrial DNA
503 variants in Parkinson disease. *Neurology* **80**, 2042–2048 (2013).
- 504 23. Chinnery, P. F., Elliott, H. R., Syed, A., Rothwell, P. M. & Oxford Vascular Study.
505 Mitochondrial DNA haplogroups and risk of transient ischaemic attack and ischaemic
506 stroke: a genetic association study. *Lancet Neurol.* **9**, 498–503 (2010).
- 507 24. Nishigaki, Y. *et al.* Mitochondrial haplogroup N9b is protective against myocardial
508 infarction in Japanese males. *Hum. Genet.* **120**, 827–836 (2007).
- 509 25. Kofler, B. *et al.* Mitochondrial DNA haplogroup T is associated with coronary artery
510 disease and diabetic retinopathy: a case control study. *BMC Medical Genetics* vol. 10
511 (2009).
- 512 26. Chinnery, P. F. *et al.* Mitochondrial DNA haplogroups and type 2 diabetes: a study of
513 897 cases and 1010 controls. *J. Med. Genet.* **44**, e80 (2007).
- 514 27. Poulton, J. *et al.* Type 2 diabetes is associated with a common mitochondrial variant:
515 evidence from a population-based case-control study. *Hum. Mol. Genet.* **11**, 1581–1583
516 (2002).
- 517 28. Hudson, G., Gomez-Duran, A., Wilson, I. J. & Chinnery, P. F. Recent mitochondrial
518 DNA mutations increase the risk of developing common late-onset human diseases.
519 *PLoS Genet.* **10**, e1004369 (2014).
- 520 29. Marom, S., Friger, M. & Mishmar, D. MtDNA meta-analysis reveals both phenotype
521 specificity and allele heterogeneity: a model for differential association. *Sci. Rep.* **7**,
522 43449 (2017).
- 523 30. Wallace, D. C. A mitochondrial paradigm of metabolic and degenerative diseases, aging,
524 and cancer: a dawn for evolutionary medicine. *Annu. Rev. Genet.* **39**, 359–407 (2005).
- 525 31. Chinnery, P. F. & Gomez-Duran, A. Oldies but Goldies mtDNA Population Variants and
526 Neurodegenerative Diseases. *Frontiers in Neuroscience* vol. 12 (2018).
- 527 32. Trounce, I., Neill, S. & Wallace, D. C. Cytoplasmic transfer of the mtDNA nt 8993 T-
528 >G (ATP6) point mutation associated with Leigh syndrome into mtDNA-less cells
529 demonstrates cosegregation with a decrease in state III respiration and ADP/O ratio.
530 *Proceedings of the National Academy of Sciences* vol. 91 8334–8338 (1994).
- 531 33. Mattiazzi, M. The mtDNA T8993G (NARP) mutation results in an impairment of
532 oxidative phosphorylation that can be improved by antioxidants. *Human Molecular*
533 *Genetics* vol. 13 869–879 (2004).
- 534 34. Goto, Y., Nonaka, I. & Horai, S. A mutation in the tRNA(Leu)(UUR) gene associated
535 with the MELAS subgroup of mitochondrial encephalomyopathies. *Nature* **348**, 651–653
536 (1990).
- 537 35. van den Ouweland, J. M. *et al.* Maternally inherited diabetes and deafness is a distinct
538 subtype of diabetes and associates with a single point mutation in the mitochondrial
539 tRNA(Leu(UUR)) gene. *Diabetes* **43**, 746–751 (1994).
- 540 36. Brown, M. D., Trounce, I. A., Jun, A. S., Allen, J. C. & Wallace, D. C. Functional
541 analysis of lymphoblast and cybrid mitochondria containing the 3460, 11778, or 14484
542 Leber’s hereditary optic neuropathy mitochondrial DNA mutation. *J. Biol. Chem.* **275**,
543 39831–39836 (2000).
- 544 37. GTEx Consortium *et al.* Genetic effects on gene expression across human tissues.
545 *Nature* **550**, 204–213 (2017).
- 546 38. Moore, C. *et al.* The INTERVAL trial to determine whether intervals between blood
547 donations can be safely and acceptably decreased to optimise blood supply: study
548 protocol for a randomised controlled trial. *Trials* **15**, 363 (2014).
- 549 39. Poplin, R. *et al.* Scaling accurate genetic variant discovery to tens of thousands of
550 samples. doi:10.1101/2011178.
- 551 40. Sun, B. B. *et al.* Genomic atlas of the human plasma proteome. *Nature* **558**, 73–79

- 552 (2018).
- 553 41. Shin, S.-Y. *et al.* An atlas of genetic influences on human blood metabolites. *Nat. Genet.*
554 **46**, 543–550 (2014).
- 555 42. Speed, D. *et al.* Reevaluation of SNP heritability in complex human traits. *Nat. Genet.*
556 **49**, 986–992 (2017).
- 557 43. Aschard, H. *et al.* Covariate selection for association screening in multiphenotype
558 genetic studies. *Nat. Genet.* **49**, 1789–1795 (2017).
- 559 44. Derenko, M. *et al.* Western Eurasian ancestry in modern Siberians based on
560 mitogenomic data. *BMC Evol. Biol.* **14**, 217 (2014).
- 561 45. Day, N. *et al.* EPIC-Norfolk: study design and characteristics of the cohort. European
562 Prospective Investigation of Cancer. *Br. J. Cancer* **80 Suppl 1**, 95–103 (1999).
- 563 46. Rajbhandary, U. L. & Ming Chow, C. Initiator tRNAs and Initiation of Protein
564 Synthesis. *tRNA* 511–528 doi:10.1128/9781555818333.ch25.
- 565 47. Tucker, E. J. *et al.* Mutations in MTFMT Underlie a Human Disorder of Formylation
566 Causing Impaired Mitochondrial Translation. *Cell Metabolism* vol. 14 428–434 (2011).
- 567 48. King, M. P. & Attardi, G. Human cells lacking mtDNA: repopulation with exogenous
568 mitochondria by complementation. *Science* **246**, 500–503 (1989).
- 569 49. Picard, M. *et al.* Progressive increase in mtDNA 3243A>G heteroplasmy causes abrupt
570 transcriptional reprogramming. *Proc. Natl. Acad. Sci. U. S. A.* **111**, E4033–42 (2014).
- 571 50. King, M. P., Koga, Y., Davidson, M. & Schon, E. A. Defects in mitochondrial protein
572 synthesis and respiratory chain activity segregate with the tRNA(Leu(UUR)) mutation
573 associated with mitochondrial myopathy, encephalopathy, lactic acidosis, and strokelike
574 episodes. *Mol. Cell. Biol.* **12**, 480–490 (1992).
- 575 51. Arguello, T., Köhrer, C., RajBhandary, U. L. & Moraes, C. T. Mitochondrial methionyl-
576 formylation affects steady-state levels of oxidative phosphorylation complexes and their
577 organization into supercomplexes. *J. Biol. Chem.* **293**, 15021–15032 (2018).
- 578 52. Hinttala, R. *et al.* An N-terminal formyl methionine on COX 1 is required for the
579 assembly of cytochrome c oxidase. *Hum. Mol. Genet.* **24**, 4103–4113 (2015).
- 580 53. Neeve, V. C. M. *et al.* Clinical and functional characterisation of the combined
581 respiratory chain defect in two sisters due to autosomal recessive mutations in MTFMT.
582 *Mitochondrion* **13**, 743–748 (2013).
- 583 54. Nijtmans, L. G. J., Henderson, N. S. & Holt, I. J. Blue Native electrophoresis to study
584 mitochondrial and other protein complexes. *Methods* **26**, 327–334 (2002).
- 585 55. Wek, R. C., Jiang, H.-Y. & Anthony, T. G. Coping with stress: eIF2 kinases and
586 translational control. *Biochem. Soc. Trans.* **34**, 7–11 (2006).
- 587 56. Ameri, K. & Harris, A. L. Activating transcription factor 4. *Int. J. Biochem. Cell Biol.*
588 **40**, 14–21 (2008).
- 589 57. Su, N. & Kilberg, M. S. C/EBP homology protein (CHOP) interacts with activating
590 transcription factor 4 (ATF4) and negatively regulates the stress-dependent induction of
591 the asparagine synthetase gene. *J. Biol. Chem.* **283**, 35106–35117 (2008).
- 592 58. Quirós, P. M. *et al.* Multi-omics analysis identifies ATF4 as a key regulator of the
593 mitochondrial stress response in mammals. *J. Cell Biol.* **216**, 2027–2045 (2017).
- 594 59. Walter, P. & Ron, D. The unfolded protein response: from stress pathway to homeostatic
595 regulation. *Science* **334**, 1081–1086 (2011).
- 596 60. Kim, J.-M. *et al.* Formyl-methionine as an N-degron of a eukaryotic N-end rule pathway.
597 *Science* vol. 362 eaat0174 (2018).
- 598 61. Eldeeb, M. A., Fahlman, R. P., Esmaili, M. & Fon, E. A. Formylation of Eukaryotic
599 Cytoplasmic Proteins: Linking Stress to Degradation. *Trends in biochemical sciences*
600 vol. 44 181–183 (2019).
- 601 62. Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for

- 602 differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139–
603 140 (2010).
- 604 63. McCarthy, D. J., Chen, Y. & Smyth, G. K. Differential expression analysis of
605 multifactor RNA-Seq experiments with respect to biological variation. *Nucleic Acids*
606 *Res.* **40**, 4288–4297 (2012).
- 607 64. Pacheu-Grau, D., Gómez-Durán, A., López-Pérez, M. J., Montoya, J. & Ruiz-Pesini, E.
608 Mitochondrial pharmacogenomics: barcode for antibiotic therapy. *Drug Discov. Today*
609 **15**, 33–39 (2010).
- 610 65. Pello, R. *et al.* Mitochondrial DNA background modulates the assembly kinetics of
611 OXPHOS complexes in a cellular model of mitochondrial disease. *Hum. Mol. Genet.* **17**,
612 4001–4011 (2008).
- 613 66. Bianchetti, R., Lucchini, G., Crosti, P. & Tortora, P. Dependence of mitochondrial
614 protein synthesis initiation on formylation of the initiator methionyl-tRNA^f. *J. Biol.*
615 *Chem.* **252**, 2519–2523 (1977).
- 616 67. Sorrentino, V. *et al.* Enhancing mitochondrial proteostasis reduces amyloid- β
617 proteotoxicity. *Nature* **552**, 187–193 (2017).
- 618 68. Eldeeb, M. A., Fahlman, R. P., Esmaili, M. & Ragheb, M. A. Regulating Apoptosis by
619 Degradation: The N-End Rule-Mediated Regulation of Apoptotic Proteolytic Fragments
620 in Mammalian Cells. *Int. J. Mol. Sci.* **19**, (2018).
- 621 69. Hudson, G., Gomez-Duran, A., Wilson, I. J. & Chinnery, P. F. Recent mitochondrial
622 DNA mutations increase the risk of developing common late-onset human diseases.
623 *PLoS Genet.* **10**, e1004369 (2014).
- 624 70. Shemorry, A., Hwang, C.-S. & Varshavsky, A. Control of protein quality and
625 stoichiometries by N-terminal acetylation and the N-end rule pathway. *Mol. Cell* **50**,
626 540–551 (2013).
- 627 71. Wrobel, L. *et al.* Mistargeted mitochondrial proteins activate a proteostatic response in
628 the cytosol. *Nature* **524**, 485–488 (2015).
- 629 72. Couvillion, M. T., Soto, I. C., Shipkovenska, G. & Churchman, L. S. Synchronized
630 mitochondrial and cytosolic translation programs. *Nature* **533**, 499–503 (2016).
- 631 73. Suhre, K. *et al.* Connecting genetic risk to disease end points through the human blood
632 plasma proteome. *Nat. Commun.* **8**, 14357 (2017).
- 633 74. Yao, C. *et al.* Genome-wide mapping of plasma protein QTLs identifies putatively
634 causal genes and pathways for cardiovascular disease. *Nat. Commun.* **9**, 3268 (2018).
- 635 75. Astle, W. J. *et al.* The Allelic Landscape of Human Blood Cell Trait Variation and Links
636 to Common Complex Disease. *Cell* **167**, 1415–1429.e19 (2016).
- 637 76. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data.
638 *Nature* **562**, 203–209 (2018).
- 639 77. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-
640 generation DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).
- 641 78. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for
642 analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
- 643 79. Weissensteiner, H. *et al.* mtDNA-Server: next-generation sequencing data analysis of
644 human mitochondrial DNA in the cloud. *Nucleic Acids Res.* **44**, W64–9 (2016).
- 645 80. Van der Auwera, G. A. *et al.* From FastQ data to high confidence variant calls: the
646 Genome Analysis Toolkit best practices pipeline. *Curr. Protoc. Bioinformatics* **43**,
647 11.10.1–11.10.33 (2013).
- 648 81. 1000 Genomes Project Consortium *et al.* A global reference for human genetic variation.
649 *Nature* **526**, 68–74 (2015).
- 650 82. Manichaikul, A. *et al.* Robust relationship inference in genome-wide association studies.
651 *Bioinformatics* **26**, 2867–2873 (2010).

- 652 83. Kielbasa, S. M., Wan, R., Sato, K., Horton, P. & Frith, M. C. Adaptive seeds tame
653 genomic sequence comparison. *Genome Res.* **21**, 487–493 (2011).
- 654 84. Stegle, O., Parts, L., Piipari, M., Winn, J. & Durbin, R. Using probabilistic estimation of
655 expression residuals (PEER) to obtain increased power and interpretability of gene
656 expression analyses. *Nat. Protoc.* **7**, 500–507 (2012).
- 657 85. Lippert, C., Casale, F. P., Rakitsch, B. & Stegle, O. LIMIX: genetic analysis of multiple
658 traits. *Genetics* **833** (2014).
- 659 86. Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for
660 interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. U. S. A.* **102**,
661 15545–15550 (2005).
- 662 87. Mootha, V. K. *et al.* PGC-1alpha-responsive genes involved in oxidative
663 phosphorylation are coordinately downregulated in human diabetes. *Nat. Genet.* **34**,
664 267–273 (2003).
- 665 88. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler
666 transform. *Bioinformatics* **25**, 1754–1760 (2009).
- 667 89. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,
668 2078–2079 (2009).
- 669 90. Koboldt, D. C. *et al.* VarScan: variant detection in massively parallel sequencing of
670 individual and pooled samples. *Bioinformatics* **25**, 2283–2285 (2009).
- 671 91. Koboldt, D. C. *et al.* VarScan 2: somatic mutation and copy number alteration discovery
672 in cancer by exome sequencing. *Genome Res.* **22**, 568–576 (2012).
- 673 92. Wilm, A. *et al.* LoFreq: a sequence-quality aware, ultra-sensitive variant caller for
674 uncovering cell-population heterogeneity from high-throughput sequencing datasets.
675 *Nucleic Acids Res.* **40**, 11189–11201 (2012).
- 676 93. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic
677 variants from high-throughput sequencing data. *Nucleic Acids Res.* **38**, e164 (2010).
- 678 94. Torroni, A. *et al.* Classification of European mtDNAs from an analysis of three
679 European populations. *Genetics* **144**, 1835–1850 (1996).
- 680 95. van Oven, M. & Kayser, M. Updated comprehensive phylogenetic tree of global human
681 mitochondrial DNA variation. *Hum. Mutat.* **30**, E386–94 (2009).
- 682 96. Weissensteiner, H. *et al.* HaploGrep 2: mitochondrial haplogroup classification in the era
683 of high-throughput sequencing. *Nucleic Acids Res.* **44**, W58–W63 (2016).
- 684 97. Pereira, L., Soares, P., Radivojac, P., Li, B. & Samuels, D. C. Comparing phylogeny and
685 the predicted pathogenicity of protein variations reveals equal purifying selection across
686 the global human mtDNA diversity. *Am. J. Hum. Genet.* **88**, 433–439 (2011).
- 687 98. Levin, L., Zhidkov, I., Gurman, Y., Hawlena, H. & Mishmar, D. Functional recurrent
688 mutations in the human mitochondrial phylogeny: dual roles in evolution and disease.
689 *Genome Biol. Evol.* **5**, 876–890 (2013).
- 690 99. Bustin, S. A. *et al.* The MIQE guidelines: minimum information for publication of
691 quantitative real-time PCR experiments. *Clin. Chem.* **55**, 611–622 (2009).
- 692 100. Bradford, M. M. A rapid and sensitive method for the quantitation of microgram
693 quantities of protein utilizing the principle of protein-dye binding. *Anal. Biochem.* **72**,
694 248–254 (1976).
- 695 101. Wittig, I., Braun, H.-P. & Schägger, H. Blue native PAGE. *Nat. Protoc.* **1**, 418–428
696 (2006).
- 697 102. Yarnall, A. J. *et al.* Characterizing mild cognitive impairment in incident Parkinson
698 disease: the ICICLE-PD study. *Neurology* **82**, 308–316 (2014).
- 699 103. Gibb, W. R. & Lees, A. J. The relevance of the Lewy body to the pathogenesis of
700 idiopathic Parkinson’s disease. *J. Neurol. Neurosurg. Psychiatry* **51**, 745–752 (1988).
- 701 104. Rothwell, P. M. *et al.* Population-based study of event-rate, incidence, case fatality, and

- 702 mortality for all acute vascular events in all arterial territories (Oxford Vascular Study).
703 *Lancet* **366**, 1773–1783 (2005).
- 704 105. Floros, V. I. *et al.* Segregation of mitochondrial DNA heteroplasmy through a
705 developmental genetic bottleneck in human embryos. *Nat. Cell Biol.* **20**, 144–151
706 (2018).
- 707 106. Calabrese, C. *et al.* MToolBox: a highly automated pipeline for heteroplasmy annotation
708 and prioritization analysis of human mitochondrial variants in high-throughput
709 sequencing. *Bioinformatics* **30**, 3115–3117 (2014).
- 710 107. Xu, H. *et al.* FastUniq: a fast de novo duplicates removal tool for paired short reads.
711 *PLoS One* **7**, e52249 (2012).

712

713 **Acknowledgments:**

714 Participants in the INTERVAL randomised controlled trial were recruited with the active
715 collaboration of NHS Blood and Transplant England (www.nhsbt.nhs.uk), which has
716 supported field work and other elements of the trial. DNA extraction and genotyping was co-
717 funded by the National Institute for Health Research (NIHR), the NIHR BioResource
718 (<http://bioresource.nihr.ac.uk>) and the NIHR Cambridge Biomedical Research Centre (BRC-
719 1215-20014)*.

720 The academic coordinating centre for INTERVAL was supported by core funding from:
721 NIHR Blood and Transplant Research Unit in Donor Health and Genomics (NIHR BTRU-
722 2014-10024), UK Medical Research Council (MR/L003120/1), British Heart Foundation
723 (SP/09/002; RG/13/13/30194; RG/18/13/33946) and the NIHR Cambridge BRC (BRC-1215-
724 20014)*. A complete list of the investigators and contributors to the INTERVAL trial is
725 provided in Di Angelantonio E, et al. *Lancet*. 2017 Nov 25;390(10110):2360-2371. The
726 academic coordinating centre would like to thank blood donor centre staff and blood donors
727 for participating in the INTERVAL trial.

728 This work was supported by Health Data Research UK, which is funded by the UK Medical
729 Research Council, Engineering and Physical Sciences Research Council, Economic and
730 Social Research Council, Department of Health and Social Care (England), Chief Scientist
731 Office of the Scottish Government Health and Social Care Directorates, Health and Social
732 Care Research and Development Division (Welsh Government), Public Health Agency
733 (Northern Ireland), British Heart Foundation and Wellcome.

734 NC is supported by the EBI-Sanger Postdoctoral Fellowship. AGD receives funding from the
735 National Institute for Health Research (NIHR) Biomedical Research Centre based at
736 Cambridge University Hospitals NHS Foundation Trust. EY-D was funded by the Isaac
737 Newton Trust/Wellcome Trust ISSF / University of Cambridge Joint Research Grants
738 Scheme. JMMH was funded by the NIHR Cambridge BRC (BRC-1215-20014)*. PFC is a
739 Wellcome Trust Principal Research Fellow (212219/Z/18/Z), and a UK NIHR Senior
740 Investigator, who receives support from the Medical Research Council Mitochondrial
741 Biology Unit (MC_UU_00015/9), the Evelyn Trust, and the National Institute for Health
742 Research (NIHR) Biomedical Research Centre based at Cambridge University Hospitals
743 NHS Foundation Trust and the University of Cambridge. CHWG is supported by a
744 RCUK/UKRI Research Innovation Fellowship awarded by the Medical Research Council
745 (MR/R007446/1) and by the Cambridge Centre for Parkinson-Plus. NS is supported by the
746 Wellcome Trust, the BHF and the NIHR. JD holds a British Heart Foundation Professorship
747 and a NIHR Senior Investigator Award*. RAL is supported by grants from Parkinson's UK.
748 Sequencing of the INTERVAL samples was supported by the Wellcome Trust grant number
749 206194. Metabolon Metabolomics assays were funded by the NIHR BioResource, the
750 Wellcome Trust (grant number 206194) and the NIHR Cambridge BRC (BRC-1215-
751 20014)*. Nightingale Health NMR assays were funded by the European Commission
752 Framework Programme 7 (HEALTH-F2-2012-279233). SomaLogic assays were funded by
753 Merck and the NIHR Cambridge BRC (BRC-1215-20014)*. ICICLE-PD was funded by
754 Parkinson's UK (J-0802, G-1301, G-1507), the Lockhart Parkinson's Disease Research Fund,
755 and supported by the National Institute for Health Research (NIHR) Newcastle Biomedical
756 Research Unit and the National Institute for Health Research (NIHR) Cambridge Biomedical
757 Research Centre (146281). The EPIC-Norfolk study
758 (<https://doi.org/10.22025/2019.10.105.00004>) has received funding from the Medical
759 Research Council (MR/N003284/1 and MC-UU_12015/1) and Cancer Research UK
760 (C864/A14136). The genetics work in the EPIC-Norfolk study was funded by the Medical

761 Research Council (MC_PC_13048). Metabolite measurements in the EPIC-Norfolk study
762 were supported by the MRC Cambridge Initiative in Metabolic Science (MR/L00002/1) and
763 the Innovative Medicines Initiative Joint Undertaking under EMIF grant agreement no.
764 115372. We are grateful to all the participants who have been part of the project and to the
765 many members of the study teams at the University of Cambridge who have enabled this
766 research. The collection of samples and assay work done on the Oxford Vascular Study
767 cohort was funded by the National Institute for Health Research (NIHR) Biomedical
768 Research Centre, Oxford, and the Wellcome Trust.

769 *The views expressed are those of the authors and not necessarily those of the NHS, the
770 NIHR or the Department of Health and Social Care.

771 The authors thank all members of the ICICLE-PD Consortium for their support of this work.
772 A full list of members of the ICICLE-PD Consortium is included in the Supplementary
773 Materials.

774

775 **Author contributions:** NC, AGD, OS, NS and PFC conceived the study; NC, AGD, OS, NS
776 and PFC designed the analyses and experiments; NC performed the quantitative genetics
777 analyses; AGD performed the wet lab experiments; EYD and JMMH obtained replication
778 data; AGD, AIB, LL and ZG performed the wet lab experiments on cell lines and patient
779 cohorts; NC and CC performed the analysis on the patient cohorts; IDS and MP performed
780 the analysis on the prospective cohort; NC, AGD, OS, NS and PFC interpreted the results;
781 NS, PFC, OS, JMMH, JD, PMR, AIB, MC, RAL, CWG, NJW and CL acquired the resources
782 and datasets; NC, KK, MJB, and KW curated the data and wrote the bioinformatics methods;
783 NS, PFC, JMMH and OS supervised the project; NS and PFC acquired the funding and
784 administered the project; NC, AGD, OS, NS, and PFC wrote the paper with input from all
785 authors, who approved the final version of the manuscript. NC and AGD contributed equally
786 as first authors; EYD and KK contributed equally as second authors; JMMH, OS, NS and
787 PFC jointly supervised the work.

788

789 **Competing interests:** JMMH and EKYD are full time employees of Novo Nordisk Ltd. ASB
790 has received grants outside of this work from AstraZeneca, Biogen, BioMarin, Bioverativ,
791 Merck and Sanofi and personal fees from Novartis outside of this work. JD sits on the
792 International Cardiovascular and Metabolic Advisory Board for Novartis (since 2010); the
793 Steering Committee of UK Biobank (since 2011); the MRC International Advisory Group
794 (ING) member, London (since 2013); the MRC High Throughput Science ‘Omics Panel
795 Member, London (since 2013); the Scientific Advisory Committee for Sanofi (since 2013);
796 the International Cardiovascular and Metabolism Research and Development Portfolio
797 Committee for Novartis; and the AstraZeneca Genomics Advisory Board (2018).

798

799 **Figure Legends**

800 **Figure 1: Overview of analysis and population structure on the nuclear and**
801 **mitochondrial genomes**

802 **A.** Overview of analyses. A flowchat overview of the analyses we performed in this study,
803 summarizing the number of samples and phenotypes used in each cohort we analysed, and
804 the number of mtDNA associations with metabolites and cis-eQTLs we found. **B.** A
805 simplified mtDNA haplogroup tree with haplogroups present in INTERVAL and GTEx
806 participants individually coloured. Haplogroups not present in INTERVAL and GTEx
807 participants are coloured grey. Each haplogroup is consistently represented by the same
808 colours throughout this manuscript. **C.** A plot of principal component (PC) 1 and 2 from a
809 principal component analysis (PCA) performed using 187 mtDNA SNPs (MAF \geq 1%) in
810 16,220 participants in INTERVAL, coloured by haplogroups identified for each individual
811 using Haplogrep v2. Haplogroups U, K (Uk), and H are labelled. **D.** A plot of PC1 and PC2
812 from a PCA performed using 5,511,276 nDNA SNPs (MAF \geq 5%) in the same participants
813 in INTERVAL, coloured by their haplogroups, as previously described. **E.** A plot of PC1 and
814 PC2 from a PCA performed using 215 common mtDNA SNPs (MAF \geq 1%) in 456
815 participants in GTEx, coloured by haplogroups identified for each individual using
816 Haplogrep v2. Haplogroups U, K (Uk), and H are labelled. **F.** A plot of PC1 and PC2 from a
817 PCA performed using 5,451,305 common nucDNA SNPs (MAF \geq 5%) in the same
818 participants in GTEx, coloured by their haplogroups

819

820 **Figure 2. Metabolites and their associations with mtDNA SNPs in INTERVAL.**

821 **A.** Manhattan plot summarizing results of association between 183 common mtDNA variants
822 (MAF \geq 5%) and 896 metabolites. Each dot corresponds to the association between a mtDNA
823 variant and a metabolite. Its x coordinate represents its position along the mitochondrial
824 genome and its y coordinate represents the $-\log_{10}$ (p value) for the association from a Wald
825 Test between LMMs with and without genotypes at a mtDNA SNP as a predictor for each
826 metabolite, implemented in LDAK v5. The red dotted horizontal line at $P=1.04 \times 10^{-6}$
827 represents the significance threshold upon correcting for 896 metabolites and an estimate of
828 53.56 independent mtDNA SNPs (**Supplementary Discussion**). Red dots are associations
829 where fMet is the metabolite tested, labelled with their position and genotype. The x axis is
830 annotated with ranges of positions in the mtDNA with function, including the D-LOOP, the
831 mtDNA encoded rRNAs (MT-rRNA), the mtDNA encoded tRNAs (MT-tRNA) and the
832 mtDNA encoded protein-coding genes. **B.** Table of association statistics between 15 mtDNA
833 SNPs associated with fMet in the discovery cohort INTERVAL and replication cohort EPIC-
834 Norfolk; for each mtDNA SNP we show its position (BP), the gene it is in (GENE), the allele
835 whose effects we test (A1), the other allele (A0), the frequency of the tested allele (A1FREQ)
836 and functional annotations of the variant (ANNOT), the standardised effect size (BETA) of
837 its association with fMet, its standard error (SE) and P values (P). Association statistics from
838 both Wald tests on LMM and linear regressions in CMS are shown for the discovery cohort,
839 while association statistics from linear regression are shown for the replication cohort. **C.**
840 Pearson correlation coefficient r^2 between the mtDNA SNPs significantly associated with
841 fMet; the red, light blue and dark blue squares denote the most significantly associated
842 variants at each of the three independent signals. **D.** The haplogroup lineage tree for super-
843 haplogroup Uk (on the left) and H (on the right). SNPs that are part of this tree and
844 significantly associated with fMet are in bold and coloured according to their haplogroups.
845 All except one SNP (mt.3992C>T) associated with fMet are on the branch for the super-
846 haplogroup Uk.

847 **Figure 3. fMet-associated genes regulate mtDNA gene expression.**
848 **A.** Manhattan plot of association between common nDNA SNPs ($MAF \geq 5\%$) with fMet; the
849 x coordinate represents positions for each nDNA SNP tested; the y coordinate represents the -
850 $\log_{10}(P)$ for the Wald Test associations for nDNA SNP effects on fMet levels; red horizontal
851 line indicates the significance threshold of $P = 5 \times 10^{-8}$; red dots represent SNPs with
852 significant associations, orange SNPs represent SNPs with $P < 5 \times 10^{-7}$. **B.** Boxplot of fMet
853 levels in INTERVAL participants with each genotype at rs550045 and mt.1811. Centre of the
854 boxplots show median fMet levels, upper and lower limits of boxplots show interquartile
855 ranges, while the whiskers show values within 1.5 times the interquartile range. Outliers
856 show values beyond 1.5 times the interquartile range. **C.** The variance decomposition model
857 for quantifying the variance in fMet levels explained by variation in nDNA SNPs, mtDNA
858 SNPs, other metabolites and blood cell counts, and those that cannot be accounted for by all
859 the above. The pie chart shows the relative contribution from all four components. **D.** A
860 Manhattan plot of association between 13 protein-coding mtDNA genes and their top eQTL
861 on the mtDNA for 41 GTEx tissues. X axis represents the position of the SNP along the
862 mtDNA and Y axis represents the $-\log_{10}(P)$ from log-likelihood ratio (LRT) tests for mtDNA
863 SNP effects on expression of mtDNA encoded genes. Dots are coloured by genes, and the x
864 axis is annotated with ranges of positions in the mtDNA with function. **E.** This figure shows
865 the standardized effect size (BETA) of the mtDNA SNPs on $\log_{10}(TPM+1)$ expression levels
866 of mtDNA encoded *MT-ND3* in 41 GTEx tissues. Colour of the boxes corresponds to the
867 tissue type; centre of the boxplots show median BETA values from all mtDNA SNPs; upper
868 and lower limits show the interquartile range, while the whiskers show values within 1.5
869 times the interquartile range. Outliers show values beyond 1.5 times the interquartile range.
870 The upper panel shows the results for the 15 fMet-associated SNPs, while the bottom panel
871 shows the results for the non-fMet associated mtDNA SNPs.

872
873 **Figure 4: fMet regulates mitochondrial protein synthesis and oxidative phosphorylation**
874 **function.**

875 **A.** Schematic representation of transmitochondrial cybrids. Black, red and blue dots represent
876 the absence of mtDNA, haplogroup H and haplogroup Uk respectively. **B.** Quantification of
877 fMet levels in cybrids of different haplogroups. Statistical testing was performed by unpaired
878 t-test. Normality was assessed using the Kolmogorov–Smirnov test. The average raw value of
879 fMet is 3.06 pg/mg. **C.** Schematic representation of mitochondrial protein synthesis and fMet.
880 **D.** Effect of fMet on mitochondrial translation. Electrophoretic patterns of the synthesized
881 mitochondrial products and fragments of the gel stained with Coomassie (used as a loading
882 control); molecular weight marker (left) and each mitochondrial protein (right) are shown.
883 Quantification of the bands was corrected by the loading control in each cell line. **E.** Effect of
884 fMet on the levels of the mitochondrial transcript *MT-CO3*. **F.** One-dimensional Blue Native
885 Gel Electrophoresis (1D-BNGE) and Western blot quantification (see **Extended Data Figure**
886 **3F**); values are corrected with relative levels between loading control (CII) in each cell line
887 and untreated samples from haplogroup H. **G.** Complex I (left gel) and IV (right gel) *in gel*
888 activity assays (IGA) after 1D-BNGE analysis of digitonin treated cybrid cell lines with and
889 without fMet treatment. Super-complexes (SC) composition is indicated. **H.** Basal
890 respiration. **I.** Glycolytic ATP levels. The average raw value of ATP of all the cybrids is
891 19367878,79 luminescence units /mg of protein. **J.** Cytoplasmic ROS levels. The average
892 raw value of ROS of all the cybrids is 17434,4/20000 cells fluorescence units. In all the
893 graphs (**B-I**) bars/lines represent the mean \pm SD of the biological replicates ($n=4$) of -
894 (Control) and + (fMet treated) cell lines of each haplogroup that were measured in 3-5
895 independent technical replicates each. Colors red and blue represent haplogroup H and Uk
896 respectively. The values are represented as relative to the average of untreated samples from

897 haplogroup H, unless otherwise indicated. Statistical testing was performed by using a 2-
898 way-ANOVA test followed by Holm-Sidak's multiple comparison test unless stated
899 otherwise. *P*-values corrected for multiple comparisons are indicated. Unprocessed S35 Blots
900 and loadings can be found in **Source Data Figure 4**.

901 **Figure 5. fMet modulates cytosolic protein homeostasis.**

902 **A.** Effect of fMet on cytosolic translation products. Electrophoretic patterns of the
903 synthesized proteins and fragments of the gel stained with Coomassie (loading control) and
904 molecular weight marker (left) are shown. Quantification of the bands was corrected by the
905 loading control in each cell line. **B.** Effect of fMet and mitochondrial haplogroup on EIF2A
906 activation (**Extended Data Figure 4A**). Quantification of the immune detected bands for
907 p.EIF2A^{Ser51} and EIF2A corrected by loading control (B-actin) in each cell line. Activation of
908 EIF2A is calculated as ratio p.EIF2A^{Ser51}/ EIF2A. Values are represented as relative to the
909 average of untreated samples from haplogroup H. **C.** Effect of fMet and haplogroup on the
910 expression of EIF2A downstream targets ATF4 and CHOP. Box plots represent minimum,
911 maximum, sample median, and the first and third quartiles. All data points are plotted. **D.**
912 Effect of fMet on protein ubiquitination. Immunoblot detection with anti-ubiquitin and anti-
913 B-actin as a loading control for untreated (-/-), proteasome inhibition with MG132 (+, -) and
914 proteasome inhibition plus fMet (+, +). The quantification of the bands for ubiquitin smear
915 was corrected by loading control (B-actin). Statistical testing was performed with a 3-way-
916 ANOVA followed by Holm-Sidak's multiple comparisons. **E.** Consensus Pathway Analysis
917 of all the differentially expressed genes between haplogroup H and Uk in all tissues in GTEx.
918 Grey represents Reactome pathways and green wikipathways. **F.** Gene Set Enrichment
919 Analysis (GSEA) analysis of all the differentially expressed genes between haplogroup H and
920 Uk in all tissues in GTEx. Enrichment score is shown. **G.** Schematic representations of the
921 effects of fMet. In all the plots (**B-D**) bars/lines represent the mean \pm SD of the biological
922 replicates (n= 4) of - (Control) and + (fMet treated) cell lines that were performed in 3-5
923 independent technical replicates each. Colors red and blue represent haplogroup H and Uk
924 respectively. The values are represented as relative to the average of untreated samples from
925 haplogroup H, unless indicated. Statistical testing was performed with a 2-way-ANOVA
926 followed by Holm-Sidak's multiple comparisons, unless otherwise stated. *P*-values corrected
927 for multiple comparisons are shown. Unprocessed blots and loadings can be found in **Source**
928 **Data Figure 5**.

929
930 **Figure 6: fMet as a biomarker for IS and other late-onset disorders.**

931 **A.** Rank normalised residuals levels of fMet in IS and controls, separated by their mtDNA
932 haplogroups, after regressing out batch for quantification of fMet and site of data collection
933 as covariates; *P*-values are from logistic regression of IS disease status with haplogroup. The
934 centre of the boxplots show the median normalised fMet levels; upper and lower limits of
935 boxplots show the interquartile range, while the whiskers show values within 1.5 times the
936 interquartile range. Outliers beyond whiskers show values beyond 1.5 times the interquartile
937 range. All data points are plotted. **B.** Hazard ratio of fMet levels measured at baseline for 24
938 health outcomes, including mortality, over a 20-year follow-up. Points and error bars shown
939 in blue represent the point estimates and 95% confidence intervals of the hazard ratio in
940 11,966 EPIC-Norfolk participants. **C.** Relationship between Z scores of fMet measured at
941 baseline of a 20-year longitudinal study with 11,966 individuals, and their ages at baseline
942 (mean 60 years, SD 6 years). Significant Spearman correlations were found between fMet
943 levels and age at measurement in participants of both mtDNA haplogroup Uk and other
944 haplogroups (Spearman correlation $P = 3.51 \times 10^{-25}$ and 1.49×10^{-284} respectively). **D.** Hazard
945 ratio of fMet levels measured at baseline for 24 health outcomes, including mortality, over a

946 20-year follow-up. Points and error bars shown in blue represent the point estimates and 95%
947 confidence intervals of the hazard ratio in participants of mtDNA haplogroup Uk, while those
948 in red represent those in participants of other haplogroups.

949

950

951 **Online Methods**

952

953 **Sample filtering in INTERVAL cohort**

954

955 In the INTERVAL dataset³⁸, 12,395 participants were sequenced across the whole genome
956 (WGS, mean coverage = 26.8x, SD = 3.1x; mtDNA mean coverage = 2022.6x, SD = 566.5x,
957 **Supplementary Figure 1A**), and 4,502 participants were sequenced across the whole exome
958 (WES, mean coverage = 48.0x, SD = 8.6x; mtDNA mean coverage = 30.6x, SD = 13.5x,
959 **Supplementary Figure 1B**), including 60 participants on whom WGS was also performed.
960 Of the 12,395 participants with WGS in the INTERVAL cohort, we identified 32 participants
961 who were sequenced in duplicate. All 32 participants had their blood samples taken and
962 sequenced at two time points, so we removed one sample per participant sequenced at the
963 later of the two time points. Second, 12,112 out of the remaining 12,363 participants with
964 WGS and 4,471 out of 4,502 participants with WES can be linked to both phenotypic data
965 and genotypes assayed with the UK Biobank Affymetrix Axiom array^{75,76}. We retain only
966 these participants for further analysis. Third, we compared genotype calls from WES against
967 those from WGS in 56 overlapping participants at 307 overlapping polymorphic mtDNA
968 SNPs between WES and WGS. One individual showed a particularly high rate of discordance
969 between WES and WGS genotype calls (14 sites discordant), and was hence removed from
970 further analysis. In the remaining 55 samples with both WGS and WES, we obtained a mean
971 per variant Pearson r^2 of 0.994 (SD = 0.039) between genotypes at all 307 mtDNA SNPs
972 called in both WGS and WES, and a mean per sample Pearson r^2 of 0.980 (SD = 0.134). We
973 retained all 55 participants for further evaluation of the genotype qualities of variants called
974 from WGS and WES, giving us 12,111 participants with WGS and 4,470 participants with
975 WES (55 of whom also have WGS).

976

977 **Genotype quality control in INTERVAL cohort**

978

979 For mtDNA, we extracted reads mapping to the rCRS mitochondrial reference genome
980 (NC_012920) from WGS in 10,704 individuals and WES in 4,502 individuals in
981 INTERVAL, and called mtDNA variants using GATK v4³⁹, obtaining 396 high quality
982 variants, of which 187 have MAF \geq 1%. We use the 187 variants for assigning Haplogroups
983 to each individual in INTERVAL using Haplogrep v2 and all further analyses. For nDNA,
984 we obtained imputed genotypes from array genotypes at 5,511,276 autosomal, biallelic SNPs
985 in 43,059 unrelated individuals from European descent from the INTERVAL project⁷⁵,
986 filtering raw imputation results with information score (INFO > 0.9), minor allele frequency
987 (MAF > 5%), P-value of violation of Hardy Weinberg Equilibrium (HWE > 10⁻⁶), and
988 missingness (< 0.1).

989

990 **mtDNA variant calling from sequencing data**

991

992 To utilize the maximum number of samples for association testing, we called mtDNA
993 variations in both WGS and WES using GATK HaplotypeCaller v4. Using --ploidy 1 in
994 GATK HaplotypeCaller v4^{77,78}, we called 4,696 variants from WGS (of which 4,602 are
995 SNPs and 255 are present on the UK Biobank Affymetrix Axiom array, **Supplementary**
996 **Figure 1C**), and 3,618 variants from WES (of which 3,546 are SNPs and 254 are present on
997 the UK Biobank Affymetrix Axiom array, **Supplementary Figure 1D**), giving a union of
998 5,247 variants (of which 5,161 are SNPs). We then performed the following checks and
999 filters for sample and mtDNA variants. As mtDNA coverage is extremely high on the WGS
1000 (mean coverage = 2015.2x, SD = 578.7x), homoplasmic variations in the mtDNA would be

1001 supported by thousands of reads and are therefore of high confidence. This confidence
1002 however cannot be extended to WES, due to the low coverage (mean coverage = 30.6, SD =
1003 13.5x) and off-target nature of mtDNA reads on the WES (no mtDNA probes are present on
1004 the Agilent SureSelect Human All Exon v.5 kit).

1005

1006 **mtDNA heteroplasmy check with mtdna-server**

1007

1008 We checked the heteroplasmy levels at all WGS variant calls using the local version of
1009 mtdna-server (v1.1.11), a specialized software for variant calling in mtDNA that is
1010 particularly optimized for identification of heteroplasmic mutations in the mtDNA⁷⁹, so as to
1011 ensure the post-VQSR variants we obtained from WGS are not likely due to misidentification
1012 of heteroplasmic mutations as inherited homoplasmic variants. We were unable to perform
1013 this check directly on the WES data because mtdna-server was able to call only variations at
1014 only coverage of 30x and above - as the average coverage on the mtDNA in WES is 30.6x,
1015 we were only able to call 2,341 variants from 540 out of 4,502 samples, much fewer than
1016 present in the WGS samples. We first checked for correlation between genotypes called with
1017 GATK HaplotypeCaller and mtdna-server at 4,544 SNPs called with both methods. While
1018 4,407 SNPs showed high correlation between both methods (Pearson $r^2 \geq 0.9$), 137 SNPs
1019 showed lower correlation. We checked the level of heteroplasmy at all levels of correlation
1020 and at all minor allele frequencies (MAF) determined with the mtdna-server genotype calls,
1021 and found that levels of heteroplasmy are higher for SNPs with lower correlation between
1022 GATK HaplotypeCaller and mtdna-server at all levels of MAF (**Supplementary Figure 1E**),
1023 and SNPs with high correlation between the two methods have lower heteroplasmy (mean
1024 heteroplasmy = 5.54×10^{-5} , SD = 2.9×10^{-4}) than those with low correlation (mean
1025 heteroplasmy = 0.018, SD = 0.08, **Supplementary Figure 1F**). This suggests that GATK
1026 may have mis-called these heteroplasmic sites as homoplasmic variants. We removed these
1027 137 SNPs from both WGS and WES variant calls from all further analyses.

1028 **Sequencing-based mtDNA quality**

1029

1030 Using 12,111 and 4,470 samples with WGS and WES, we compared genotypes at biallelic,
1031 polymorphic, and non-strand-ambiguous mtDNA SNPs called in WGS (209 SNPs) and WES
1032 (206 SNPs) against genotypes called in the Affymetrix Axiom array (out of a total of 235
1033 biallelic, non strand-ambiguous SNPs genotyped on the array). We find that WES showed
1034 similar correlation and concordance with array genotypes (mean Pearson $r^2 = 0.964$, SD =
1035 0.131, mean concordance = 0.999, SD = 0.004, 5 sites with more than 1% participants
1036 discordant) as WGS (mean Pearson $r^2 = 0.961$, SD = 0.131, mean concordance = 0.999, SD =
1037 0.004, with the same 5 sites with more than 1% participants discordant as WES). Both
1038 analyses showed that WES can produce high quality genotypes and variant calls in the
1039 mtDNA, despite having two orders of magnitude lower coverage on the mtDNA than the
1040 WGS. In addition, we found that the same sites were discordant between Affymetrix and both
1041 sequencing platforms: 13 sites had Pearson r^2 of lower than 0.9 between Affymetrix and both
1042 WGS and WES, inclusive of the 5 that had greater than 1% discordant rate between
1043 Affymetrix and both sequencing platforms. This indicates that these 13 sites (magenta points
1044 in **Supplementary Figure 2**) likely represent errors in the Affymetrix array rather than either
1045 sequencing platforms.

1046 **mtDNA variant quality recalibration**

1047

1048 To assess the quality of variant and genotyping calling at sites that are polymorphic in all
1049 participants in the WES and WGS cohorts, we used Gaussian mixture models in Variant
1050 Quality Score Recalibration (VQSR) in GATK (version 4.0.3.0)⁸⁰ to cluster all SNP calls
1051 from WES with SNPs of high concordance between the three platforms based on their variant
1052 call metrics. We first identified 92 SNPs of high concordance between all three platforms
1053 (Pearson $r^2 > 0.9$), and designated them as the "training" set for training the Gaussian mixture
1054 model (orange points in **Supplementary Figure 2**). We then restricted our "known" set to a)
1055 the 92 SNPs in the training set, b) the 189 SNPs genotyped on the array that showed good
1056 correlation between WES and array (Pearson $r^2 > 0.9$), and c) the 231 common SNPs with
1057 high minor allele frequency (MAF > 1%) in WGS, WES or array calls. From c), we removed
1058 10 SNPs with low Affymetrix genotype quality (Pearson r^2 with WGS and WES < 0.9), 9
1059 SNPs with low WES quality (Pearson r^2 with WGS < 0.9), and 1 SNP with low WGS quality
1060 (Pearson r^2 with Affymetrix < 0.9). We took the union of the remaining 211 SNPs from c)
1061 with a) and b) to arrive at 314 SNPs to use as the "known" set. The correlation between
1062 platforms and MAF of all SNPs from the three platforms is shown in **Supplementary Figure**
1063 **2**. We then applied VQSR separately on the WES callset with "--trust-all-polymorphic" and
1064 "--max-gaussians 2" in SNP mode using the following annotations: QD, FS, MQ,
1065 MQRankSum, ReadPosRankSum, BaseQRankSum, SOR, and MLEAF. We obtained
1066 VQSLOD scores for all variants based on the clustering of their annotations with "known"
1067 and "training" sets, and selected variants not in the above sets whose transition to
1068 transversion ratio (Ti/Tv) most closely matched that of "known" variants in WES (known
1069 Ti/Tv = 63.3). After filtering WES variant calls with VQSLOD scores of lower than 1.91
1070 (TruthSensitivity = 50%, Ti/Tv = 59.0), we obtained 212 high-quality SNPs from WES. We
1071 took the union of this set with a) 92 polymorphic SNPs from the "training" set and b) 189
1072 SNPs with high correlation between WES and array, to obtain a total of 403 high-quality
1073 SNPs from the WES. 7 SNPs were strand-ambiguous or multi-allelic (mt.373, mt.1766,
1074 mt.3308, mt.7960, mt.13816, mt.14605, mt.15625), and were hence removed, leaving 396
1075 high-quality SNPs, all of which were called from WGS. Of this, 187 are common at MAF >
1076 1%, and we used these SNPs for association analyses as well as quantification of the total
1077 mtDNA contribution to molecular phenotypes.

1078 **Metabolite data quality control**

1079
1080 Metabolome profiling was performed in two batches on plasma samples extracted from the
1081 blood of 9000 INTERVAL participants using the Metabolon HD4 mass spectroscopy
1082 discovery platform. This platform quantifies plasma metabolites using the Ultrahigh
1083 Performance Liquid Chromatography-Tandem Mass Spectroscopy (UPLC-MS/MS) method,
1084 and produces ion-counts for specific fragments that identify specific metabolites. Raw data
1085 was extracted, peak-identified and QC processed using Metabolon's hardware and software.
1086 As such, the raw data from this platform corresponds to but is not a direct measurement of
1087 plasma metabolite concentrations. Where metabolite levels were below the lower limit of
1088 detection, they were set to "missing" rather than 0 or the lowest detectable value, in order to
1089 prevent skewing of the data. The following steps were carried out for the filtering of samples
1090 and metabolites for ensuring only high-quality metabolite quantification was used in all
1091 subsequent analyses. First, only 7,778 out of 9,000 INTERVAL participants who were
1092 previously found to be unrelated and of European ancestry and had not withdrawn from the
1093 INTERVAL study are included in this study. Second, 68 participants had metabolite data
1094 quantified from repeat blood samples (taken either at the same or different times). For these
1095 participants, metabolite data quantified from the first blood samples (baseline) were kept,
1096 where available, and repeats were dropped. Where samples were from the same survey time,

1097 both samples were dropped. Third, metabolites which were measured in only one batch (n=22)
1098 were excluded from the analysis. Fourth, participants with missing data for at least 300
1099 metabolites (n=9) were excluded (arbitrary cut-point decided based on histogram of
1100 missingness) from further analysis. Quantification measures of all remaining 995 metabolites
1101 were then transformed by taking the natural logarithm, and then winsorized where the value
1102 was 5 or more standard deviations away from the mean metabolite value. The transformed
1103 metabolites were then regressed (linear regression) against the following covariates: age, sex,
1104 batch, INTERVAL recruitment centre, plate number, appointment month, the lag time
1105 between the blood donation appointment and sample processing, and the first 5 ancestry
1106 principal components obtained from genome-wide genotyping data from the Affymetrix UK
1107 Biobank Axiom array⁷⁵. Following regression, the residuals were generated and inverse rank
1108 normalized. As metabolite levels below detection limits set to “missing” did not factor into
1109 any of the above transformations and remained “missing”, we finally removed all metabolites
1110 that were missing in more than 50% (4000) of the samples, leaving 896 metabolites for
1111 association analysis with mtDNA variants.

1112

1113 **Partitioning of contributions to metabolite levels**

1114

1115 We tested for single mtDNA or nDNA SNP association with levels of each metabolite in
1116 LDAK, which implements a linear mixed model: for each of the common mtDNA or nDNA
1117 variants x_j of $MAF \geq 5\%$ in INTERVAL, we tested for their effect β on metabolite level
1118 y_i in a linear mixed model, controlling for all autosomal genetic effects and population
1119 structure in the random effect term g mean metabolite level (using an intercept term) in the
1120 fixed effect term F_i : $y_i = \alpha F_i + \beta x_j + g + \phi$ where $g \sim N(0, \sigma_g^2 K_g)$ and $\phi \sim N(0, \sigma_e^2 I)$,
1121 g is the random effect of the LD-weighted relatedness matrix K_g constructed using
1122 LDAK v5⁴² with 5,511,276 common ($MAF \geq 5\%$) autosomal, biallelic SNPs nuclear DNA
1123 variants using --power -0.25 as recommended, and ϕ is the residual variance assuming an
1124 independent identically distribution (i.i.d) matrix for noise and uncaptured environmental
1125 effects. For each mtDNA SNP-metabolite association, we asked if we could increase the
1126 power to detect true associations by controlling environmental factors indexed by levels of
1127 selected metabolites or blood cell counts using the CMS framework⁴³ (**Supplementary**
1128 **Discussion**). To quantify relative levels of contribution of autosomal, mtDNA, and known
1129 environmental contribution to each metabolite, we performed a variance decomposition
1130 analysis with LDAK v5, using combinations of the following variance components: i) LD-
1131 weighted and MAF-adjusted (with --power 0.25 as recommended) relatedness matrix from
1132 5,511,276 common SNPs on the autosomes only, ii) non-weighted relatedness matrix from
1133 187 common SNPs on the mtDNA only, and iii) non-relatedness matrix from levels of 896
1134 metabolites and 36 blood cell types, and assessed significance of contribution from each of
1135 the three variance components using likelihood ratio tests.

1136

1137 **GTEEx sample selection and genotype quality control**

1138

1139 We obtained 5,696,458 autosomal, biallelic SNPs in 635 GTEEx samples from the variant call
1140 set from WGS data in version 7³⁷, filtering on minor allele frequency ($MAF \geq 5\%$), P-value
1141 of violation of Hardy Weinberg Equilibrium ($HWE > 10^{-6}$), and missingness (< 0.1). We
1142 then extracted reads mapping to the rCRS mitochondrial reference genome (NC_012920),
1143 and called mtDNA variants using mtdna-server⁷⁹, obtaining 1,756 variants, all of which are
1144 SNPs, and of which 46 are multi-allelic. For multi-allelic SNPs, we retained the two alleles
1145 with highest frequencies for analysis, hence obtaining a total of 1,714 SNPs, 64 of which are
1146 common ($MAF \geq 5\%$) for use in association testing. To identify unrelated individuals, we

1147 used KING9 to identify related samples among the 635 GTEx samples. Two pairs of
1148 individuals in GTEx are related up to third-degree (kinship ≥ 0.04419), though only
1149 marginally (kinship among pairs = 0.0477 and 0.0657 respectively), so we did not remove
1150 them from analyses. To identify individuals of European ancestry, we selected 4,812,475
1151 common SNPs (MAF > 5%, P value for HWE > 10^{-6}) from the autosomes that overlap
1152 between GTEx and 1000 Genomes Project Phase 3 (1000G)⁸¹, built an LD-weighted kinship
1153 matrix using LDAK v5 and performed PCA on 1000G samples, projected the GTEx samples
1154 onto PCs from 1000G samples, and selected 491 GTEx samples that cluster with 1000G
1155 samples from the EUR superpopulation. We then built a kinship matrix using 5,696,456
1156 common SNPs (MAF $\geq 5\%$, missingness < 0.1, P value for HWE > 10^{-6}) in 491 European
1157 samples in GTEx for use in testing for association between mtDNA encoded genes and
1158 mtDNA SNPs with a linear mixed model.

1159 **Population structure on nuclear DNA and mtDNA in GTEx**

1160
1161 We obtained 5,696,456 autosomal, biallelic SNPs in 635 GTEx participants from the variant
1162 call set from WGS data in version 7³⁷, filtering on minor allele frequency (MAF > 5%), P-
1163 value of violation of Hardy Weinberg Equilibrium (HWE > 10^{-6}), and missingness (< 0.1).
1164 To identify unrelated participants, we used KING⁸² to identify related participants among the
1165 635 GTEx participants. Two pairs of participants in GTEx are related up to third-degree
1166 (kinship ≥ 0.04419), though only marginally (kinship among pairs = 0.0477 and 0.0657
1167 respectively), so we did not remove them from our analyses. To identify participants of
1168 European ancestry, we selected 4,812,475 common SNPs (MAF > 5%, P value for HWE >
1169 10^{-6}) from the autosomes that overlap between GTEx and 1000 Genomes Project Phase 3
1170 (1000G)⁸¹ and built an LD-weighted kinship matrix using LDAK v5⁴², and projected the
1171 GTEx participants onto PCs obtained from 1000G samples using the same SNPs
1172 (**Supplementary Figure 5A**). We selected 491 GTEx participants that cluster with 1000G
1173 participants from the EUR superpopulation (**Supplementary Figure 5B**). Of these 491
1174 participants, we were able to obtain WGS reads on the mtDNA of 456 participants. We then
1175 extracted reads mapping to the rCRS mitochondrial reference genome (NC_012920) from the
1176 456 participants, and called mtDNA variants using both GATK HaplotypeCaller, using the
1177 same settings as we did in INTERVAL, as well as mtdna-server (v1.1.11)⁷⁹ to check for
1178 heteroplasmy. We obtained a total of 1,180 SNPs. Of the 1,180 SNPs, 38 are likely mis-
1179 identified as inherited homoplasmic SNPs due to heteroplasmy (mean heteroplasmy = 0.23,
1180 SD = 0.35, **Supplementary Figure 5C,D**). Of these, 12 are also found to be potential
1181 heteroplasmic sites in INTERVAL, indicating that these sites are consistently mis-identified
1182 as homoplasmic SNPs from WGS data due to heteroplasmy across studies. We removed all
1183 38 potentially heteroplasmic sites, leaving us with a total of 1,142 SNPs (mean heteroplasmy
1184 = 2.48×10^{-4} , SD = 9.90×10^{-4} , **Supplementary Figure 5D**). Of these, we use the 56 SNPs
1185 that are common (MAF > 5%), for use in association testing. To assess the diversity on
1186 nuclear DNA and mtDNA in GTEx in these 456 participants, we built a LD-weighted kinship
1187 matrix using LDAK v5⁴² with 5,451,305 common SNPs (MAF > 5%, missingness < 0.1, P
1188 value for HWE > 10^{-6} in just the EUR participants) to obtain PCs specifically in these
1189 participants to compare against PCs obtained from the mtDNA (**Figure 1E,F**).

1190 **Obtaining PEER factors as covariates for eQTL analysis**

1191
1192 Inaccuracies in quantification of expression levels of some nDNA genes attributable to mis-
1193 mapping of 100bp RNAseq reads originating from mtDNA encoded genes to the nuclear

1194 mitochondrial sequence (NUMT) regions of the nucDNA, sequences on the nDNA that are
1195 highly similar to the mtDNA. As such, we aligned the sequence of each nuclear gene (57,820
1196 genes) in GENCODE v19 on the rCRS mtDNA reference sequence NC012920 using lastal
1197 (version 744)⁸³ and found 651 genes with extensive sequence similarity (≥ 100 bp) where
1198 the total fraction of genes with such alignment $\geq 5\%$ (**Supplementary Table 6**). Of these,
1199 84.0% are pseudogenes, 5.68% are lincRNAs, 2.30% are antisense RNAs, and 7.93% are
1200 protein coding genes, of which 0.1% (of total) are in introns. We excluded all of them in the
1201 calculation of PEER factors⁸⁴ for capturing unknown confounding factors in the RNAseq
1202 data and the association analysis. The numbers of PEER factors we used for correction of
1203 gene expression levels in each tissue increased with the number of participants with gene
1204 expression data in each tissue, following suggestions from GTEx release 7³⁷, and are shown
1205 in **Supplementary Table 5**.

1206

1207 **GTEx mtDNA eQTLs and multiple testing correction**

1208

1209 We tested for single mtDNA SNP association with nDNA gene expression levels for every
1210 gene i (mt-eQTL) using LIMIX⁸⁵, a linear mixed model package in python. For each of the
1211 common mtDNA variants x_j of MAF $> 5\%$ in the GTEx participants, we tested for its effect β
1212 on PEER-factor⁸⁴ corrected gene expression y_i in a linear mixed model, controlling for all
1213 autosomal genetic effects and population structure in the random effect term g , and age, sex
1214 and mean gene expression level (using an intercept term) in the fixed effect term F_i : $y_i = \alpha F_i +$
1215 $\beta x_j + g + \varphi$, where $g \sim N(0, \sigma_g^2 K_g)$ and $\varphi \sim N(0, \sigma_e^2 I)$, g is the random effect of the LD-
1216 weighted relatedness matrix K_g constructed using LDAK with 5,696,456 common (MAF \geq
1217 5%) autosomal, biallelic SNPs nDNA variants, and φ is the residual variance assuming an
1218 independent identically distributed (i.i.d) matrix for noise and uncaptured environmental
1219 effects. We obtained nominal p-values for each variant-gene pair by testing the alternative
1220 hypothesis that the β between genotype and expression deviates from 0. We then calculated
1221 Beta distribution-adjusted (using the beta distribution model of the minimum P value
1222 distribution) empirical P values for the top cis-eQTL per gene using P values generated from
1223 100 permutations of SNPs. These empirical P values were used to calculate Q values, the
1224 false discovery rate (FDR), using the “qvalue” package in R. A FDR threshold of 0.05 was
1225 applied to identify genes with a significant eQTL (“eGenes”), and the maximum empirical P
1226 value with FDR smaller than or equal to 0.05 was the gene-level threshold for identifying
1227 significant eQTLs (equivalent to empirical P value < 0.004).

1228 **Differential expression and pathway enrichment**

1229

1230 We downloaded counts of reads mapping onto each gene quantified in each tissue in GTEx
1231 from the GTEx portal (version 2016-01-15_v7_RNASeQCv1.1.8), and performed
1232 normalization of the read counts by effective library size using calcNormFactors with the
1233 default trimmed-mean of M values (TMM) method using R package “edgeR”^{62,63}. We then
1234 estimated the common, trended and tagwise dispersions over all genes using estimateDisp in
1235 edgeR, before computing the log(fold change), P values and false discovery rate (FDR) for
1236 differential expression in each gene between donors with mtDNA haplogroup H and
1237 haplogroup Uk using the exactTest function in edgeR. Pathway and GSEA analysis was
1238 performed using the WEB-based GEne SeT AnaLysis Toolkit^{86,87} following their instructions
1239 online. Results are shown in **Supplementary Data 12,13**.

1240 **mtDNA sequencing in cytoplasmic hybrid cell lines**

1241

1242 We verified the mtDNA sequence from the cybrid cell lines. mtDNA from cybrid cell lines
1243 was enriched using long-range PCR. To eliminate the potential for error and nDNA
1244 contamination (nuclear-mitochondrial sequences, NUMTs), amplicons were polymerised
1245 using PrimeSTAR GXL DNA polymerase (error rate = 0.00108 %, Takara Bio, Saint-
1246 Germain-en-Laye, France) in two overlapping fragments, using primer set-1: CCC TCT CTC
1247 CTA CTC CTG-F (m.6222-6239) and CAG GTG GTC AAG TAT TTA TGG-R (m.16133-
1248 16153), and set-2: CAT CTT GCC CTT CAT TAT TGC-F (m.15295-15315) and GGC AGG
1249 ATA GTT CAG ACG-R (7773-7791). Primer efficiency and specificity was assessed as
1250 successful after no amplification of DNA from rho⁰ cell lines, minimising the unintended
1251 amplification of nuclear pseudogenes. Amplified products were assessed by gel
1252 electrophoresis, against DNA+ve and DNA-ve controls, and quantified using a Qubit 2.0
1253 fluorimeter (Life Technologies, Paisley, UK). Each amplicon was individually purified using
1254 Agencourt AMPure XP beads (Beckman-Coulter, USA), pooled in equimolar concentrations
1255 and re-quantified. For the mtDNA sequencing pooled amplicons were ‘tagmented’, amplified,
1256 cleaned, normalised and pooled into 48 sample multiplexes using the Illumina Nextera XT
1257 DNA sample preparation kit (Illumina, CA, USA). Multiplex pools were sequenced using
1258 MiSeq Reagent Kit v3.0 (Illumina, CA, USA) in paired-end, 250 bp reads. Post run data,
1259 limited to reads with QV >= 30, were exported for analysis. Post-run FASTQ files were
1260 analysed using an in-house developed bioinformatic pipeline. Reads were aligned to the
1261 rCRS (NC_012920) using BWA v0.7.10, invoking `-mem`⁸⁸. Aligned reads were sorted and
1262 indexed using Samtools v0.1.18⁸⁹, duplicate reads were removed using Picard v1.85
1263 (<http://broadinstitute.github.io/picard/>). Variant calling (including somatic calling) was
1264 performed in tandem using VarScan v2.3.8^{90,91} (minimum depth = 1,500, supporting reads =
1265 10, base-quality (BQ) => 30, mapping quality (MQ) => 20 and variant threshold = 1.0 %) and
1266 LoFreq v0.6.1⁹². Concordance calling between VarScan and LoFreq was > 99.5%.
1267 Concordant variants were annotated using ANNOVAR v529⁹³. In-house Perl scripts were
1268 used to extract base/read quality data and coverage data. The mtDNA haplogroup was
1269 determined through in-house algorithms based upon existing phylogenetic data^{94,95} and
1270 through Haplogrep⁹⁶. The pathogenicity Score was given to each of the SNPs as previously
1271 described^{97,98}.

1272 **Cell lines and formyl-methionine treatment**

1273
1274 Cell lines were grown in Dulbecco’s modified eagle medium (DMEM) containing glucose
1275 (4.5 g/l), pyruvate (0.11 g/l) and fetal bovine serum (FBS) (5 %) without supplemented fMet
1276 and/or antibiotics at 37°C and 5 % CO₂ conditions. 8 cell lines H (4 lines), Uk (4 lines) from
1277 8 independent healthy control subjects were used. All the cybrids were obtained from cybrid
1278 pools after the selection process¹⁰. All the experiments were performed in at least 3 cell lines
1279 derived from different donors per haplogroup. The mtDNA sequences of all the cell lines can
1280 be found on GenBank and their mtDNA accession numbers¹⁰ are included in **Supplementary**
1281 **Table 9**. In fMet supplementation experiments, cells were incubated with 1ng /ml of N-
1282 Formyl-L-methionine (F3377, SigmaAldrich) during 2 days prior to the experiments.

1284 **Cell growth and doubling time experiments**

1285
1286 For growth experiments 30000 cells were seeded in a 6-well plate and counted every 4h using
1287 an Incucyte® Live-Cell Analysis system. Three to five growth curves were performed for
1288 every cell line, and each time point was counted in triplicate. Time 0h was used for correction
1289 of each well. For the doubling time analysis, the data was analyzed on an exponential curve
1290 and only those curves with R² ≥ 0.9 were considered.

1291

1292 **Formyl-methionine quantification**

1293

1294 fMet was quantified using the Formyl-methionine ELISA Kit from Elabscience (E-EL-0063)
1295 following manufacturer's conditions.

1296 **Real-time PCR quantification of transcripts.**

1297

1298 Total RNA was isolated from cells exponentially growing using an RNA isolation kit from
1299 Qiagen® according to the manufacturer's protocol. Quantification of mRNA by real-time
1300 PCR (RT-PCR) was carried out using the High capacity cDNA reverse transcription kit
1301 (Applied Biosystems) following the manufacturer's conditions. The mRNA levels were
1302 determined using probes from Applied Biosystems and following MQIE guidelines⁹⁹. The
1303 expression levels were normalized with *GADPH* and *B-ACTIN* as housekeeping genes. The
1304 codes of each of the probes are included in **Supplementary Table 17**. The comparative Cq
1305 method was used for relative quantification of gene expression. Differences in the Cq values
1306 (dCq) of the transcript of interest and the reference gene were used to determine the relative
1307 expression of the gene in each sample. The dCq method was used to calculate the number of
1308 copies.

1309 **Mitochondrial and Cytoplasmic Translation Assay**

1310

1311 Protein translation was assessed by labeling with 35S-methionine/35S-cysteine
1312 (EXPRE35S35S Protein Labeling Mix; Perkin Elmer Life Sciences) in cells seeded at 80% of
1313 confluence. For cytoplasmic translation assessment, cells were washed twice with Met and
1314 Cys-free DMEM (21013024, ThermoFisher) followed by an incubation on the same media
1315 for one hour. Then, the cells were incubated with a labelling media containing Met and Cys-
1316 free DMEM, 2 mM Glutamine (25030081, ThermoFisher), 1mM Sodium Pyruvate, 96ug/ml
1317 Cystein (DOC0122, ForMedium) and 5% dialyzed Fetal Bovine Serum (30067334,
1318 ThermoFisher) for 10 minutes at 37 °C, followed by the addition of 100uCi 35S L-
1319 Methionine and incubation for 30 minutes at 37 °C. Mitochondrial translation assay was
1320 performed similarly than in the cytoplasmic assay with some adaptations. Cells were
1321 incubated during 20 minutes at 37 °C with Labeling Medium including 100ug/ml emetine
1322 (E2375, Sigma) followed by the addition of 100μCi 35S L-Methionine and incubation for 60
1323 min at 37 °C. In both cases, cells were trypsinized and collected with PBS (and washed twice)
1324 and pellets kept at -80°C. Proteins extraction was performed using a buffer containing 0.1%
1325 DDM (D4641, SigmaAldrich), 1% Sarkosyl (L9150, SigmaAldrich) and 50 units of
1326 Benzonase (Novagen 70664 25U/ul) vortexed vigorously and left on ice for 30 minutes.
1327 Protein quantities were assessed by DC assay (Biorad 500-0113) following manufacturer's
1328 instructions. Next, 15ug of protein was loaded onto 15-well-12% Tris-Glycine gels
1329 (Invitrogen NP0343BOX) using MES buffer and ran for 3 hours at 70 volts. Total protein
1330 levels were assessed by Coomassie blue staining (0.1% Coomassie blue in 7% acetic acid and
1331 40% methanol and de-stained with (20% methanol, 7% Acetic acid) solution for 2-3h. Images
1332 of the gel were collected using a Scanner. Gels were dried at 80°C for 2 hours. Dried gels
1333 were then exposed for several days and imaged using a phosphor imaging screen on an
1334 Amersham™ Typhoon™ Biomolecular Imager. The bands were quantified, aligned and
1335 cropped using the Fiji program and the OD was used as a value for statistical purposes.

1336 **Mitochondrial bioenergetics characterization**

1337

1338 Oxygen consumption modifications. Briefly, 20×10^4 cells/well were seeded 8-12 hours
1339 before the measurement basal respiration, leaking respiration (LR), maximal respiratory
1340 capacity (MRC) and not mitochondrial respiration (NMR) were determined by adding 1 μM
1341 oligomycin (LR), 0.75 and 1.5 μM of carbonyl cyanide-p-trifluoromethoxyphenylhydrazine
1342 (FCCP) (MRC) and 1 μM rotenone/antimycin (NMR), respectively. Data were corrected by
1343 the NMR and expressed as pmol of oxygen/min/mg of protein. The quantity of protein in
1344 each well was measured by Bradford method¹⁰⁰.

1345 **Determination of MIMP and cytoplasmic and mitochondrial ROS,**

1346
1347 The determination of mitochondrial inner membrane potential (MIMP) was carried using
1348 Tetramethylrhodamine, methyl ester (TMRM) at 20nM (DMSO) in parallel to the
1349 mitochondrial mass detection using Mito-Traker Green (20nM in DMSO). Mitochondrial
1350 superoxide content was measured using MitoSOX Red at 5 μM in DMSO. Cytosolic ROS
1351 were measured using 2',7'-dichlorofluorescein-diacetate at 9 μM in DMSO. All the reagents
1352 were purchased in Invitrogen®. Fluorescence activated detection was carried using a BD
1353 LSRFortessa™ cell analyzer from BD. 20000 events were recorded and doublet
1354 Discrimination was carried using the FCS-Height and Area FlowJo Software. An example of
1355 the gating strategy is shown in **Supplementary Figure 7**. The data is expressed as intensity
1356 of fluorescence.

1357

1358 **Determination of ATP levels**

1359

1360 ATP levels were measured four times in three independent experiments using the CellTiter-
1361 Glo® Luminiscent Cell Viability Assay (Promega) according to the manufacturer's
1362 instructions. Briefly, 10,000 cells/well were seeded and the media was changed 48h before
1363 the measurement. After that time cells were lysed, and lysates were incubated with the
1364 luciferin/luciferase reagents. Samples were measured using a NovoStar MBG Labtech
1365 microplate luminometer, and the results referred to the protein quantity measured in a parallel
1366 plate.

1367 **Electrophoresis and Western blot analysis.**

1368

1369 Samples for blue-native gel electrophoresis (BNGE) and *in gel* activities were prepared as
1370 previously described^{54,101}. Native samples were run through precast NativePAGE 3–12%
1371 Bis–Tris gels during 6-9 hours. Total protein extracts were prepared according to each
1372 protein's solubilities. Mitochondrial proteins were prepared using 2% dodecyl-maltoside in
1373 PBS including protease inhibitors. Protein extracted for kinase phosphorylation analysis was
1374 extracted using PathScan® Sandwich ELISA Lysis Buffer from Cell signaling. In any case
1375 protein extracts were loaded on NuPAGE® Bis-Tris Precast Midi Protein Gels with MES
1376 (Invitrogen®) with 20 or 26 wells depending on the experiment. Electrophoresis was carried
1377 out following the manufacturer's conditions. SeeBlue® Plus2 Pre-stained Protein Standard
1378 from Invitrogen® was used in each electrophoresis as protein size markers. The separated
1379 proteins were transferred to polyvinylidene fluoride membranes using the iBLOT system
1380 (Invitrogen®) or Mini Trans-Blot® transfer system from Biorad®. The resulting blots were
1381 probed overnight at 4 °C with primary antibodies with the appropriate concentration
1382 following manufacturer's condition with small adaptations (Antibodies, and concentrations
1383 are attached in **Supplementary Table 18**). After the primary antibody, blots were incubated
1384 for 1 h with secondary antibodies conjugated with horseradish peroxidase (HRP) and
1385 immuno-detected using an Amersham Imager 600. The bands for each antibody were

1386 quantified, aligned and cropped using the Fiji program and the OD was used as a value for
1387 statistical purposes. In order to avoid inter-blot variation one cell line was used as an internal
1388 control and the values of the OD corrected by β -Actin were relative to it in each case.

1389 **Reproducibility of the experiments and statistical analysis**

1390

1391 All of the experiments present in this work were performed in 3 independent biological
1392 replicates (unless noted otherwise) and statistical analyses were derived from these data
1393 (Prism 8.0.1). Normal distributions were validated by the Kolmogorov–Smirnov test. One-
1394 way ANOVA followed by the Holm-Sidak test for multiple comparisons Kruskal-Wallis was
1395 applied for group comparison tests.

1396 **Mitochondrial disease patients harboring the m.3243A>G variant**

1397

1398 Serum samples from patients carrying the m.3243A>G variant were obtained from Prof
1399 Chinnery's neurogenetic/mitochondrial clinic at through the study: Genotype and Phenotype
1400 in Inherited Neurodegenerative Diseases (REC ID: 13/YH/0310, IRAS ID: 136697)
1401 Cambridge University Hospitals NHS Trust. Age and gender-matched controls were obtained
1402 in the NIHR BioResource and the Blood and Stem Cell Biobank (Cambridge, UK) Ethics ID:
1403 13/YH/0310.

1404 **Oxford Vascular Study**

1405

1406 OXVASC is a longitudinal population-based incidence cohort of all acute vascular events in
1407 a defined population of 92,728 people, covered by around 100 primary care physicians in
1408 nine primary care practices in Oxfordshire, UK. An estimated 97% of the true study
1409 residential population is registered with a primary care practice; most unregistered people are
1410 young students. The study area contains a mix of urban and rural populations. The OXVASC
1411 population is 94% Caucasian, 3% Asian, 2% Chinese, and 1% Afro-Caribbean. Written
1412 informed consent or assent from relatives is obtained in all participants for study, interview
1413 and follow-up, including ongoing review of primary care and hospital records and death
1414 certificate data. OXVASC was approved by the Oxfordshire research ethics committee
1415 (OREC A: 05/Q1604/70). Multiple overlapping methods are used for ascertainment of all
1416 participants with TIA and stroke, approaching 100% of events reaching medical attention.
1417 These include the following: (a) a daily, rapid access clinic to which participating general
1418 practitioners and the local emergency department refer participants with suspected TIA or
1419 minor stroke; (b) daily searches of admissions to the medical, stroke, neurology, and other
1420 relevant wards; (c) daily searches of the local emergency department attendance register; (d)
1421 daily searches of in-hospital death records via the Bereavement Office; (e) monthly searches
1422 of all death certificates and coroner's reports for out-of-hospital deaths; (f) monthly searches
1423 of general practitioner diagnostic coding and hospital discharge codes; and (g) monthly
1424 searches of all brain and vascular imaging referrals. Demographic data and stroke risk factors
1425 are collected from face-to-face interviews by study physicians as soon as possible after
1426 referral or hospital admission and cross-referenced with primary care records. Detailed
1427 clinical history was recorded in all patients and assessments were made for stroke severity
1428 using the National Institute of Health Stroke Scale (NIHSS) as recorded on assessment.
1429 Cause of ischaemic events was classified according to the Trial of Org 10172 in Acute Stroke
1430 Treatment (TOAST) criteria. Stroke and TIA were defined according to WHO criteria (acute
1431 onset of neurological deficit, persisting for >24 hours in case of a stroke, or for <24 hours in

1432 case of a TIA), with review of all cases as soon as possible after presentation by the same
1433 senior neurologist throughout the study. Non-fasting blood samples were taken as soon as
1434 possible after the event, usually within one day. These included serum, 3.2% buffered tri-
1435 sodium citrate plasma and lithium heparin plasma (Vacutainer tubes; Becton Dickinson,
1436 United Kingdom). Samples were centrifuged at 3000 g for 10 minutes, and aliquots of serum
1437 and plasma were stored at -80°C before analysis when they were thawed for use at 37°C. All
1438 times from sampling to freezing were documented, typically within 4 hours of taking.

1439

1440 **ICICLE-PD Cohort**

1441

1442 Plasma samples from PD patients and controls were obtained from the ‘Incidence of
1443 Cognitive Impairment in Cohorts with Longitudinal Evaluation-PD’ (ICICLE-PD) study,
1444 which includes newly-diagnosed PD cases and unrelated control subjects of a similar age
1445 recruited from the community and outpatient clinics in Newcastle and Cambridge, UK¹⁰².
1446 Idiopathic PD was diagnosed according to UKPDS Brain Bank criteria¹⁰³. The study was
1447 approved by the Newcastle and North Tyneside Research Ethics Committee. All patients
1448 provided written informed consent. Venous blood samples were collected in EDTA tubes at
1449 baseline study visits (between 2009 and 2011), and centrifuged within 30 minutes at 2000rpm
1450 for 15 minutes. Plasma was removed and stored in 200µl aliquots at -80C until assays were
1451 performed.

1452 **mtDNA sequencing of patients cohorts**

1453

1454 We selected 282 participants with ischaemic stroke (IS) and 181 age-matched controls from
1455 the Oxford Vascular Study (OXVASC)¹⁰⁴ for sequencing of mtDNA using the Illumina
1456 Hiseq 2000 using an amplicon-based paired-end library preparation; both groups are enriched
1457 for individuals with Haplogroup Uk (32% haplogroup Uk in cases, 23% haplogroup Uk in
1458 controls). We also sequenced 123 participants with Parkinson’s disease (PD) and 40 age-
1459 matched controls from the Incidence of Cognitive Impairment in Cohorts with Longitudinal
1460 Evaluation–PD (ICICLE-PD) cohort¹⁰² with the same platform; both groups represent
1461 population samples and are not enriched for any mtDNA haplogroups. The Fluidigm Access
1462 Array™ technology was used to generate tagged and indexed amplicons (on average 100
1463 per sample of 150-200bp), with sample-specific barcodes and Illumina adaptor sequences.
1464 The resulting PCR products were checked for quality using the Agilent 2100 Bioanalyzer and
1465 then pooled together in equal volumes. The PCR product library was purified using AMPure
1466 XP beads and quantified with PicoGreen prior to loading for Illumina sequencing. 183 age-
1467 matched controls from the OXVASC cohort were sequenced with Illumina Miseq using a
1468 paired-end library preparation. Mitochondrial DNA MiSeq libraries were prepared by
1469 amplification of two overlapping fragments¹⁰⁵. After individual purification and
1470 quantification, the amplicons from each sample were pooled in equal amounts. Libraries were
1471 prepared with NEBNext Ultra library prep reagents (New England BioLabs, MA) according
1472 to manufacturer’s instructions and sequenced using a 2 × 250-cycle MiSeq Reagent kit v3.0
1473 (Illumina, CA).

1474 **Mitochondrial variant calling and haplogroup prediction**

1475

1476 Quality of raw sequencing fastq files was checked with FastQC
1477 (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) prior to mapping and eventually
1478 trimmed using TrimGalore!

1479 (https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/) to remove low-quality
1480 read ends (--q 20), remove adapters (--stringency 5) and Ns from either side of the read (--
1481 trim-n). Trimmed reads below 35bp were removed (--length 35). Reads were mapped using the
1482 MToolBox pipeline (v.1.1) ¹⁰⁶ which performs a two-step reads mapping, first on the rCRS
1483 mitochondrial reference genome and then simultaneously on the hg19 nuclear reference and
1484 rCRS reference, to remove possible nuclear-mitochondrial DNA sequences (NumtS)
1485 contaminations. PCR duplicates were removed with MarkDuplicates in the picard package
1486 (<https://gatk.broadinstitute.org/hc/en-us/articles/360037052812-MarkDuplicates-Picard->)
1487 from sequencing generated with Illumina Miseq and with FastUniq ¹⁰⁷ from sequencing
1488 generated with amplicon-based library preparations. The average coverage obtained
1489 (percentage of mtDNA molecules covered by at least one read) was 99.5% for IS samples,
1490 98.7% for PD samples and 100% for controls. Average mitochondrial read depth was 1220X
1491 for IS samples, 1299X for PD samples and 2282X for controls. Mitochondrial variant calling
1492 was performed with the MToolBox pipeline, using the default options (minimum read depth
1493 per alternative allele ≥ 5 and minimum quality score per base ≥ 25). Haplogroup predictions
1494 were generated with the Haplogrep 2 software ⁹⁶, using VCF files with homoplasmic and
1495 nearly homoplasmic (i.e. with heteroplasmic fraction ≥ 0.8) variants generated with
1496 MToolBox.

1497

1498 **mtDNA haplogroup association in patient cohorts**

1499

1500 We measured fMet levels 631 participants from both the IS and PD cohorts (Stroke N = 282,
1501 Parkinson's disease N = 124, Control N = 225) in three batches (**Supplementary Table 14**).
1502 Haplogroup predictions were available for 95% of the samples (N=601) whose mtDNA were
1503 deep-sequenced. In the IS cohort, 92 IS cases and 17 controls were of haplogroup Uk, while
1504 190 IS cases and 139 controls were of other haplogroups. fMet levels measured in pg/ml
1505 from the samples were controlled for batch and data collection site (Cambridge or Newcastle)
1506 of fMet measurement using linear regression, and residuals were rank normalised for further
1507 analysis. Associations between normalized fMet levels with Uk haplogroup was tested using
1508 a logistic regression implemented with the R glm function (family = "binomial"). In the PD
1509 cohort, 2 PD cases and 4 controls were of haplogroup Uk, and 118 PD cases and 39 controls
1510 were of other haplogroups. fMet levels measured in pg/ml from the samples were controlled
1511 for batch of fMet measurement (all PD samples were collected at the same site) using linear
1512 regression, and residuals were rank normalised for further analysis in the same fashion as in
1513 the IS cohort.

1514 **EPIC Norfolk cohort**

1515

1516 We obtained incidences of late-onset diseases from 11,966 men and women from the EPIC-
1517 Norfolk prospective cohort EPIC Norfolk cohort. Participants were identified as having
1518 experienced an event if the corresponding ICD-10 code was registered on the death certificate
1519 (as the underlying cause of death or as a contributing factor), or as the cause of
1520 hospitalization. Participants were on average 60 years (standard deviation: 6 years) old and
1521 46.3% were men. All participants were flagged for mortality at the UK Office of National
1522 Statistics, and vital status was ascertained for the entire cohort. Death certificates were coded
1523 by trained nosologists according to the International Classification of Diseases (ICD), 10th
1524 revision. Hospitalization data were obtained using National Health Service numbers through
1525 linkage with the East Norfolk Health Authority (ENCORE) database, which contains
1526 information on all hospital contacts throughout England and Wales. Participants were
1527 identified as having experienced an event if the corresponding ICD-10 code was registered on

1528 the death certificate (as the underlying cause of death or as a contributing factor), or as the
1529 cause of hospitalization, **Supplementary Data 11**). The current study is based on follow-up
1530 to 31st March 2016.

1531

1532 **fMet measurement and mtDNA genotyping in EPIC-Norfolk**

1533

1534 fMet was measured from plasma samples stored in liquid nitrogen since baseline in 1993-97
1535 from a total of 11,966 men and women from the EPIC-Norfolk prospective cohort as part of
1536 an untargeted metabolomic profiling using Metabolon's DiscoveryHD4™ platform
1537 (Metabolon Inc., Morrisville, North Carolina, USA). Measurements were undertaken in two
1538 sub-cohorts of 5,989 and 5,977 participants, respectively, quasi-randomly selected from the
1539 full cohort. Prior to statistical analyses, fMet levels were transformed using the natural
1540 logarithm and values at the tail of the distribution, defined by mean \pm 5 x standard deviation,
1541 were replaced by the respective lower/upper bound. They were then rescaled to a mean of
1542 zero and standard deviation of one. Processing steps were performed for each batch
1543 separately. mtDNA haplogroups in 10,782 participants were obtained with Haplogrep v2,
1544 using the --chip option with genotype data at 262 (2 of which were multi-allelic) mtDNA
1545 variants on the Affymetrix UK Biobank Axiom genotyping array (895 participants fall under
1546 haplogroup Uk, 9,887 individuals were of other haplogroups).

1547 **Cox-proportional hazards models**

1548

1549 We used Cox-proportional hazards models to estimate hazard ratios for the association of
1550 fMet levels (log-transformed and standardized) with first incidences of 24 diseases and health
1551 outcomes during a 20-year follow up period, with age as the underlying time scale adjusting
1552 for sex. For each incident outcome, we excluded participants reporting an instance of the
1553 outcome at baseline. For cancer outcomes, we additionally excluded all participants with the
1554 onset of any cancer within six months after baseline. mtDNA haplogroups in 10,782
1555 participants were obtained with Haplogrep v2, using the --chip option with genotype data at
1556 262 (2 of which were multi-allelic) mtDNA variants on the Affymetrix UK Biobank Axiom
1557 genotyping array (895 participants fall under haplogroup Uk, 9,887 individuals were of other
1558 haplogroups). Cox-proportional hazard models for fMet effects on outcome incidence were
1559 calculated for all participants, as well as separately within each haplotype group on 24
1560 outcomes, all of which had more than 10 incidences in both haplogroup Uk or otherwise, to
1561 obtain both cohort-based and haplotype-specific hazard ratios (**Supplementary Table 16**).
1562 As fMet levels are positively associated with age, we note that outcomes with significant
1563 fMet hazard ratios were late-onset diseases or outcomes. Further, we checked for potential
1564 violations of the proportional hazard assumption in Cox-proportional hazard models with
1565 Schoenfeld residual tests for each outcome, and found no violations except for cataracts, on
1566 which fMet does not have a significant effect (**Supplementary Table 16**). We further tested
1567 for interaction effects between mtDNA haplogroups and fMet levels to formally test for
1568 differences in effect estimates, and found no significant interaction effects (**Supplementary**
1569 **Table 16**).

1570

1571 **Data availability:**

1572 The Genotype-Tissue Expression (GTEx) Project was supported by the Common Fund of the
1573 Office of the Director of the National Institutes of Health, and by NCI, NHGRI, NHLBI,
1574 NIDA, NIMH, and NINDS. The data used for the analyses described in this manuscript were
1575 obtained from the GTEx Portal (GTEx_Analysis_2016-01-15_v7_RNASeQCv1.1.8) and

1576 dbGaP accession number phs000424.v7.p2. All data is available in the main text or the
1577 supplementary materials, or available upon request to the authors.

1578

1579 **Code Availability:**

1580 We conducted our analyses using the following published and publicly available software: 1)
1581 calling mtDNA variants: GATK v4.0.3.0 HaplotypeCaller
1582 (<https://gatk.broadinstitute.org/hc/en-us/articles/360037225632-HaplotypeCaller>) and mtdna-
1583 server local version (<https://github.com/seppinho/mutserve>), 2) for mtDNA association
1584 analysis using a linear mixed model and variance decomposition analysis: LDAK v5
1585 (<http://dougspeed.com/downloads2/>); 3) for improving power of mtDNA association: CMS
1586 v1.0 (<https://github.com/haschard/CMS>); 4) for eQTL analyses: limix v3.0
1587 (<https://github.com/limix/limix>); 5) for identifying pseudogenes in nuclear genome with high
1588 sequence similarity to the mtDNA: lastal 744 (<http://last.cbrc.jp/doc/lastal.html>); 6) for
1589 identifying PEER factors that capture unknown confounding in gene expression data: PEER
1590 v1.3 (<https://github.com/PMBio/peer>); 7) differential expression analysis: edgeR v3.11
1591 (<http://bioconductor.org/packages/release/bioc/html/edgeR.html>); 8) gene set enrichment
1592 analysis: GSEA v4.1.0 (<https://www.gsea-msigdb.org/gsea/index.jsp>); 9) Flow cytometry
1593 analysis: FlowJo v10.2

1594

1595 **Supplementary Materials:**

1596 Members of the ICICLE-PD Consortium, Supplementary Discussion, Supplementary Figures
1597 1-7, Supplementary Tables 1-18, References (1-10)

1598

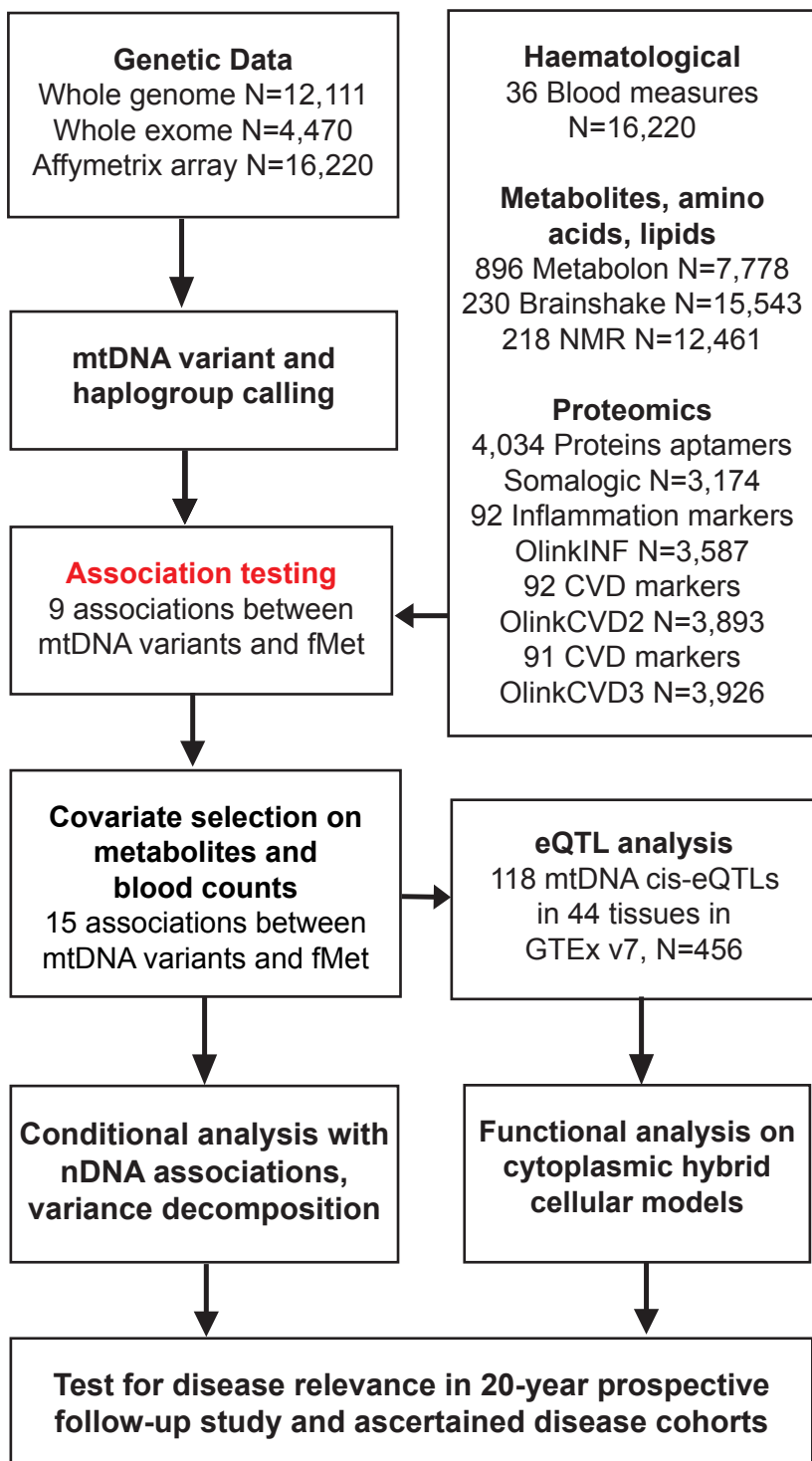
1599 **Members of the ICICLE-PD Consortium included in the author list:**

1600 Caroline H Williams-Gray³

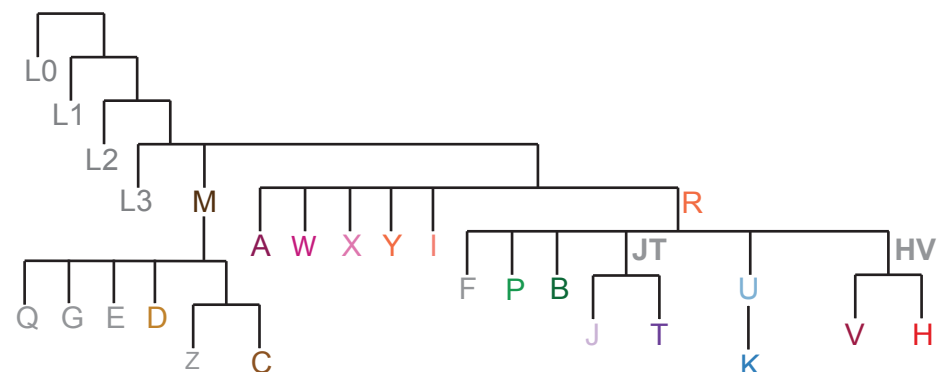
1601 3. Department of Clinical Neurosciences, School of Clinical Medicine, University of
1602 Cambridge, Cambridge Biomedical Campus, Cambridge, UK

1603 A full list of members and their affiliations appears in the Supplementary Information.

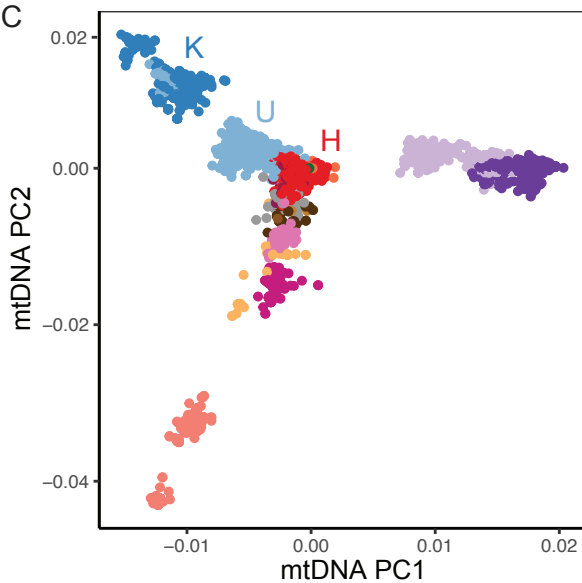
A



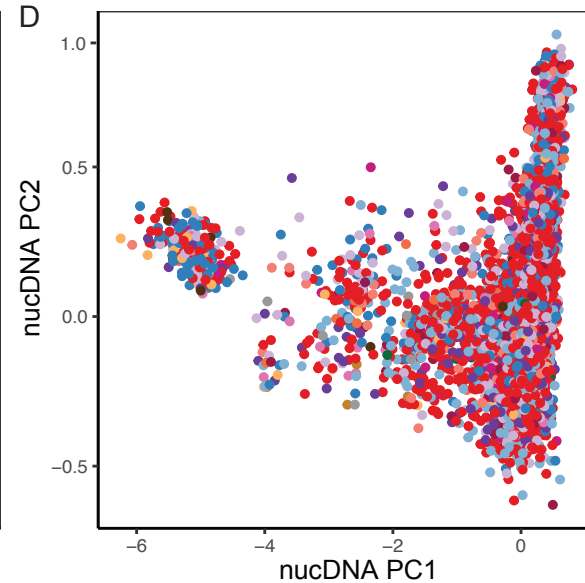
B



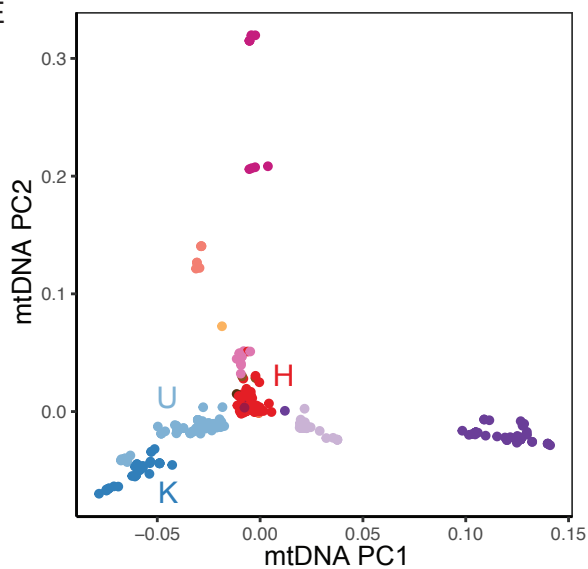
C



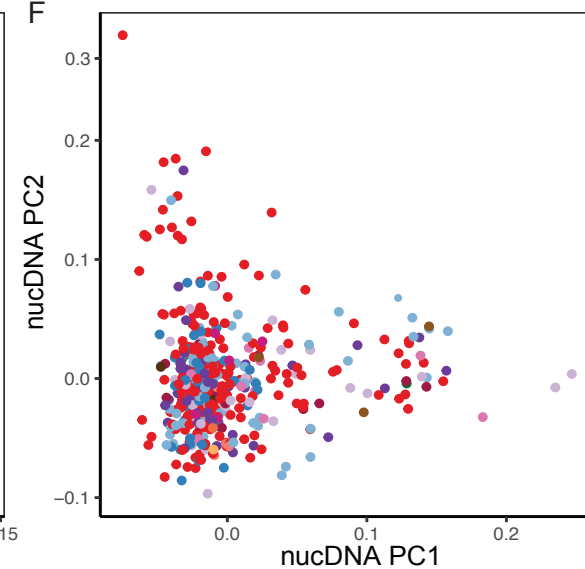
D

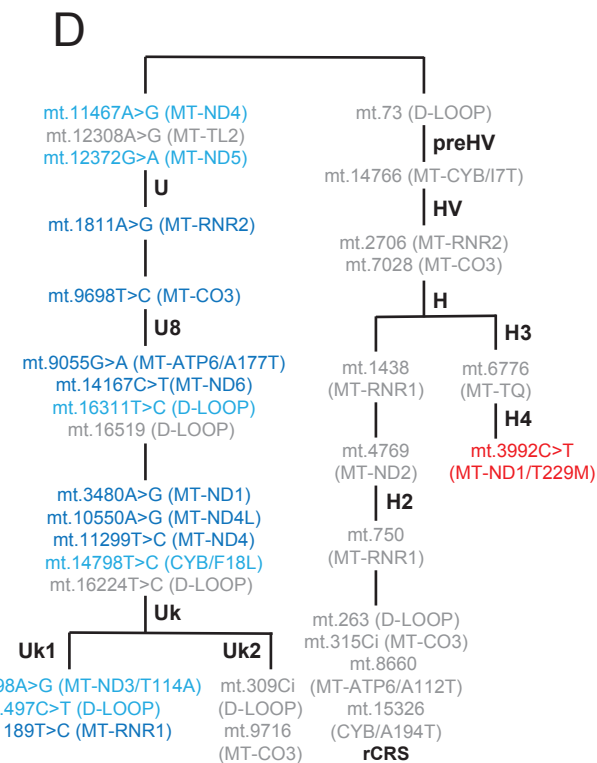
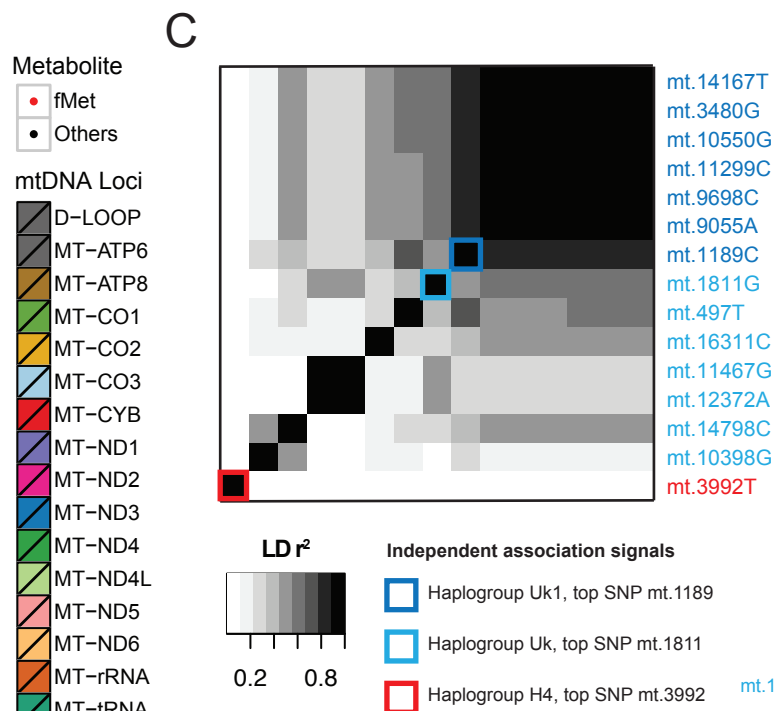
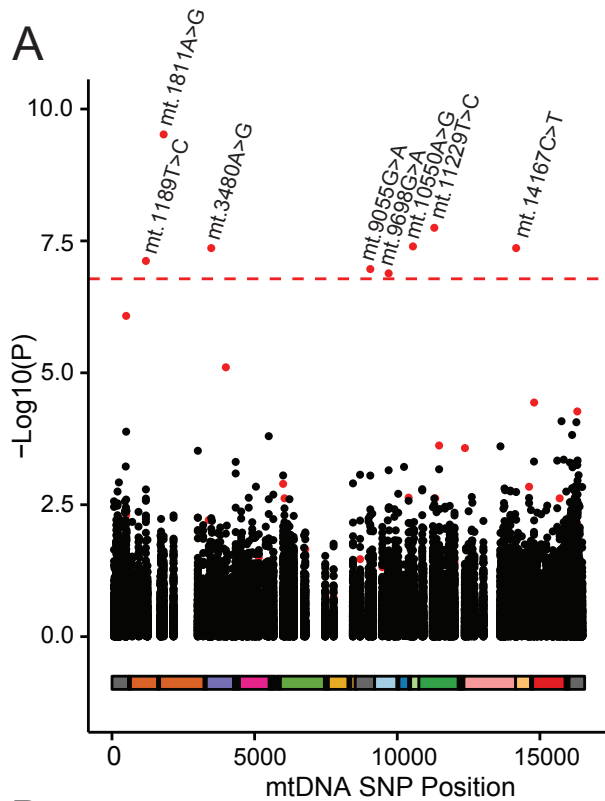


E



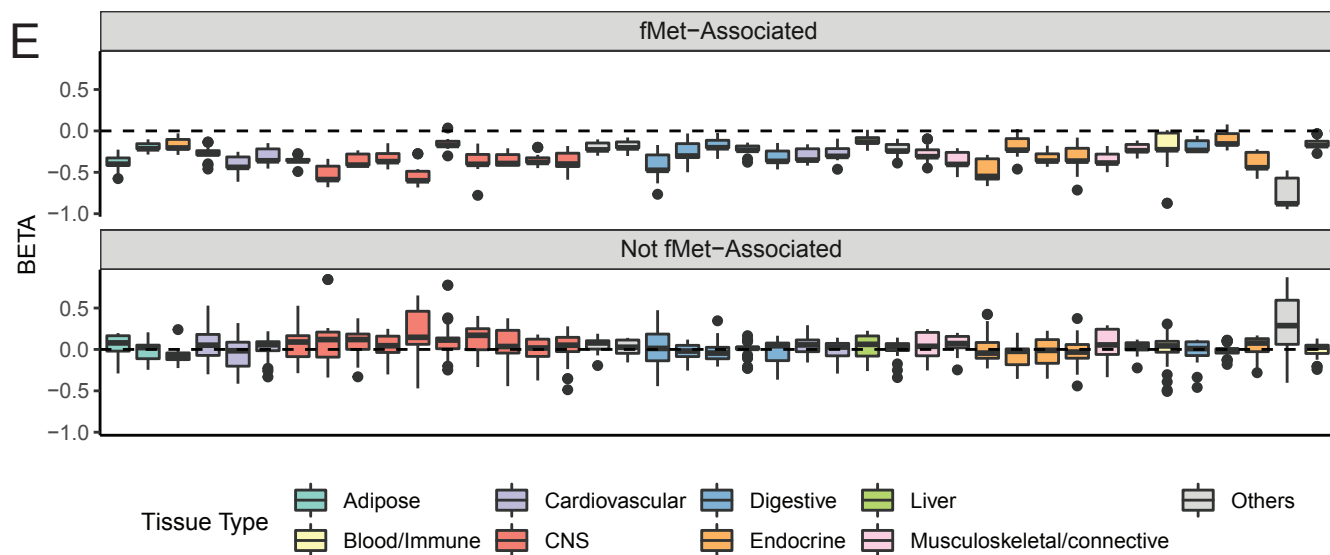
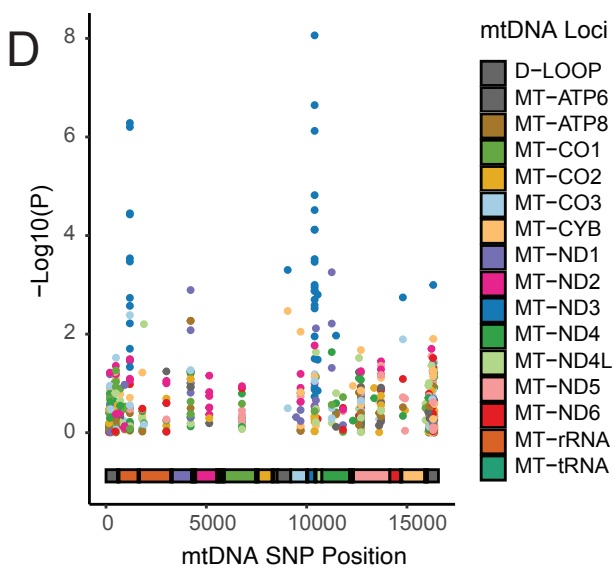
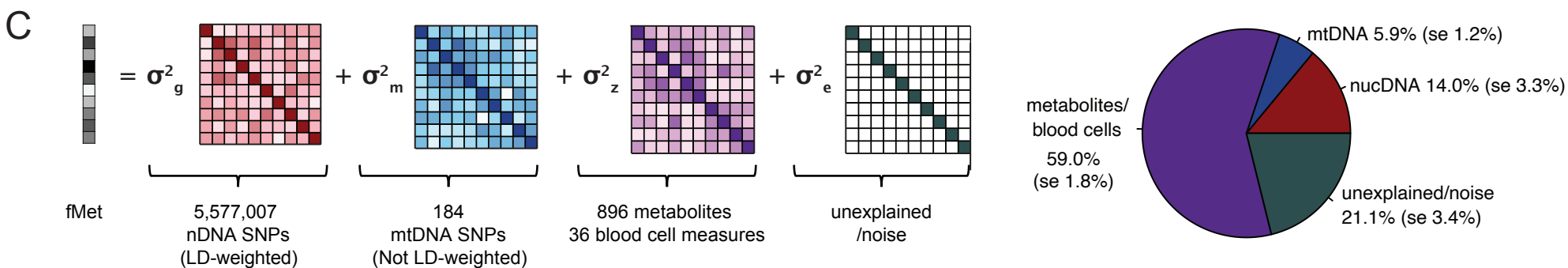
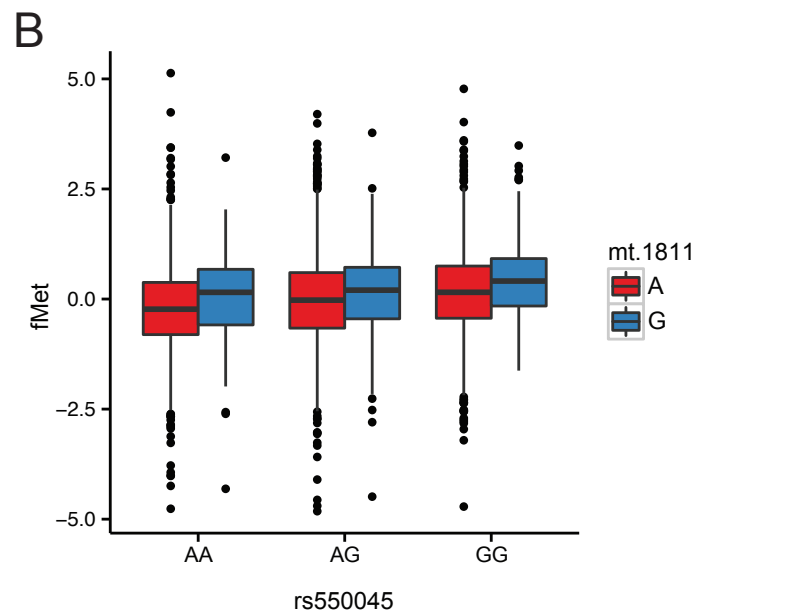
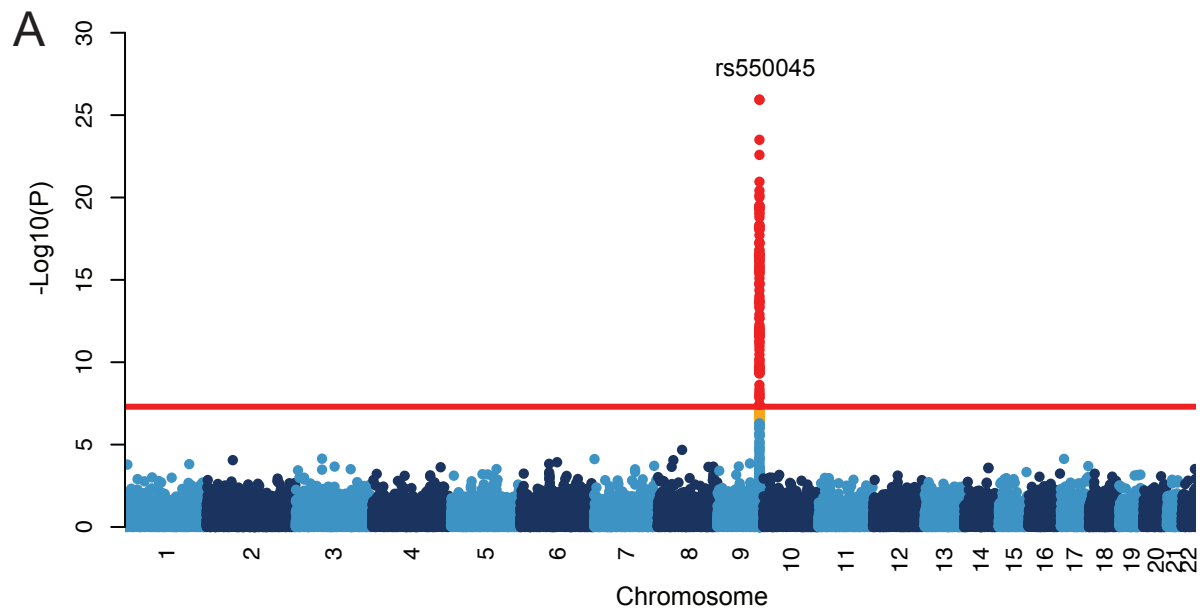
F

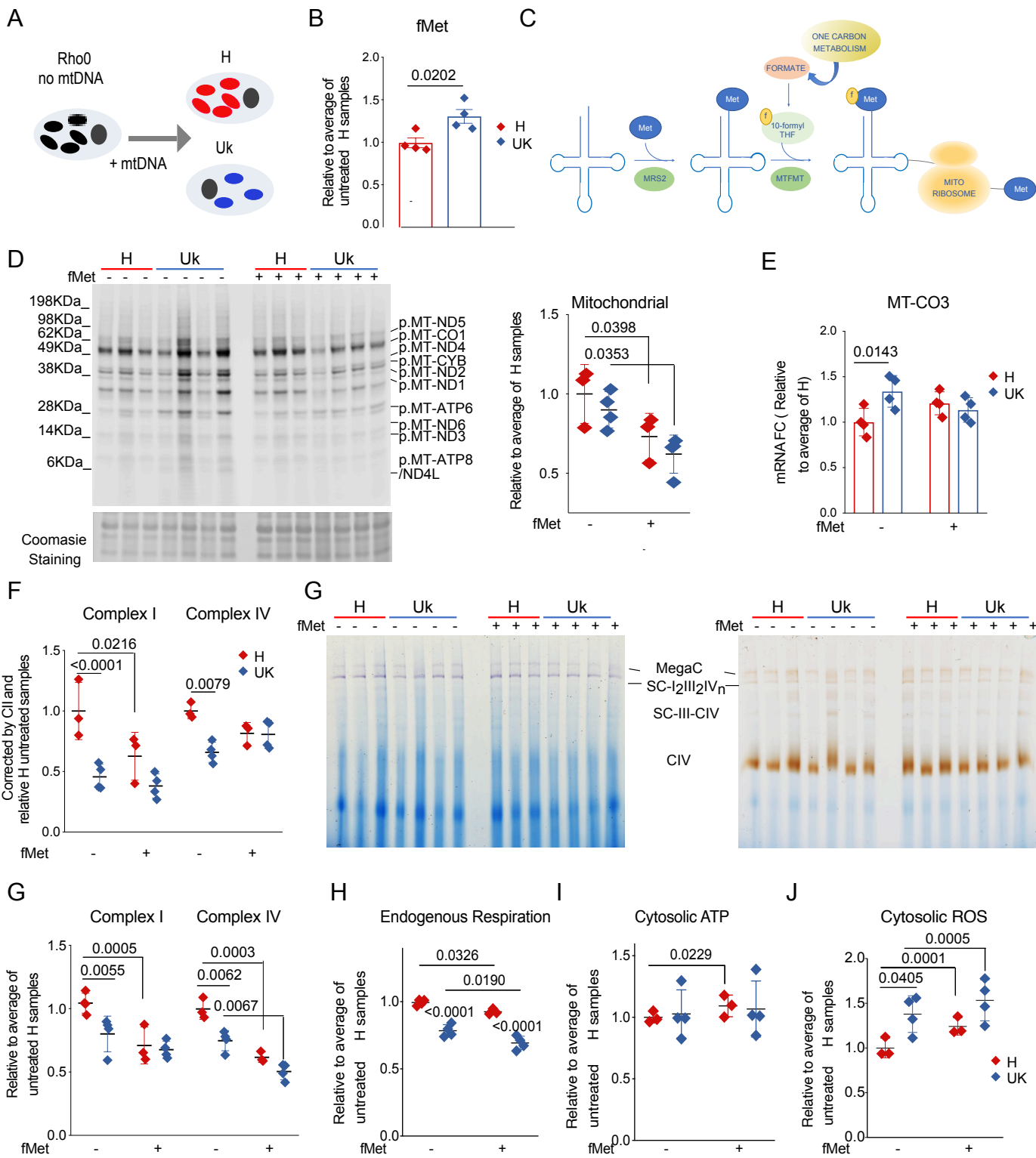


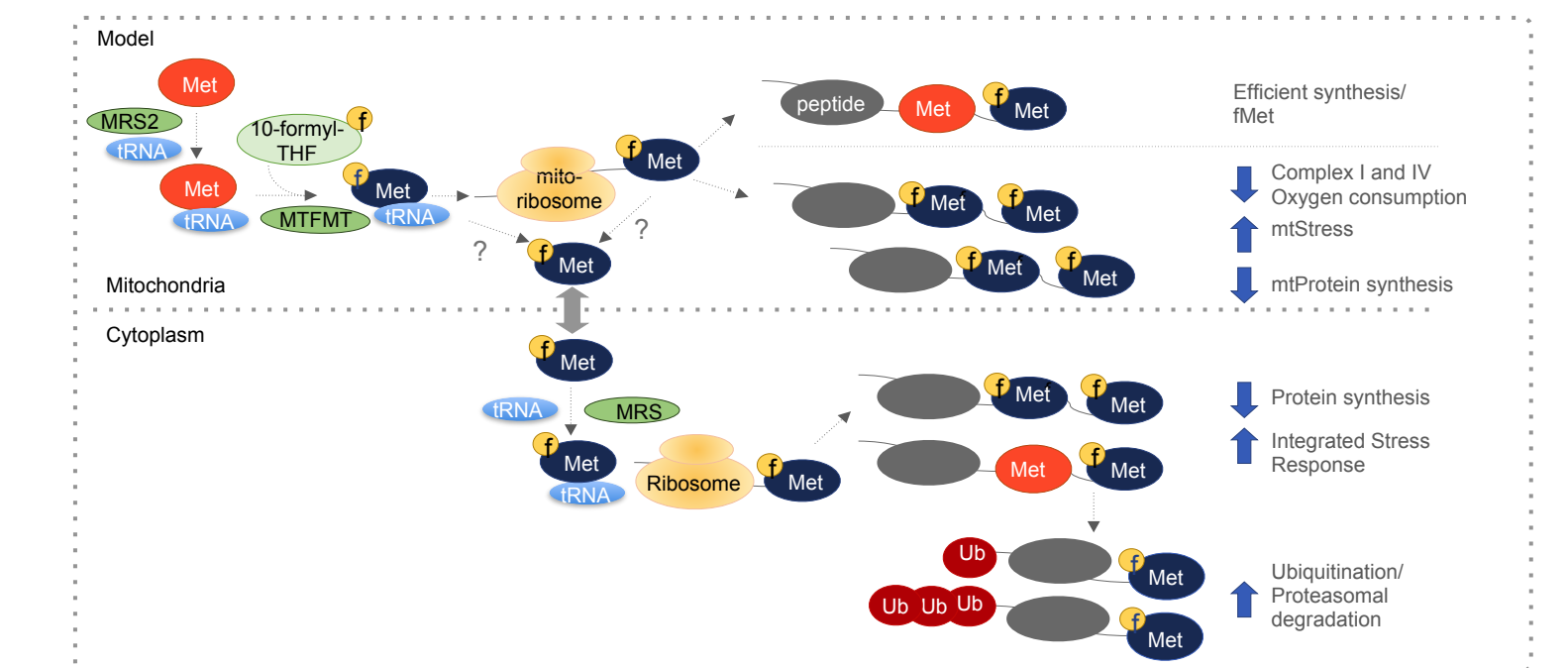
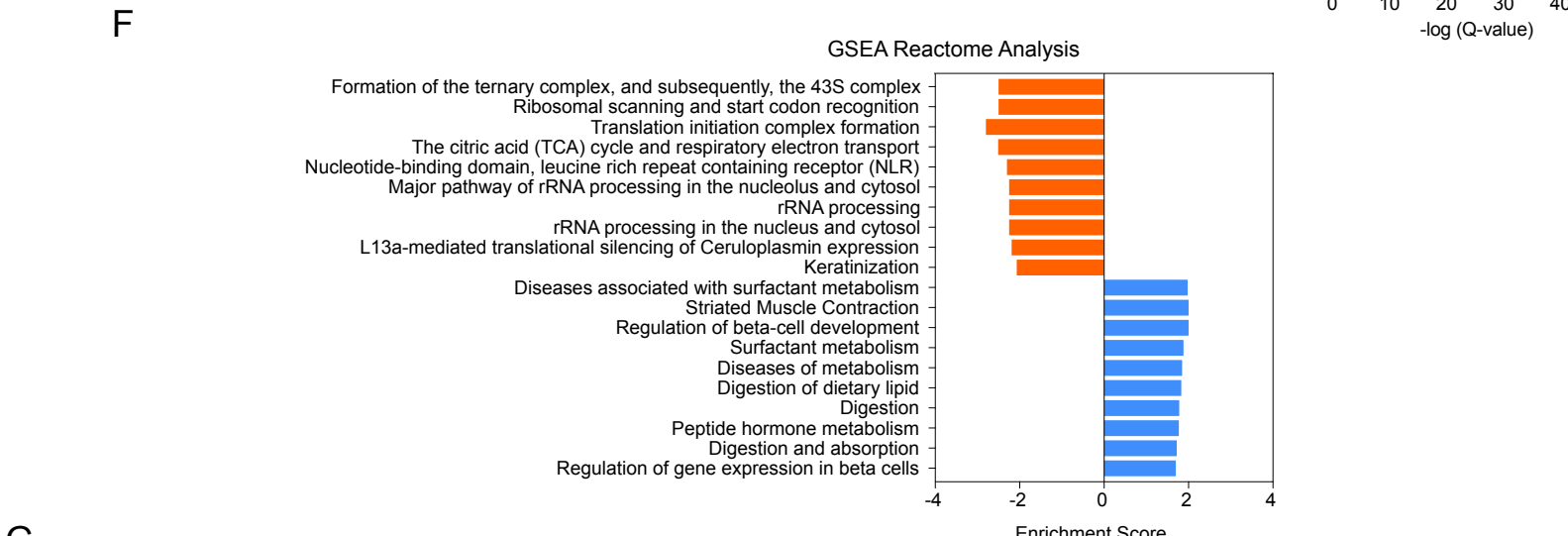
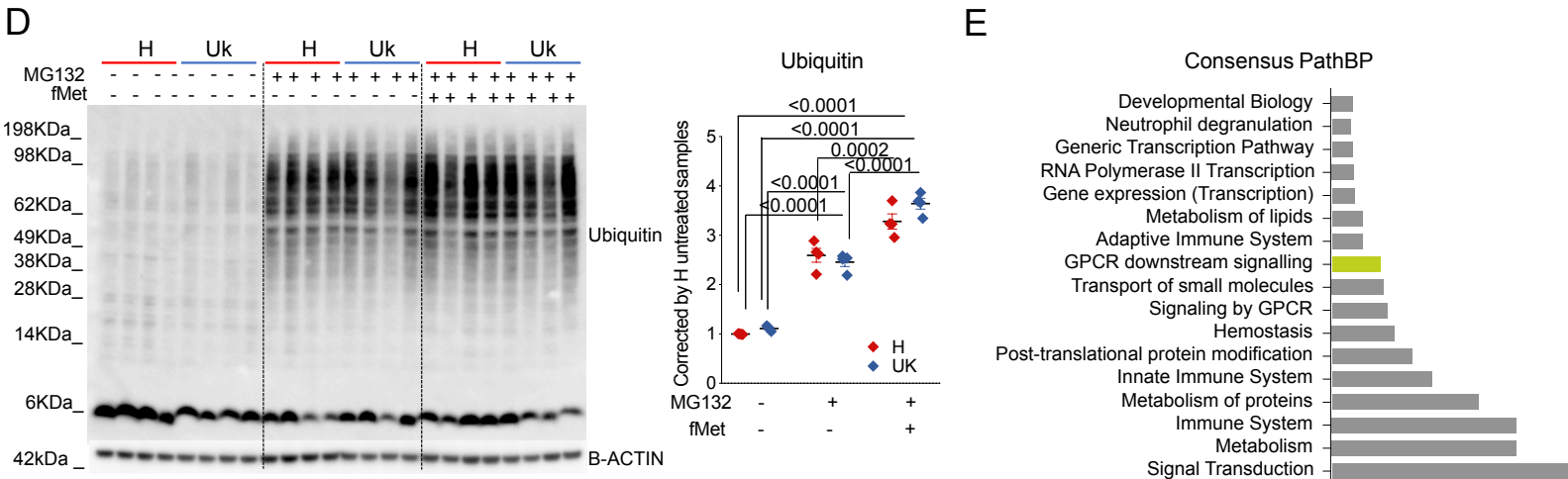
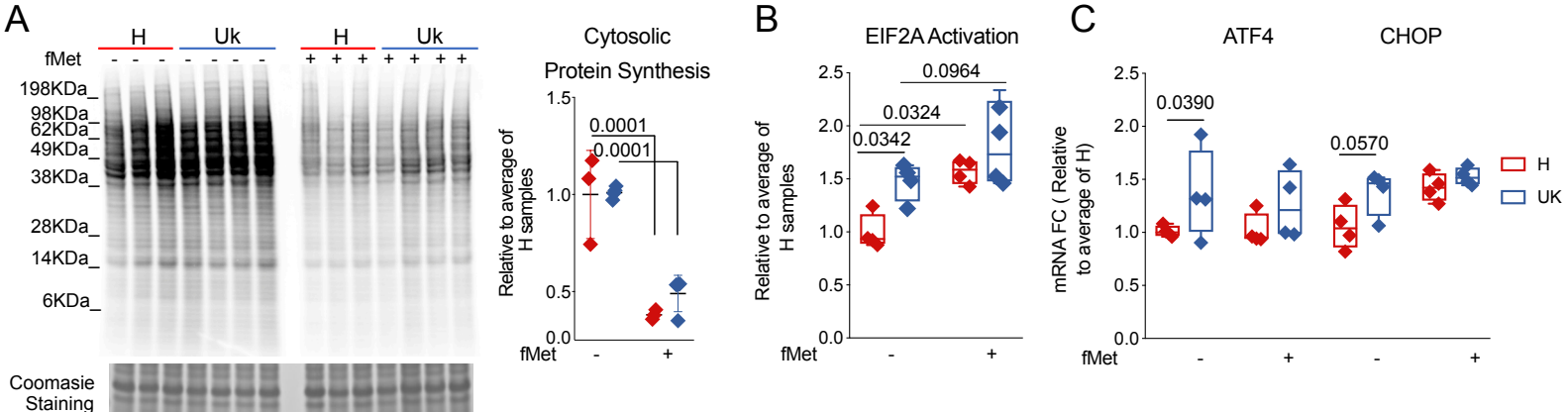


B

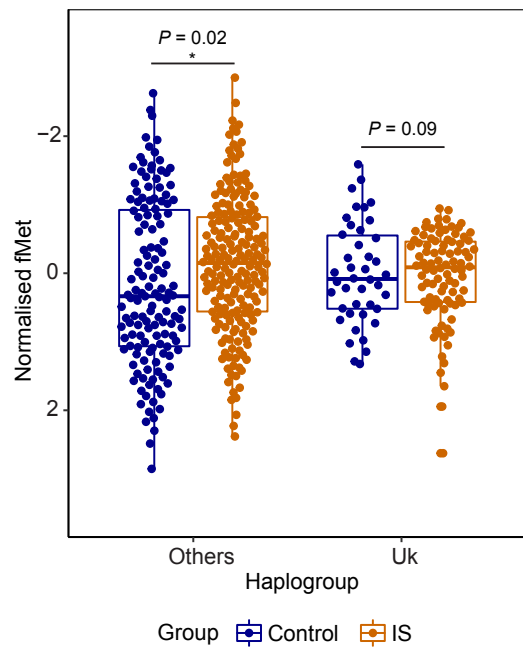
SNP						DISCOVERY (INTERVAL)						REPLICATION (EPIC)		
						LMM			CMS			LMM		
BP	GENE	A1	A0	A1FREQ	ANNOT	BETA	SD	P	BETA	SD	P	BETA	SD	P
497	D-LOOP	T	C	0.050	D-LOOP	0.144	0.029	8.33×10^{-7}	0.039	0.007	5.67×10^{-9}	0.100	0.022	4.74×10^{-6}
1189	MT-RNR1	C	T	0.072	12S rRNA	0.133	0.025	7.57×10^{-8}	0.057	0.007	7.36×10^{-18}	Not genotyped		
1811	MT-RNR2	G	A	0.129	16S rRNA	0.120	0.019	3.03×10^{-10}	0.057	0.007	4.63×10^{-18}	Not genotyped		
3480	MT-ND1	G	A	0.087	syn:K58K	0.125	0.023	4.32×10^{-8}	0.057	0.007	1.44×10^{-17}	0.078	0.017	3.33×10^{-6}
3992	MT-ND1	T	C	0.023	non-syn:T229M	-0.190	0.042	7.86×10^{-6}	-0.054	0.007	6.05×10^{-16}	-0.137	0.033	3.93×10^{-5}
9055	MT-ATP6	A	G	0.089	non-syn:A177T	0.119	0.022	1.08×10^{-7}	0.054	0.007	2.47×10^{-16}	0.070	0.016	2.13×10^{-5}
9698	MT-CO3	C	T	0.089	syn:L164L	0.119	0.023	1.30×10^{-7}	0.056	0.007	2.05×10^{-17}	0.079	0.017	3.23×10^{-6}
10398	MT-ND3	G	A	0.216	non-syn:T114S	0.048	0.016	2.32×10^{-3}	0.038	0.007	8.68×10^{-9}	Not genotyped		
10550	MT-ND4L	G	A	0.086	syn:M27M	0.126	0.023	4.02×10^{-8}	0.058	0.007	3.59×10^{-18}	0.079	0.017	2.82×10^{-6}
11229	MT-ND4	C	T	0.087	syn:T180T	0.128	0.023	1.78×10^{-8}	0.058	0.007	1.80×10^{-18}	Not genotyped		
11467	MT-ND4	G	A	0.229	syn:L236L	0.056	0.015	2.39×10^{-4}	0.040	0.007	1.50×10^{-9}	Not genotyped		
12372	MT-ND5	A	G	0.229	syn:L12L	0.056	0.015	2.67×10^{-4}	0.040	0.007	1.82×10^{-9}	Not genotyped		
14167	MT-ND6	T	C	0.086	syn:E169E	0.125	0.023	4.32×10^{-8}	0.057	0.007	1.44×10^{-17}	0.076	0.017	7.62×10^{-6}
14798	MT-CYB	C	T	0.166	non-syn:F18L	0.071	0.017	3.66×10^{-5}	0.039	0.007	3.92×10^{-9}	0.044	0.013	6.83×10^{-4}
16311	D-LOOP	C	T	0.088	D-LOOP	0.070	0.017	5.39×10^{-5}	0.050	0.007	3.60×10^{-12}	Not genotyped		



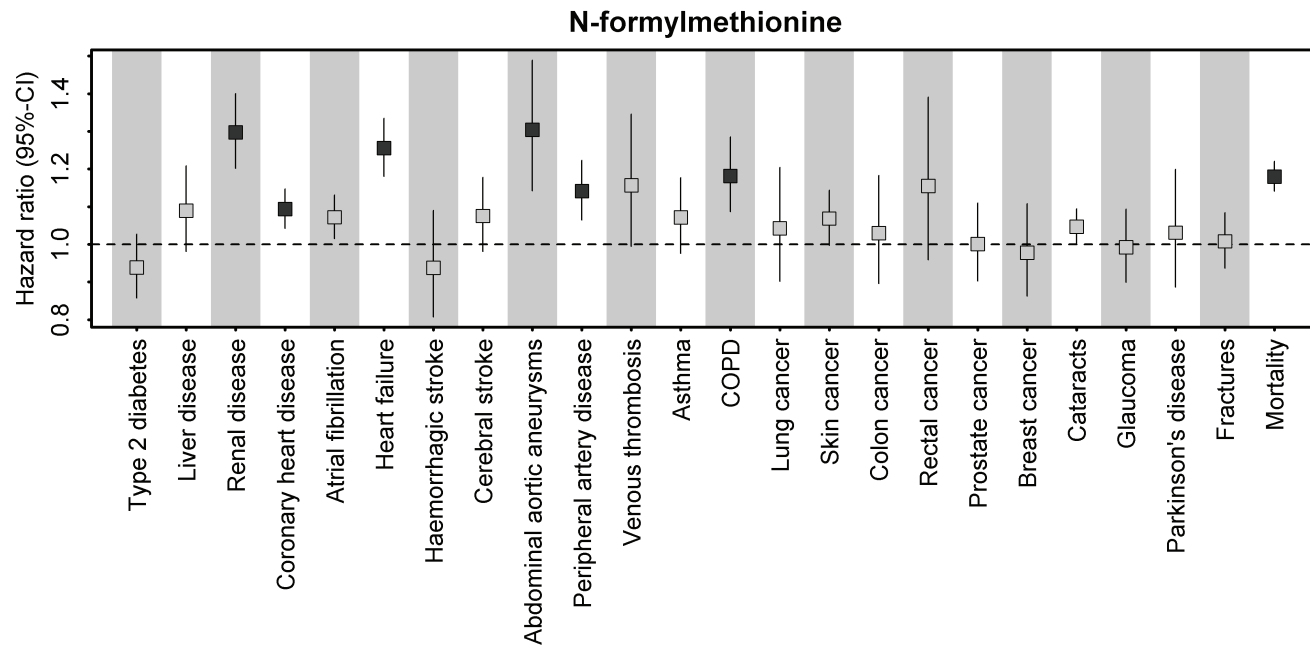




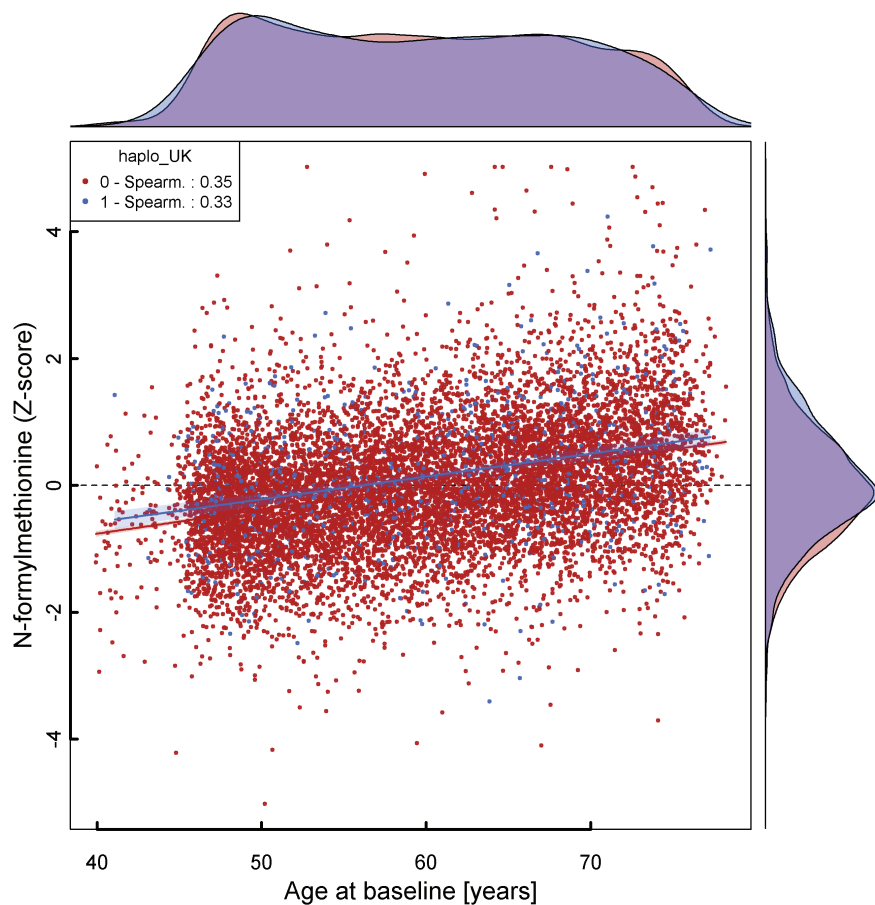
A



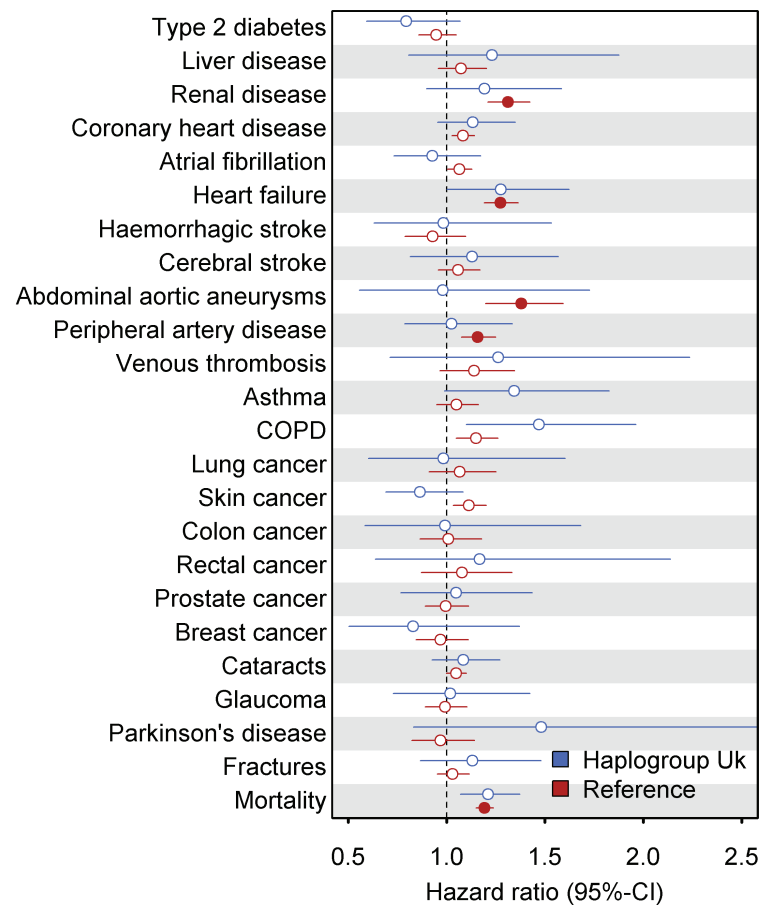
B

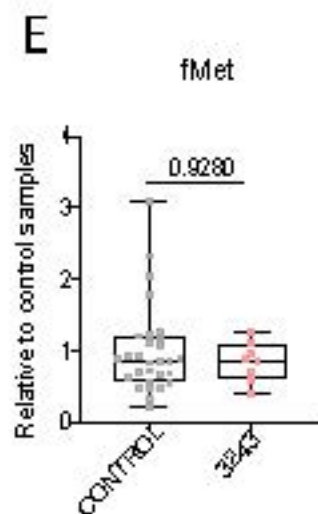
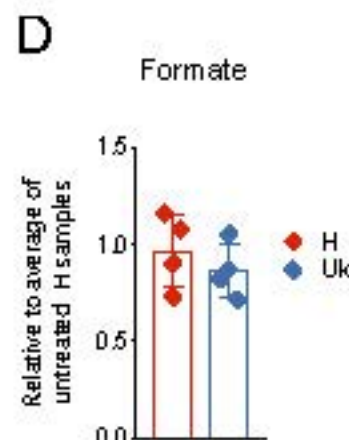
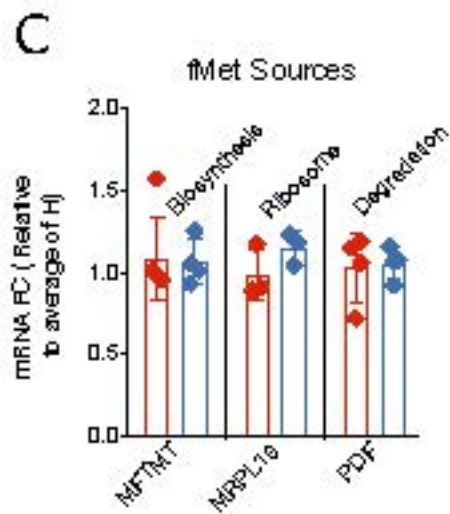
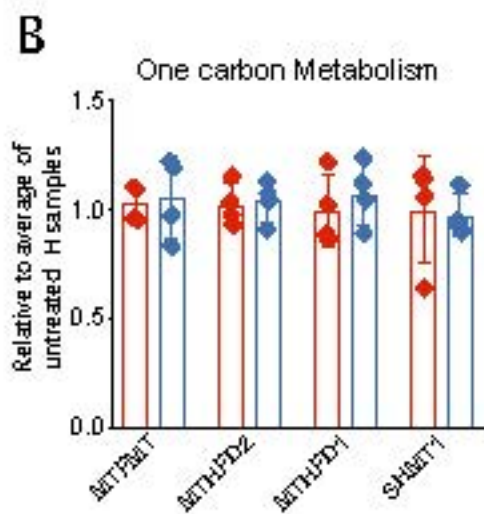
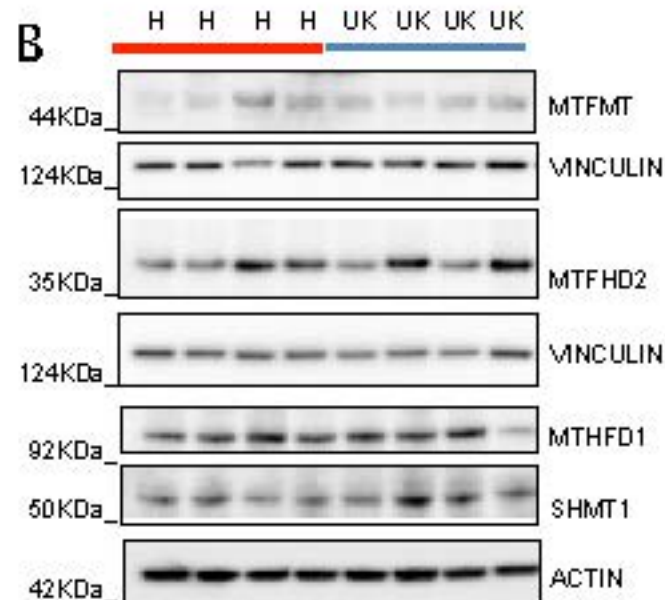
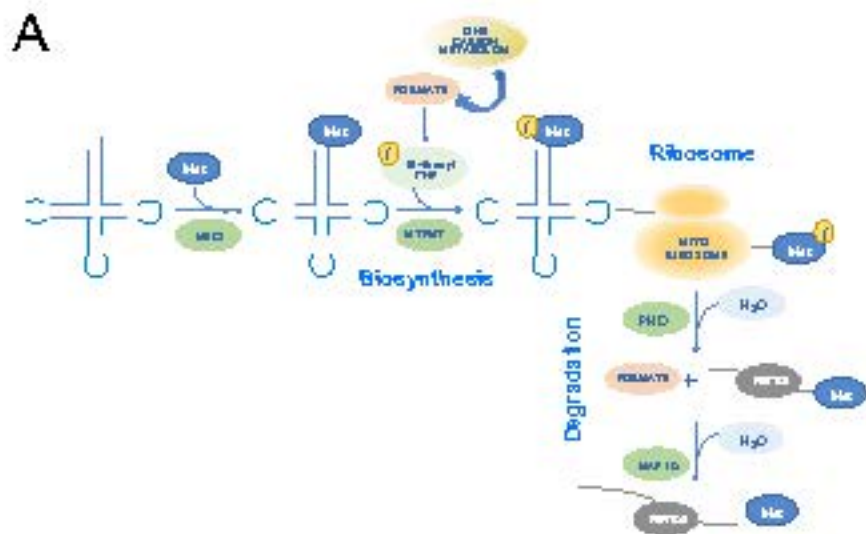


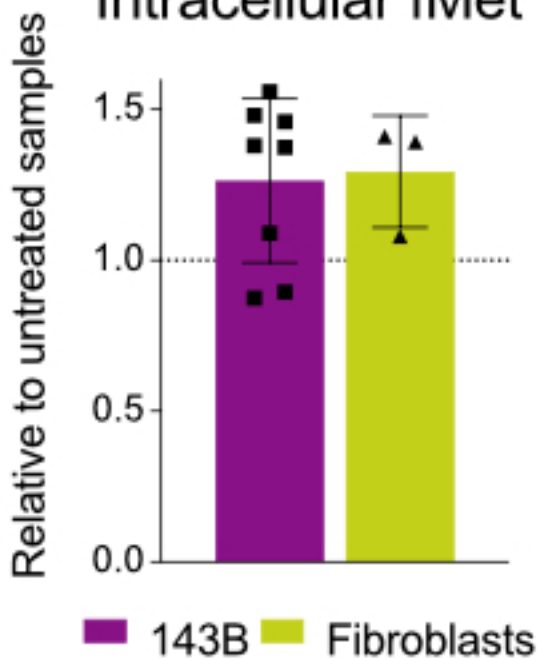
C



D





A**Intracellular fMet****B**