1  **Respiratory complex and tissue lineage drive mutational patterns in the tumor**

2  **mitochondrial genome**

3

4  Alexander N. Gorelick[1,5], Minsoo Kim[1], Walid K. Chatila[1,4,5], Konnor La[2], A. Ari Hakimi[3], Barry S.

5  Taylor[1,4,5], Payam A. Gammage[6,7,*], Ed Reznik[1,4,*,†]

6  [1] Computational Oncology Service, Memorial Sloan Kettering Cancer Center

7  [2] Laboratory of Metabolic Regulation and Genetics, Rockefeller University

8  [3] Urology Service, Memorial Sloan Kettering Cancer Center

9  [4] Marie-Josee and Henry R. Kravis Center for Molecular Oncology, Memorial Sloan Kettering Cancer

10  Center

11  [5] Human Oncology and Pathogenesis Program, Memorial Sloan Kettering Cancer Center

12  [6] CRUK Beatson Institute, Glasgow, UK

13  [7] Institute of Cancer Sciences, University of Glasgow, Glasgow, UK

14  [*] Correspondence to: Ed Reznik (reznike@mskcc.org) and Payam Gammage

15  (payam.gammage@glasgow.ac.uk)

16  [†] Lead Contact

17

18

19  **Abstract**

20  Mitochondrial DNA (mtDNA) encodes essential protein subunits and translational machinery for

21  four distinct complexes of oxidative phosphorylation (OXPHOS). Using repurposed whole-

22  exome sequencing data, we demonstrate that pathogenic mtDNA mutations arise in tumors at a

23  rate comparable to the most common cancer driver genes. We identify OXPHOS complexes as

24  critical determinants shaping somatic mtDNA mutation patterns across tumor lineages. Loss-of-

25  function mutations accumulate at an elevated rate specifically in Complex I, and often arise at

26  specific homopolymeric hotspots. In contrast, Complex V is depleted of all non-synonymous

27  mutations, suggesting that mutations directly impacting ATP synthesis are under negative

28  selection. Both common truncating mutations and rarer missense alleles are associated with a

29  pan-lineage transcriptional program, even in cancer types where mtDNA mutations are

30  comparatively rare. Pathogenic mutations of mtDNA are associated with substantial increases in

31  overall survival of colorectal adenocarcinoma patients, demonstrating a clear functional

32  relationship between genotype and phenotype. The mitochondrial genome is therefore

33  frequently and functionally disrupted across many cancers, with significant implications for

34  patient stratification, prognosis and therapeutic development.

35

36

37

38

39 **Introduction**

40 Somatic mutations are the underlying drivers of malignancy in cancer, and the identification and

41 characterization of recurrent, functional somatic events has been the capstone goal of cancer

42 genomics. Genomic searches for recurrent driver mutations have focused on the nuclear exome

43 or subsets thereof, motivated by the observation that recurrent mutations are concentrated in

44 the coding regions of a subset of nuclear-DNA-encoded genes. This targeted approach has

45 powered the discovery of common and rare driver mutations in exonic regions, but by corollary

46 has also left underexplored the overwhelming majority of the genome and the driver events it

47 may harbor. Numerous examples now exist of the prevalence and function of oncogenic

48 mutations beyond the nuclear exome, including mutations to the *TERT* promoter, non-coding

49 RNAs including ribosomal RNA and snRNAs, and enhancers [1]. A fundamental challenge is

50 therefore to discover new functional somatic alterations beyond the exome with a fixed and

51 limited sequencing capacity.

52

53 Somatic mutations in tumors commonly target human mitochondrial DNA (mtDNA)[2–6], affecting

54 both the thirteen essential protein components of four distinct complexes (CI, CIII, CIV, and CV)

55 in oxidative phosphorylation (OXPHOS) as well as the non-coding RNA (22 tRNAs, 2 rRNAs)

56 necessary for mtDNA translation. (**Fig. 1a**). Despite abundant pharmacological, genetic, and

57 clinical data demonstrating that perturbation of different OXPHOS complexes (referred to in

58 shorthand as complexes) produce distinct cellular adaptations [7,8], the importance of each

59 complex in shaping mtDNA mutation patterns in cancer is unknown. Because mtDNA is not

60 commonly targeted by exome sequencing panels, prior analyses of mtDNA mutations have

61 relied on cohorts profiled with whole genome sequencing,  with consequently diminished

62 statistical power to detect recurrent patterns of mutations relative to exome sequencing studies[8].

63 However, due to the extremely high copy number and off-target hybridization rate of mtDNA,

64 mtDNA reads are abundant in widely-available exome sequencing of tumors[9]. Mitochondrial

65 DNA therefore represents an opportunity for discovery through repurposing of existing exome

66 sequencing data.

67

68 Here, by utilizing existing exome sequencing data to more than double statistical power of prior

69 analyses, we report that OXPHOS complex, in combination with tissue lineage and mutational

70 consequence, is a critical determinant of mtDNA mutation patterns in cancer. We find that

71 NADH:ubiquinone oxidoreductase (complex I, CI) mutations are strongly enriched for highly

72 pathogenic mutations in specific tissue lineages, whereas ATP synthase (complex V, CV) is

73    broadly depleted of all non-synonymous mutations. We further identify six highly recurrent

74    mtDNA mutation hotspots at specific homopolymer sequence contexts, which collectively

75    account for over 40% of all truncating mutations to mtDNA, as well as recurrent mutations in

76    both protein-coding genes and non-coding RNA elements. These mutations produce a defined,

77    lineage-agnostic transcriptional program and, in specific tumor lineages, associate with both

78    underlying molecular subtypes and clinical outcomes. Our results argue that specific

79    components of mitochondrial respiration are broadly perturbed across many tissue lineages,

80    and that re-analysis of existing genomic data can yield new discoveries in underexplored

81    genomic terrain.

82

83    **Results**

84    *mtDNA Mutations in Tumors from Off-target Reads*

85    To study patterns of mtDNA mutations in tumors, we reasoned that the sheer amount of off-

86    target reads aligning to mtDNA in whole-exome sequencing data would be sufficient to call

87    somatic mtDNA mutations in a large proportion of samples. We therefore assembled a dataset

88    of pan-cancer paired tumor and matched-normal exome sequencing samples from the TCGA,

89    *n*=10,132 (**Supplementary Fig. 1a**). Inconsistent sequencing coverage between samples is an

90    inherent limitation to this approach, as variants located in regions without adequate sequencing

91    coverage are not identifiable, and we therefore developed our methodology to be cognizant of

92    the sequencing coverage at each position in each sample (see **Methods**). We focused our

93    analysis on regions of mtDNA in protein-coding genes and genes coding for mitochondrial

94    ribosomal RNAs (rRNAs) and transfer RNAs (tRNAs), excluding the control region (pos 1-576,

95    16,024-16,569), several known hypervariable loci (pos 302-314; 514-524; 3,106-3,109), and 89

96    remaining positions not within any genic region from all further analyses (excluded positions

97    listed in **SI table 1**). The combination of an increase in sample size (in TCGA, relative to whole

98    genome cohorts) and high off-target read coverage effectively doubled the number of tumor-

99    associated mtDNA genomes sequenced compared to the largest published dataset of whole-

100    genome sequenced tumor mitochondrial genomes [3]: On average 6,100 tumors were sequenced

101    at sufficient depth to call mutations at each mtDNA position (mean+/-SD: 5,399-6,800 samples

102    covered at a given position, **Fig. 1b, Supplementary Fig. 1b**), compared to 2,836 whole-

103    genome tumor sequences from the PCAWG dataset. When further restricted to regions

104    sequenced at sufficient depth in both tumor and matched-normal samples, each position was

105    covered in 4,769 tumor/normal pairs on average (mean+/-SD: 4,148-5,390 samples).

106

107   We implemented a conservative variant calling approach modeled after state-of-the-art

108   methodologies for exome sequencing, in which we took the intersection of two variant callers

109   (MuTect2 [10] and an in-house variant caller based on the SAMtools mpileup utility [11], see

110   **Methods**).  Consistent with prior work, mtDNA variants exhibited a strand-specific enrichment

111   for C>T mutations on the heavy strand and T>C mutations on the light strand (**Supplementary**

112   **Fig. 1c**).  Based on 789 tumor samples from TCGA with whole genome sequences in the

113   PCAWG cohort [3], 95.6% of mutation calls from whole exome sequencing validated against

114   published mutation calls from the PCAWG data (**Fig. 1c**). We also evaluated the possibility that

115   nuclear-encoded mitochondrial pseudogenes (NUMTs) could corrupt variant calling. Although

116   both mtDNA and NUMTs are not targeted by exome sequencing, mtDNA is unique in that it

117   exists at orders of magnitude higher copy number in each cell and, critically, is expressed at

118   extremely high levels, whereas NUMTs do not show evidence of significant transcription [12]. We

119   therefore determined the fraction of somatic mtDNA variants from exome sequencing which

120   could be recapitulated in matched RNA sequencing from the same sample, revealing that

121   96.9% of such variants were validated. Finally, we observed a strong correlation between DNA

122   and RNA heteroplasmy overall (Pearson's r = 0.918) (**Fig. 1d**), confirming that the vast majority

123   of observed mutations are expressed and providing further evidence that the mutations called

124   by our approach are not attributable to NUMTs. In total, we identified 4,381 mtDNA mutations

125   from 10,132 tumor samples which were either protein-truncating (*i.e.* frame-shift indels or

126   nonsense mutations); or non-truncating variants (missense, in-frame indels, translation start-

127   site, non-stop, or mutations to tRNA/rRNAs) which were detected in tumor and absent from

128   matched-normal samples. Among a subset of 3,264 paired tumor/normal samples with sufficient

129   coverage to call mtDNA mutations in at least 90% of the mitochondrial genome (32% of

130   tumor/normal pairs in our dataset overall, referred to throughout as "well-covered" samples),

131   57% (95%CI 56-59%) had at least one mtDNA variant, in agreement with previous estimates for

132   mtDNA mutation incidence in pan-cancer sequencing data [2].  Consistent with independent

133   mutagenic processes operating in the nuclear and mitochondrial genomes, we observed no

134   correlation between nuclear and mitochondrial mutation burdens pan-cancer or within individual

135   cancer types (**Fig. 1e, Supplementary Fig. 1d)**. Furthermore, in colorectal and stomach

136   cancers where microsatellite instability (MSI) is common, the presence of MSI affected mutation

137   burden in the nuclear but not in the mitochondrial genome (**Supplementary Fig. 1e**).

138

139   The mutation rate in the coding region of mtDNA is roughly 67.8 mut/Mb, roughly 6-fold higher

140   than the rate in 468 cancer-associated genes in the MSK-IMPACT panel [13] of 11.3 mut/Mb ($P <$

141    $10^{-308}$ (computational limit of detection), two-sided Poisson test). When calculated for each

142    gene, we also observed mtDNA-encoded genes to have enriched mutation rates compared to

143    nuclear-DNA-encoded MSK-IMPACT genes ($P$=2x10$^{-22}$, two-sided Wilcoxon rank sum test):

144    only 2 MSK-IMPACT genes (*TP53, KRAS*) exhibited rates higher than that of the most mutated

145    mtDNA-encoded genes (**Fig. 1f**). Furthermore, the 13 protein-coding mtDNA genes exhibited a

146    4.2-fold higher rate of truncating variants which disrupt the reading frame (*i.e.* nonsense

147    mutations and frameshift indels) compared to truncating mutations among 185 known tumor

148    suppressor genes in the MSK-IMPACT cohort (which also included splice-site variants which

149    cannot arise in the mitochondrial genome due to the lack of introns) ($P$=9x10$^{-5}$, two-sided

150    Wilcoxon rank sum test) (**Fig. 1g**), and a 6.7-fold higher rate of non-truncating, non-synonymous

151    mutations (collectively referred to here as "missense" mutations) than 168 MSK-IMPACT

152    oncogenes ($P$=6x10$^{-9}$) (**Fig. 1h**). In total, 11.9% of tumors across all cancers (95% CI: 11.0-

153    12.9%) harbored a truncating mtDNA variant absent in the patient's matched normal sample. In

154    contrast, only 0.15% of normal blood samples exhibited a truncating variant (95% CI: 0.13-

155    0.17%) based on a recent analysis of ~200,000 mtDNA genomes [14] (**Fig. 1i)**. The rate of

156    truncating mutations in mtDNA genes in tumors therefore represents an 80-fold increase

157    compared to truncating mutations observed in normal human genomes (**SI table 2.)** Of the 619

158    truncating mutations we observed, 196 (32%, 95% binomial CI: 28-35%) had >80%

159    heteroplasmy despite underlying infiltration of the bulk tumor by normal stromal and immune

160    cells, indicating that a significant number of tumors are dominated by a highly dysfunctional

161    mitochondrial genotype. Furthermore, high heteroplasmy truncating variants were significantly

162    more common than high-heteroplasmy silent mutations (139/555, 25%, 95% CI: 21-29%)

163    ($P$=0.01, two-sided Fisher's exact test) under predominantly neutral selection.

164

165    *Truncating Mutations Preferentially Target Complex I at Homopolymeric Hotspots*

166    The physiologic response to genetic and pharmacologic inhibition of mitochondrial respiration

167    depends strongly on which mtDNA-encoded complex (CI, CIII, CIV, CV) is disrupted, implicating

168    OXPHOS complex as a potential determinant of selective pressure for mutation. We therefore

169    investigated the somatic mutation rate according to the OXPHOS complex, controlling for the

170    relative length of mtDNA coding for genes in each complex and uneven coverage within each

171    sequenced sample. This revealed a striking dichotomy in the relative enrichment of mutations in

172    each complex. Truncating variants (nonsense mutations and frame-disrupting indels) arose at a

173    2-fold or greater rate in complex I relative to the other complexes ($P$=0.001 for least significant

174    comparison, two-sided Poisson test) (**Fig. 2a**). No difference in mutation rate between

175    complexes was observed for silent mutations, consistent with a lack of differential selective

176    pressure for synonymous protein changes (*P*=0.5 for most-significant comparison). Unlike

177    variants in other complexes, truncating variants in CI demonstrated higher heteroplasmy

178    (variant allele frequency) than silent variants ($P=1\times10^{-6}$, CI; most significant for other complexes,

179    *P*=0.4, two-sided Wilcoxon rank sum test) (**Fig. 2c**). Finally, complex V genes (*MT-ATP6* and

180    *MT-ATP8*) demonstrated significantly lower rates of truncating but not synonymous mutations.

181    The findings above were recapitulated in an independent cohort, composed of a distinct mixture

182    of cancer types, of N=1,951 whole-genome sequenced tumors from the PCAWG dataset, after

183    excluding samples overlapping with our own cohort (**Fig. 2b**). Tumors of different lineages

184    exhibited wide variability in the incidence of truncating mutations, with ≤5% of some cancer

185    types affected by truncating mutations (sarcomas, gliomas), to 20% or greater of other cancer

186    types (renal cell, colorectal, thyroid) (**Fig. 2d**). In renal, thyroid, and colorectal cancers, the high

187    burden of truncating variants was driven by a specific enrichment for mutations to complex I (*Q*-

188    value < 0.01, two-sided McNemar's test) (**Fig. 2e)**. Truncating variants in these three cancers

189    affected between 20-30% of all samples, corresponding to a prevalence on the same order or

190    exceeding that of common tumor suppressors in these diseases. Taken together, these data

191    indicate that the functional consequence of mtDNA variants and the complex they target are key

192    determinants of the pattern of somatic mtDNA mutations. Additionally, they suggest that

193    disruption of complex V, which would fundamentally impair mitochondrial ATP production

194    independent of the activity of all other OXPHOS complexes, is not tolerated.

195

196    Unexpectedly, we observed that truncating mutations frequently arose at the same genomic

197    locus, analogous to well-described hotspot mutations that accumulate in nuclear cancer driver

198    genes and often reflect selective pressure [15,16]. These apparently recurrent alleles were

199    exclusively indels, rather than nonsense mutations, characterized by a homopolymeric

200    sequence context. We therefore developed an approach to detecting recurrent mutations at

201    homopolymeric loci by modeling incidence of frame-shift indels at each locus as a function of

202    their base-pair length (see **Methods**). Six single-nucleotide repeat loci (out of seventy three loci

203    of 5 or more base-pairs in length) in *MT-ND1* (c.3,566-3,571, *n* = 32), *MT-ND4* (c.10,947-

204    10,952, *n* = 25; c.11,032-11,038, *n* = 34; and c.11,867-11,872, *n* = 50), and *MT-ND5* (c.12,385-

205    12,390, *n* = 23 and c.12,418-12,425, *n* = 73) accumulated mutations at a rate above null

206    expectation (*Q*-value<0.01, **Fig. 2f**). Homopolymer hotspots only arose at single-nucleotide loci

207    of at least 6 nt in length (*P*=0.0002, two-sided Fisher's exact test), were composed of A or C

208    homopolymer repeats, and exclusively encoded subunits of complex I. Importantly, other

209 homopolymers of equivalent length (≥6) and nucleotide content exist both in complex I and

210 complex III/IV/V but did not exhibit recurrent mutations, indicating a high degree of specificity to

211 hotspot positions (**Fig. 2g**). These six homopolymeric repeat loci collectively accounted for 40%

212 of all truncating variants observed in our data (95% binomial CI: 36-44%), and 57% (95% CI:

213 52-62%) of frame-shift indels overall. Notably, recurrent loss-of-function frameshift indels have

214 been observed at these sites as early driver mutations in rare, often benign renal oncocytomas

215 [17]; however we observed mutations at these loci to be a pervasive phenomenon across tumor

216 lineages (**Supplementary Fig. 2a**). Homopolymeric hotspot mutations arose in the PCAWG

217 cohort (after excluding any samples overlapping with our cohort) at a rate highly consistent with

218 the TCGA cohort (Pearson's $r = 0.95$), indicating that the indels detected in TCGA at hotspot

219 loci were not artifacts due to calling variants in microsatellite regions with poor coverage

220 (**Supplementary Fig. 2b**). Moreover, the three most prevalently mutated of these homopolymer

221 loci in our dataset (c.11,032-11,038, c.11,867-11,872, c.12,418-12,425) intersected with 100-bp-

222 long windows enriched for frameshift indels identified in an analysis of 616 pediatric and 2,202

223 adult tumors (527 of which were from TCGA), highlighting the power of our approach to resolve

224 focal, recurrent alterations [18]. Although mutations at homopolymeric tracts have not been widely

225 described in the germline literature, the most recurrent hotspot (*MT-ND5* c.12,418-12,425) has

226 been previously reported as the site of a germline frame-shift deletion (A12425del) in a pediatric

227 patient, where the *de novo* heteroplasmic deletion resulted in mitochondrial myopathy and renal

228 failure[19].

229

230 *Non-synonymous mutations and RNA variants arise as rare recurrent alleles with elevated*

231 *pathogenicity*

232 The bulk of somatic variants we observed in mtDNA were non-truncating, non-synonymous

233 mutations, including missense mutations, in-frame indels, translation start site mutations and

234 non-stop mutations (collectively variants of unknown significance or VUS, 73.2% of $n$=4,381

235 variants, **Fig. 3a**). Interestingly, non-synonymous variants were again depleted in CV relative to

236 other complexes, suggesting that CV is intolerant both to truncating variants and to presumably

237 less-disruptive non-synonymous mutations. Using the APOGEE framework to evaluate the

238 functional consequence of mutations to protein-coding mtDNA genes [20], we found that somatic

239 VUSs were twice as likely to be predicted pathogenic compared to germline polymorphisms

240 observed among ~200K normal samples from the HelixMTdb dataset (39.5% of somatic-only

241 variants compared to 20.4% of germline-arising, $P$=6x10^{-14}, two-sided Wilcoxon rank sum test,

242 **Fig. 3b**). Furthermore, when considering all possible mtDNA variants excluding germline

243  polymorphisms (*i.e.* the complete set of all possible somatic variants), VUSs observed in tumors

244  were more pathogenic than the set of possible somatic variants which never arose in tumors,

245  suggesting that somatic VUSs are more pathogenic than expected by random chance. We next

246  evaluated the tendency for VUSs to target specific complexes of the ETC (this necessarily

247  reduced the types of VUSs to protein-coding variants, including missense, in-frame indels, and

248  a small number of translation start site and nonstop mutations).  In contrast to truncating

249  variants, protein-coding VUSs were most frequent in CIII ($P$=1x10$^{-7}$ for least significant

250  comparison, two-sided Poisson test, **Fig. 3c**), whose functional integrity as a site for ubiquinol

251  oxidation has recently been described as essential for tumor cell proliferation[21], although as with

252  truncating variants VUSs to CV subunits were still depleted compared to the other complexes

253  ($P$=0.01 for least significant comparison). These observations were validated using data from

254  PCAWG (**Fig. 3d**). Together, these findings suggest that tumors preferentially accumulate

255  somatic missense mtDNA mutations in a manner dictated by OXPHOS complex, possibly driven

256  by their capacity to disrupt mitochondrial function due to their elevated pathogenicity.

257  Furthermore, they support the hypothesis that a purifying selection exists against variants (both

258  truncating and VUSs) that compromise physiological function of complex V/ATP Synthase.

259

260  Single nucleotide variants (SNVs) were far less recurrent than homopolymer indels ($P$=0.01,

261  two-sided Wilcoxon rank sum test among distinct variants mutated in >=3 tumors, **Fig. 3e**).

262  However, we nevertheless observed a small number of loci with recurrent non-truncating

263  variants. recurrent mutant loci. We developed a statistical test for recurrence of these loci, and

264  identified 7 SNV hotspots in the mitochondrial genome ($Q$<0.01, **Fig. 3f**), including 3 in protein-

265  coding genes (all in complex I), 3 in ribosomal RNAs (all in *MT-RNR2*), and 1 in a tRNA (*MT-*

266  *TL1*) (**see Methods**). In contrast to the high fraction of truncating mutations which are explained

267  by a relatively small number of hotspot alleles, hotspot SNV mutations collectively accounted for

268  1.6% of all VUSs; the vast majority of VUSs were non-recurrent, usually arising in a single

269  sample. Furthermore, 0/33 mutations arising at the three protein-coding hotspot positions were

270  nonsense mutations introducing an early stop codon, suggesting either the mutagenic

271  mechanism generating homopolymeric indel hotspots has a high degree of specificity, or that

272  truncating hotspots themselves may engender unique phenotypes beyond conventional loss-of-

273  function.

274

275  Mitochondrial tRNAs (mt-tRNAs) are commonly mutated in the context of germline mitochondrial

276  disease. Interestingly, the somatic hotspot *MT-TL1*$^{A3243G}$ (somatically mutated in 6 patients) is

277     also the causative variant of around 80% of MELAS disease cases and approximately 30% of

278     all mtDNA disease [22,23]. We additionally observed mutations clustered in adjacent positions

279     3242 ($n = 5$) and 3244 ($n = 4$, recently described as a recurrent mutation in Hürthle cell

280     carcinoma of the thyroid [24]), suggesting that recurrent mutations in *MT-TL1* could affect a

281     common secondary structure element. Mitochondrial tRNAs adopt a relatively conserved

282     cloverleaf structure upon folding, and mutations to mt-tRNAs are known to disrupt the function

283     of specific secondary structure elements. We therefore sought to test whether any positions of

284     the tRNA cloverleaf structure were enriched for somatic mutations in tumors. We aligned all

285     tRNA mutations according to their position in the canonical mitochondrial tRNA structure and

286     developed a statistical approach to identify enrichment in specific secondary structure elements

287     (see **Methods**). This analysis identified position 31 in the anti-codon stem of the folded tRNA

288     molecule as a site of recurrent mutation across mt-tRNAs ($Q=4.7 \times 10^{-4}$, **Fig. 3g**), which we

289     further validated using the non-TCGA subset of PCAWG samples ($Q=0.014$, **Supplementary**

290     **Fig. 3a**). Interestingly, position 31 was observed to be mutated at an 8-fold higher rate in tRNAs

291     encoded on the light-strand (*e.g. MT-TC*, $n=5$; *MT-TP*, $n=4$; *MT-TA*, $n=3$) compared to heavy-

292     strand-encoded tRNAs ($P=2 \times 10^{-4}$, two-sided Fisher's exact test). As a group, mutations at

293     structural position 31 were predicted to be more pathogenic by MITOTIP relative to mutations at

294     other tRNA positions (**Fig. 3h**), and in the case of *MT-TA*$^{T5628C}$ ($n=3$) are associated with the

295     mitochondrial disease chronic progressive external ophthalmoplegia (CPEO) [25]. In analogy to

296     the recurrent mutation of conserved amino acid residues in domains of homologous proteins [26]

297     or within 3-dimensional regions of folded protein structures [27], these data suggest that specific

298     structural features of mt-tRNAs may undergo recurrent mutation and impair normal

299     mitochondrial physiology.

300

301     To understand the potential function of rare protein-coding SNV hotspots in mtDNA, we focused

302     on a recurrent mutation at *MT-ND1*$^{R25}$, which was identified somatically in 11/10,132 TCGA

303     patients (0.11%), and 5/2,836 PCAWG patients (0.18%). All 16 instances resulted in a

304     substitution of arginine (R) with glutamine (Q), encoded by a G>A substitution at position 3380.

305     *MT-ND1*$^{R25Q}$ was previously described in a case report as the causative variant in the

306     development of MELAS in a mitochondrial disease patient [28], but was never observed among

307     ~200K normal samples, where the mutant alleles at residue R25 always produced synonymous

308     mutations (A3381G, $n=57$). Residue R25 is conserved across vertebrates [28], and is part of a

309     cluster of charged residues in complex I which form a structural bottleneck in the ubiquinone

310     binding tunnel leading to the Q binding site[29]. This led us to hypothesize that the R25Q mutation

311  could potentially disrupt the site, impacting ubiquinone : complex I binding kinetics and/or Q-site

312  substrate specificity, impeding the downstream electron transport chain. We therefore modelled

313  the effect of $MT\text{-}ND1^{R25Q}$ using a recent, high resolution structure of the mammalian. This

314  analysis highlighted changes to the local charge environment due to loss of the relatively bulky,

315  positively charged arginine sidechain. Due to the location of this substitution within the Q

316  binding tunnel, this is predicted to significantly impact function (**Fig. 3i**). Focusing on colorectal

317  tumors, which demonstrated the largest numbers of tumors harboring $MT\text{-}ND1^{R25Q}$ (*n*=8 tumors

318  total), we examined whether the presence of $MT\text{-}ND1^{R25Q}$ was associated with a particular

319  transcriptional signature. Relative to mtDNA-wild-type tumors, we observed that $MT\text{-}ND1^{R25Q}$

320  tumors were characterized by upregulation of MYC targets and oxidative phosphorylation, and

321  downregulation of gene signatures associated with hypoxia, IL2/STAT5 signaling, TNFα

322  Signaling via NFκB (**Fig. 3j**). These data suggest that $MT\text{-}ND1^{R25Q}$ promotes a transcriptional

323  phenotype characterized by increased mitochondrial metabolism and suppressed expression of

324  innate immune genes.

325

326  *Mitochondrial genotype underlies a lineage-agnostic transcriptional program*

327  Given the lineage specificity underlying both truncating variants and truncating/SNV hotspots,

328  we studied the overall burden of distinct classes of mtDNA variants (*i.e.* producing a truncating,

329  missense, synonymous, tRNA or rRNA variant) across cancer types. Restricting our analysis to

330  well-covered samples including coverage over all homopolymeric hotspots (see **Methods**), we

331  found that the fraction of mutant samples across cancer types ranged from approximately 23%

332  of leukemias (95% binomial CI: 13-35%) to as high as 80% of thyroid cancers (95% CI: 63-92%)

333  (**Fig. 4a**). Moreover, we observed no correlation between the fraction of well-covered samples in

334  a cancer type and the proportion of samples with a somatic mtDNA mutation (**Supplementary**

335  **Fig. 4a)**, indicating that the highly variable incidence of different somatic variants across cancer

336  types was not biased by their differing sequencing coverages. This extensive variation suggests

337  tumor lineages may be subject to different degrees of selection for or against mtDNA mutations,

338  consistent with the extensive variability of dN/dS ratios previously described in somatic mtDNA

339  mutations derived from whole genome sequencing of the TCGA [5].

340

341  Truncating mtDNA mutations approaching homoplasmy (>90% heteroplasmy) were identified in

342  nearly all cancer types, despite the tendency for stromal and immune cell infiltration to suppress

343  apparent tumor cell heteroplasmy, suggesting that even cancers in which mtDNA mutations are

344  uncommon may still contain rare instances of individual tumors with highly mutant mitochondria.

345    In renal and thyroid tumors, truncating mtDNA mutations have historically been associated with

346    the development of oncocytic neoplasia, whereby tumor cells accumulate dysfunctional

347    mitochondria [30,31]. That truncating mutations induce a morphologically similar response in two

348    different tissue lineages suggests that cells may adopt a lineage-agnostic adaptation to the

349    presence of a truncating mutation.  To evaluate if truncating mutations induced functionally

350    similar consequences across different tumor lineages, we compared the gene expression

351    profiles of tumor samples with truncating mtDNA variants to tumor samples with wild-type

352    mtDNA (harboring no nonsynonymous somatic mutations in protein-coding or RNA genes, see

353    **Methods, Fig. 2f**). In half of all cancer types, tumors harboring truncating mutations exhibited a

354    conserved expression program characterized by upregulation of genes associated with

355    oxidative phosphorylation and downregulation of genes associated with TNFα via NFκB

356    signaling (**Fig. 4b** and **Supplementary Fig. 4b**). Critically, these expression programs were

357    evident in cancer types such as glioma and mesothelioma, where the proportion of samples with

358    a truncating variant was comparatively low (**Fig. 4c**). These data suggest that, even in cancer

359    types where mtDNA mutations are rare, truncating mtDNA mutations produce similar phenotypic

360    outcomes.

361

362    Given that the hotspot $MT\text{-}ND1^{R25Q}$ exhibited an expression program resembling truncating

363    variants, we investigated the generic transcriptional consequences of mtDNA VUSs (see

364    **Methods**). Compared to truncating variants, fewer genesets demonstrated lineage-agnostic

365    changes in samples with VUSs. As with truncating variants, the most upregulated geneset in

366    VUS-harboring tumors was Oxidative Phosphorylation (increased in 5/18 cancer types)

367    (**Supplementary Fig. 4c**), but the magnitude of this enrichment was attenuated relative to

368    truncating variants. Notably, several cancer types, such as colorectal cancer, demonstrated a

369    lineage-specific pattern of gene expression changes, suggesting that mtDNA VUSs are capable

370    of eliciting a phenotype in specific cancer types.

371

372    To examine the translational value of mtDNA genotype, we determined the association between

373    mtDNA mutation status and clinical outcome (overall survival) across cancer types. Using

374    univariate Cox proportional-hazards regression, for each cancer type we determined the effect

375    size and significance of both mtDNA truncating variants and VUSs compared to samples with

376    no somatic mtDNA variants (wild-type). Colorectal cancer demonstrated the largest (by effect-

377    size) significant association between overall survival time and mtDNA genotype (colorectal

378    patients with VUSs had a hazard ratio of 0.47 (95%CI 0.03-0.75) compared to those with wild-

379   type mtDNA, *Q*-value=0.02, Cox proportional-hazards regression) (**Fig. 4d**). Notably, VUSs in

380   colorectal cancer also associated with a unique transcriptional down-regulation of multiple

381   genesets including TNFα via NFκB, Hypoxia and Complement (**Supplementary Fig. 4c, Fig.**

382   **3j**), further suggesting a cryptic phenotype of these variants in affected tumors. We additionally

383   observed a weak association between mitochondrial genotype and underlying molecular

384   subtype [32], with some enrichment of mtDNA mutations in the canonical subtype CMS2 of

385   colorectal tumors(**Supplementary Figure 4d**). We therefore further evaluated if mtDNA

386   mutations may be prognostically meaningful in colorectal cancer, using a multivariate analysis to

387   control for known prognostic clinical and genomic covariates. Among 344 stage 1-3 colorectal

388   cancer patients, the presence of mtDNA alterations was significantly associated with better

389   overall survival compared to wild-type samples (*P*=0.002, Kaplan–Meier test), with patients

390   whose tumors harbored VUSs having the best prognosis, and those with truncating variants

391   having an intermediate improvement (**Fig. 4e**). This association remained significant after

392   controlling for clinically-relevant prognostic covariates (*i.e.* age, cancer stage, primary site, MSI-

393   status, consensus molecular subtype and the presence of established nuclear-encoded

394   genomic driver mutations) [32,33] in a multivariate analysis. VUSs again had a significantly

395   protective association compared to wild-type (Hazard ratio=0.18, 95%; CI: 0.08-0.44; *Q*-

396   value=0.001, Cox proportional-hazards model); truncating variants had an intermediate effect

397   (HR=0.38, 95% CI: 0.15-0.97; *Q*=0.18) (**Fig. 4f**). These data together suggest that somatic

398   mtDNA mutations are associated with a clinically and molecularly-distinct class of colorectal

399   tumors, and that the functional consequence of an mtDNA mutation is a determinant of its

400   clinical significance.

401

402

403   **Discussion**

404   Although recent evolutionary data suggests that mtDNA mutations may be under positive

405   selection in cancers of the kidney and thyroid [5], the broader significance of somatic mtDNA

406   mutations in cancer remains a point of confusion and debate. Drawing inspiration from analyses

407   describing hotspots of somatic mutations in the nuclear DNA of tumors, we studied the

408   recurrence of mutant mtDNA alleles. The discovery that OXPHOS complex shapes mtDNA

409   mutation patterns in a manner that produces mutation hotspots, in connection with orthogonal

410   data on the structural consequences, transcriptomic effects and clinical significance of these

411   alleles in patients with germline mtDNA disease, supports the hypothesis that mitochondrial

412   respiration is perturbed across many tumors.

413

414    Our results indicate that OXPHOS complex, tissue lineage, and mutation consequence

415    collectively shape the incidence and putative function of mtDNA mutations. Whereas previous

416    studies have demonstrated localized regions of mtDNA with elevated somatic mutation rate in

417    tumors, these works have generally been underpowered to probe phenotypic differences

418    between alleles. Our data reveal that truncating mutations preferentially impact complex I, and

419    that non-synonymous mutations of all classes are depleted in complex V. This suggests that

420    cancer cells can better tolerate, or perhaps even utilize, loss of complex I and the associated

421    metabolic consequences (e.g. NAD+:NADH changes), whereas loss of capacity for ATP

422    synthesis through complex V mutations appears to be negatively selected against. That CIII

423    demonstrates elevated rates (relative to other complexes) of missense mutations, but not

424    truncating mutations, is consistent with its essential role in ubiquinol oxidation and suggests that

425    weak disruption of CIII is preferential for clonal expansion in tumor cells [21]. Whether truncating

426    mutations in CIII and CIV promote different phenotypes in cancer cells relative to complex I loss

427    warrants further investigation.

428

429    There is substantial evidence that in particular subtypes of thyroid and kidney cancer, mtDNA

430    mutations are the root cause of metabolic adaptations and morphological (oncocytic) changes

431    associated with suppression of mitochondrial respiration [34]. What remains unclear is how to

432    extrapolate the function of truncating mutations in otherwise essential mtDNA genes to cancer

433    types where oncocytic tumors are rarely if ever observed but the fraction of samples harboring

434    these mutations is nevertheless substantial (*e.g.* colorectal cancers). Critically, our

435    transcriptional data suggests that, even in cancer types where truncating mtDNA mutations are

436    rare, they nevertheless promote a transcriptional program characterized by increased

437    expression of OXPHOS genes and downregulation of innate immune pathways. Because

438    homoplasmic loss of any gene in the mtDNA necessarily cripples the cell's ability to respire and

439    disrupts dependent metabolic pathways, these findings suggest that pathogenic and high

440    heteroplasmy mtDNA mutations potentially render a large fraction of tumors vulnerable to a

441    metabolic therapeutic intervention.

442

443    **Methods**

444

445    **Tumor and normal sample sequencing cohorts**

446    Tumor and matched normal sequencing data for TCGA samples were obtained from the GDC

447    Data Portal (https://portal.gdc.cancer.gov/). Briefly, all tumor and matched-normal barcodes

448    included in the MC3 MAF [35] (https://gdc.cancer.gov/about-data/publications/pancanatlas) file

449    were converted to UUIDs using the TCGAutils R package (v1.9.3), and these UUIDs were

450    queried for whole-exome sequencing BAM files sliced for chrM using the GDC API. We then

451    queried the GDC Data Portal for RNA-Sequencing BAM files for TCGA tumors already with

452    whole-exome sequencing data. This process yielded paired tumor and matched-normal whole-

453    exome sequencing BAMs for 10,132 TCGA patients, of which 9,455 had additional RNA-

454    sequencing data. In addition to the raw sequencing data for TCGA samples from which we

455    called mtDNA mutations (see: Calling mitochondrial variants), we additionally obtained somatic

456    mitochondrial mutation calls for 2,836 whole-exome sequenced tumors from ICGC/PCAWG [3], of

457    which 885 also had TCGA sequencing data. Nuclear somatic mutations for TCGA samples were

458    obtained from the MC3 MAF, subset for the samples for which mtDNA whole-exome

459    sequencing BAMs were available. Finally, mtDNA mutation calls for 195,983 normal samples

460    were obtained from the HelixMTdb cohort of sequenced saliva samples from healthy individuals

461    [14].

462

463    **Annotating mtDNA regions included in our analysis**

464    Each mitochondrially-encoded gene's name, start/end positions and DNA strand was obtained

465    from Biomart for human reference genome GRCh38 (release 95). Subsequently, each mtDNA

466    position (1-16569) was annotated with its associated genetic information. Any mtDNA positions

467    located at the overlap of two genes were annotated only as associated with whichever gene

468    started first in numerical genomic position. Variants in non-genic mtDNA regions were excluded

469    in our analyses. To this end, we excluded any variants in the mtDNA Control Region (positions

470    1-576, 16,024-16,569) as well as 89 other non-genic positions. We similarly excluded variants in

471    hypermutated regions of mtDNA, including 302-316, 514-524, and 3,106-3,109). Following

472    these measures, the genomic length of mtDNA retained in our analyses was 15,354bp. (The

473    complete list of 16,569 mtDNA positions and their annotated reasons for exclusion is provided in

474    **SI table 1**.)

475

476    **Calling mitochondrial variants**

477    Mutations to the mitochondrial genome were obtained from variants called by both of two

478    independent variant-calling pipelines. In the first pipeline, Mutect2 (GATK v4.1.2.0) [36] was used

479    to call variants in chrM in tumor and normal samples individually, the results of which were

480    subsequently intersected to obtain variants called supported in a given patient's tumor and

481    matched normal samples. Briefly, Mutect2 was run in mitochondrial-mode for each patient's

482    tumor and normal sample independently against human reference genome GRCh38 (with

483    minimum base quality-score 20, minimum mapping quality 10, aggressive pcr-indel model, and

484    other standard quality control arguments for paired-end reads). Artifacts were subsequently

485    removed using GATK *FilterMutectCalls* (GATK v4.1.2.0) [36] , and multi-allelic sites were split into

486    individual variants using the *norm* function from bcftools (v1.9) [37] . The resulting tumor and

487    normal VCF files were then merged using gatk *HaplotypeCaller* (GATK v4.1.2.0) [36] , to annotate

488    variants in the tumor VCF with their coverage in the normal sample. The resulting VCF was

489    converted to a MAF file using vcf2maf (v1.6.17, https://github.com/mskcc/vcf2maf). Finally,

490    variants from the generated MAF file were then filtered out unless the variant allele was

491    supported at least one read in both forward and reverse directions. In the second pipeline,

492    samtools mpileup (v1.9) [11] was used to generate a pileup file using variant-supporting reads

493    with minimum mapping quality 20 and base alignment quality 10. Reads failing quality checks or

494    marked as PCR duplicates were removed. Variants were required to contain at least 2 variant-

495    supporting reads in the forward and reverse direction. In each pipeline, variants were

496    additionally filtered to ensure ≥ 5% variant allele frequency in the tumor, and tumor coverage ≥ 5

497    reads. Variants identified by both pipelines were retained for further analysis. In rare cases,

498    multiple indels were called in a sample within a homopolymeric region (single-nucleotide

499    repeats of 5 or more basepairs), with distinct alt-read counts and VAF values, and identical

500    read-depth values. These multiple indels were collapsed to a single representative indel call.

501    Briefly, using the Mutect2 variant calls, whichever indel had the highest VAF in the tumor

502    sample was taken as the representative indel. The count of alt-reads in both tumor and normal

503    were replaced with their corresponding summed counts across the original multiple indels, and

504    the VAFs in both tumor and normal were re-calculated from the new summed alt-read counts

505    divided by the original read-depth. Finally, mutations were classified as of somatic origin

506    according to the following criteria: Non-truncating variants (that is, all variant classifications

507    other than nonsense mutations and frame-shift indels) were classified as somatic if the matched

508    normal sample had a minimum coverage of 5 reads and 0 normal reads called the alternate

509    allele. Truncating variants in the tumor sample were assumed to be of somatic origin. All other

510    variants were not classified as somatic and excluded from this study.

511

512    **Nuclear mutational data and annotation**

513     Somatic mutations in nuclear-encoded cancer-associated genes for TCGA samples were

514     obtained from the PanCanAtlas MC3 MAF file. Mutations in this file were subset for those

515     among the 468 genes on the MSK-IMPACT clinical sequencing panel [13]. The MAF file was

516     annotated for known, likely, and predicted oncogenic driver mutations using the MAF-Annotator

517     tool provided by OncoKB [38] (https://github.com/oncokb/oncokb-annotator). Mutations annotated

518     by OncoKB as "Oncogenic", "Likely Oncogenic" or "Predicted Oncogenic", previously

519     determined cancer hotspot mutations [15,16], or truncating variants to tumor suppressor genes (*i.e.*

520     frame-shift indels, splice-site and nonsense mutations) were classified as potential driver

521     alterations.

522

523     **Calculating tumor mutational burden in mtDNA or nuclear DNA**

524     Tumor mutational burden (TMB) was calculated for cohorts of tumors subset for various

525     genomic regions, including: 1) individual mitochondrial- or nuclear-encoded genes; 2) mtDNA

526     genes grouped by OXPHOS complex I, III, IV, or V; 3) the entire mitochondrial genome

527     (excluding non-genic and polymorphic regions); 4) a set of known nuclear-encoded tumor

528     suppressor genes; and 5) a set of known nuclear-encoded oncogenes. In each case, TMB was

529     calculated as the total number of somatic mutations among the relevant collection of tumors

530     divided by the total genomic length sequenced in these tumors (in Mbps). For TMBs calculated

531     from mutations called in off-target sequencing data (*i.e.* mtDNA variants in TCGA samples), the

532     total genomic length sequenced was the number of the genomic positions with sufficient

533     coverage to call somatic variants (5+ read coverage in both tumor and normal sample),

534     summed across all samples. For TMBs calculated from targeted regions (nuclear DNA; mtDNA

535     in PCAWG samples), the total genomic length sequenced was the length of the targeted region

536     (entire gene for mtDNA, exonic regions for nuclear DNA) multiplied by the number of samples.

537     Error bars for TMBs were calculated as 95% Poisson exact confidence intervals for rates, using

538     the total number of mutations as the count of events, and the genomic length sequenced in Mb

539     as the time at risk.

540

541     **Identifying hotspot positions for mitochondrial variants**

542     We identified mtDNA positions with statistically recurrent single-nucleotide variants (SNVs) by

543     comparing the observed proportion of mutations at an individual position (out of the total number

544     of mutations acquired in its gene) to a rate of mutations at the position expected by chance with

545     a one-sided binomial test. The probability for SNVs at each position of a gene $P_{pos,gene}$ was

546     modeled as a bernoulli trial, where the likelihood of a mutation arising at a given position by its

547 mutability relative to the mutability of all other bases in the gene: $P_{pos,gene} = \frac{\mu_{pos}}{\mu_{gene}}$. Consistent

548 with previous work [15], we estimated the mutability for each position as a function of its

549 trinucleotide context. That is, for each position, it's mutability $\mu_{pos}$ was calculated as the count of

550 SNVs matching the trinucleotide context of the position of interest $s_{pos}$, out of the total count of

551 SNVs anywhere in the mitochondrial genome $s_{total}$ (after excluding the control region and other

552 blacklisted regions). Due to the highly strand-specific mutation signatures we observed for

553 SNVs in mtDNA (**Supplementary Fig. 1c**), we used the complete set of 64 unique

554 trinucleotides in order to retain this information when calculating the mutability for each position,

555 rather than collapsing the central nucleotide to C or T resulting in the conventional 32 unique

556 trinucleotides. As the proportion of patients for which a given position had sequencing coverage

557 in paired tumor and normal samples linearly affects the likelihood of observing a somatic

558 mutation at the position, the mutability of a position was adjusted to control for this by

559 multiplying it by the ratio of the number of samples with paired tumor-normal sequencing

560 coverage at the position $C_{pos}$ out of the total number of samples $N_{samples}$, so that $\mu_{pos} =$

561 $\frac{s_{pos}}{s_{total}} \times \frac{C_{pos}}{N_{samples}}$. The mutability associated with the gene was calculated as the sum of each

562 position's trinucleotide mutability. Therefore, for a gene $L$ basepairs in length: $\mu_{gene} =$

563 $\sum_{pos=1}^{L} \mu_{pos}$ . The final parameter for the binomial test (i.e. the likelihood for a mutation in a gene

564 to arise at the given position by chance) was therefore $P_{pos,gene} = \frac{\mu_{pos}}{\mu_{gene}}$. Each position mutated

565 in 5 or more samples in each gene was subsequently tested for statistically enriched mutations

566 by comparing its observed number of mutations out of the total number of mutations in the gene

567 to this binomial parameter using a right-tailed binomial test. The full list of generated *P*-values

568 across all genes were then corrected for multiple hypothesis testing.

569

570 **Homopolymer hotspots for indels**

571 To identify homopolymer regions with statistically enriched rates of insertions and deletions

572 (indels), we modeled the proportion of samples with indels across all homopolymers as a

573 function of the homopolymer region's width (i.e. the number of repeated nucleotides, from 5-8).

574 To this end, all single-nucleotide repeats of 5 or more basepairs were identified in the

575 mitochondrial reference genome, resulting in N=73 unique homopolymer loci in whitelisted

576 coding mtDNA. We then modeled the fraction of frame-shift indels across 73 homopolymers

577 observed to arise at a specific homopolymer locus $h$ as a binomial process dictated by the

578 length of the homopolymer $l_h$ divided by the summed length of all homopolymers, such that the

579    expected likelihood of a frame-shift indel arising at a homopolymer by chance is given by: $p_h =$

580    $\frac{l_h}{\sum_{i=1}^{73} l_i l_i}$. We then tested each homopolymer locus for enriched mutations with a one-sided

581    binomial test. That is, for each homopolymer locus, the number of bernoulli trials was the

582    number of samples with complete sequencing coverage for the homopolymer region and two

583    flanking basepairs; the number of successes was the number of samples with frame-shift indels

584    at (or immediately adjacent to) the given homopolymer, and the fraction of successful trials was

585    compared to the expected probability $p_h$.

586

587    **Hotspot positions in tRNA cloverleaf structure**

588    Positions of the tRNA cloverleaf secondary structure were individually tested for an enriched

589    rate of SNVs at the equivalent aligned positions of the 22 mitochondrially-encoded tRNAs. A

590    map of genomic positions in mitochondrial tRNAs to cloverleaf structure positions was provided

591    by Mitotip [39]

592    (https://github.com/sonneysa/MitoTIP/blob/master/Output/tRNA%20data%20and%20scoring_sc

593    ored.xlsx) and used to assign SNVs at tRNAs to structural positions. Under the null hypothesis

594    that mutations accumulate at structurally-aligned positions randomly, the proportion of SNVs

595    aligning to a specific position in the tRNA cloverleaf should be approximately equal to the

596    number of times the aligned position was sequenced at a sufficient depth in both tumor and

597    matched normal samples to call somatic mutations, out of the total number of tRNA basepairs

598    sequenced at sufficient depth across all samples at all structural positions. Therefore for a given

599    position of the tRNA cloverleaf structure $p$, the number of SNVs observed across all tRNAs at

600    this aligned position $t_p$ out of $T$ SNVs across all positions of all tRNAs was tested for enrichment

601    using a one-sided binomial test, compared to an expected rate equal to the number of tRNA

602    bases aligned to this position sequenced at sufficient depth $b_p$ out of $B$ tRNA bases sequenced

603    at sufficient depth across all positions of all tRNAs.

604

605    **Classifying sample mtDNA variant status**

606    Each tumor sample was classified according to the presence and type of its somatic

607    mitochondrial variants. Because gaps in sequencing coverage may make existing variants

608    undetectable and result in the incorrect classification of such samples as "wild-type" for somatic

609    variants, we only attempted to classify samples with sequencing coverage in both tumor and

610    matched normal of at least 90% of the included region of mtDNA (referred to as "well-covered"

611    throughout). Furthermore, given the high incidence of truncating indels we observed at 6

612    hotspot loci, we additionally required that these 6 loci were sequenced at sufficient coverage in

613    the tumor sample, to ensure that samples potentially harboring recurrent indels would be

614    excluded and not misclassified. Samples not meeting either of these conditions were classified

615    as having 'Unknown' mtDNA mutation status. The remaining samples were then classified

616    according to a decision tree as follows: Samples with any protein-truncating variants were

617    classified as 'Truncating'; remaining samples still unclassified with multiple mtDNA variants of

618    different types (among missense, rRNA, and tRNA variants) were classified as '2+ non-

619    truncating types'; remaining samples with tRNA mutations were classified as 'tRNA'; remaining

620    samples with rRNA mutations were classified as 'rRNA'; remaining samples with non-truncating,

621    non-synonymous protein-coding mutations as 'missense'; remaining samples with silent

622    mutations as 'Silent'; and finally samples still unclassified were classified as 'wild-type'. This

623    logic prioritizes minimizing annotation bias over conserving sample size, in order to meaningfully

624    compare the incidence of different variant types across samples. However, in our analysis of the

625    effect of mtDNA variants on differential gene expression or survival, we modified the logic to

626    prioritize conservation of sample size. To this end, in RNA-Seq and survival analyses, samples

627    with any observed truncating variants were classified as truncating, regardless of their

628    sequencing coverage.

629

630    **Testing genesets for transcriptional dysregulation due to mtDNA variants**

631    A matrix of estimated gene expression counts (RSEM values normalized to correct for batch

632    effects) for TCGA samples was downloaded from the TCGA PanCanAtlas [35] supplemental data

633    (http://api.gdc.cancer.gov/data/3586c0da-64d0-4b74-a449-5ff4d9136611). Gene expression

634    estimates were rounded to integer values, and subsequently genes with zero estimated counts

635    in all samples were removed, as were genes with unknown gene symbols. To evaluate

636    differentially expressed genes between two groups of samples with different mtDNA variant type

637    (i.e. truncating vs wild-type samples colorectal samples), the rounded gene expression matrix

638    was subset for the relevant samples and input into the DESeq2 [40] package in R using the

639    DESeqDataSetFromMatrix utility, along with a table of tumor sample barcodes with their

640    associated mtDNA classification. Differentially expressed genes were tested and their log-fold

641    change (LFC) values were shrunken using the apeglm [41] package. *P*-values for all genes tested

642    were corrected for multiple-hypothesis testing with the Benjamini Hochberg method [42]. The

643    resulting data from this analysis were used to calculate a statistic for each gene equal to

644    $\log_{10}$(*Q*-value) x sign(LFC). All genesets from the mSigDB Hallmark geneset collection [43] (v7.1)

645    were then tested for significant up- or down-regulation based on this statistic for each gene

646    using the fgsea package [44] in R, with a minimum geneset size of 10 genes, a maximum size of

647    500 genes, and 100,000 permutations.

648

649    **Annotating genomic and clinical covariates in colorectal cancer survival analysis**

650    Clinical data for TCGA colorectal cancer patients including: overall survival time/status, AJCC

651    pathologic tumor stage, age at diagnosis, sex, and tumor tissue site were obtained from the

652    TCGA FIrehose legacy data on cbioportal

653    (https://www.cbioportal.org/study/summary?id=coadread_tcga). Clinical data was subset for

654    patients with sequencing data in the MC3 MAF. These data were then annotated with MSI

655    status (MSS, MSI-low, MSI-high) based on published data for patients where this was available

656    [45]. AJCC Pathologic Tumor Staging data was collapsed into Stages I, II, III, IV, and Stage-IV

657    patients were excluded. The tumor site was encoded as "Right-colon" if the primary site was:

658    ascending colon, cecum, hepatic flexure, or transverse colon; or encoded as "Left-colon" for:

659    descending colon, sigmoid colon, or splenic flexure. Patients with tumor tissue from the rectum

660    were encoded as "Rectum" for their tumor site. The clinical data for each sample was then

661    annotated for the presence of known or likely nuclear-encoded driver alterations in

662    KRAS/HRAS/NRAS, BRAF, APC, SMAD4 and TP53 as based on mutation calls from the TCGA

663    MC3 MAF [46] (see: Methods "Nuclear mutational data and annotation"). Each patient in the

664    clinical data was then annotated as having a known/likely driver alteration in each of

665    KRAS/HRAS/NRAS (grouped into RAS), BRAF, APC, SMAD4 or TP53. The complete multi-

666    variate model use in the Cox proportional-hazards regression was therefore: Overall Survival ~

667    mtDNA-status + Age + Stage + Site + RAS + RAF + APC + SMAD4 + TP53 + Sex + MSI-status

668    + CMS-type.

669

670    **Structural impact of *MT-ND1$^{R25Q}$* variant on complex I**

671    The structural impact of the *MT-ND1$^{R25Q}$* variant was investigated using an electron-microscopy

672    derived structure of mitochondrial CI in *mus musculus* (PDBID: 6G2J) [29]. The UCSF Chimera

673    software (v1.13.1) [47] was used to insert the R25Q mutation using the *swapaa* command. The

674    ubiquinone binding tunnel was predicted using the CAVER Analyst (v2.0b) [48] software run on

675    the wild-type PDB structure, starting from the side chain oxygen atom in *Ndufs2$^{Y108}$*, and using a

676    minimum probe radius of 1.4Å as described by the authors [49]. Surface electrostatic charge for

677    wild-type and mutant structures were determined using the APBS software [50]

678    (http://server.poissonboltzmann.org/pdb2pqr) using default parameters, after subsetting the

679     PDB structure for Mtnd1 (chain H), and converting the resulting PDB file to PQR using

680     PDB2PQR [51]. All structure visualizations were generated using UCSF Chimera.

681

**Statistical analyses and figures**

683     All statistical analyses were performed using the R statistical programming environment (version

684     3.6.1). Protein structure figures were generated using UCSF Chimera, Kaplan-Meier plots and

685     Cox proportional hazard forest plots were generated with the survminer library in R, ETC

686     schematic (Fig. 1a) in Adobe Illustrator. All other figures were generated using the ggplot2

687     library in R. Unless otherwise noted, error bars for proportions are 95% binomial CIs calculated

688     using the Pearson-Klopper method; error bars for rates (e.g. Mutations/Mb) are 95% Poisson

689     CIs calculated with the pois.exact function from the epitools library in R. Unless otherwise noted,

690     *P*-values for difference in proportions were calculated using Fisher's exact tests or two-sample

691     Z-tests, and for difference in rates using Poisson exact tests. *P*-values were corrected for

692     multiple comparisons using the Benjamini-Hochberg method [42] and reported as *Q*-values when

693     applicable.

694

**Data and code availability**

696     All relevant data and R code are available on GitHub with instructions to execute the code and

697     regenerate all figures (https://github.com/reznik-lab/mtdna-mutations).

698

**Acknowledgements**

700     We thank the members of the Reznik and Taylor laboratories for discussion and support. We

701     also thank Lydia Finley, Kivanc Birsoy, and Nicole Rusk for their feedback.

702

**Author contributions**

704     ANG, PAG, and ER conceived the study.MK, WKC, KL, AAH, and BST assisted with genomic

705     data analysis. ANG, PAG, and ER wrote the manuscript with input from all authors.

706

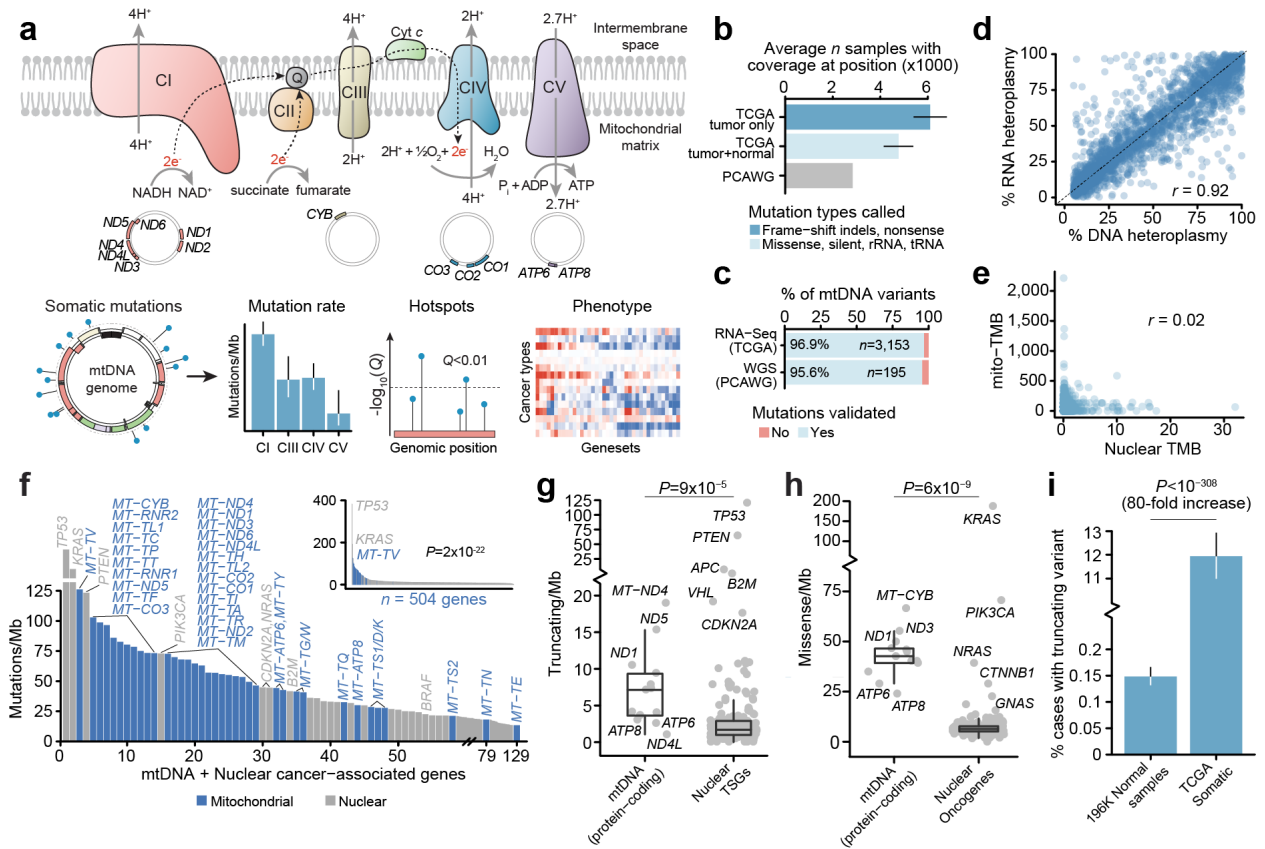**Competing financial interests**

708     The authors declare no competing financial interests

709

710

711     **Figures**

712



713

714     **Fig. 1: mtDNA mutations are among the most frequent genomic alterations in cancer. a)**

715     Schematic of oxidative phosphorylation (OXPHOS) system and project workflow. Top row,

716     complexes I-V and their reactions. Center row: mtDNA genomic regions encoding protein

717     subunits of the associated OXPHOS complex. Bottom row, overview of project workflow, in

718     which somatic mutations in mtDNA genes are used to explore inter-complex differences,

719     mutational recurrence and transcriptional phenotype associated with mitochondrial dysfunction.

720     **b)** Average number of tumors with sufficient coverage to call variants at a mtDNA position.

721     Truncating mutations were assumed to be somatic and therefore allowed for tumor-only variant-

722     calling (dark blue), whereas non-truncating (protein-coding non-truncating, tRNA and rRNA

723     mutations) required sufficient coverage in both tumor and matched normal samples (light blue).

724     Gray, the number of whole-genome sequenced (WGS) samples from PCAWG for comparison.

725     **c)** The percentages of variants called from off-target reads which were validated in either RNA-

726     Seq or WGS data from the same tumors. **d)** The correlation between variant heteroplasmy as

727     observed in RNA and DNA-sequencing ($n$=2,575 mutations with coverage ≥30 reads in both

728     DNA and RNA). **e)** The correlation between tumor mutation burden (TMB, Mutations/Mb) among

729    mtDNA (Y-axis) and nuclear-encoded cancer-associated genes (referred to simply as cancer

730    genes) (X-axis), *n*=3,624 well-covered pan-cancer tumors. **f)** Mutation rates (Mutations/Mb) of

731    individual mtDNA-encoded genes (blue) and nuclear-encoded cancer-associated genes (gray).

732    Inset plot: mutation rates among 504 genes with mtDNA genes highlighted. Outer plot: closeup

733    of the inset plot in the region containing all 37 mtDNA genes; commonly-mutated nuclear cancer

734    genes in this region are labeled for reference. **g)** Comparison of truncating mutation rates

735    (truncating variants/Mb) between 13 mtDNA-encoded protein-coding genes and 185 nuclear-

736    encoded TSGs. **h)** Comparison of non-truncating mutation rate (nonsynonymous, non-

737    truncating variants/Mb) between 13 mtDNA protein-coding genes and 168 nuclear oncogenes. **i)**

738    Percentage of patients with truncating mtDNA variants either somatically (in TCGA tumor

739    samples) or germline (among ~200K normal samples).

740

741

743 **Fig. 2: Truncating variants preferentially target complex I. a)** Comparison of truncating

744 mutation rate (truncating variants/Mb) between OXPHOS complexes I, III, IV, V. Synonymous

745 mutation rates shown below for comparison. Truncating mutations *n*=352; synonymous *n*=475.

746 *P*-values from two-sided Poisson-exact. Single asterisk denotes *P*<0.1; double asterisk *P*<0.01;

747 n.s., not significant. **b)** Validation of analysis in a) using data from n=1,951 whole-genome

748 sequenced tumors from ICGC/PCAWG after removing samples also in TCGA. Truncating

749 mutations *n*=198; synonymous *n*=263. *P*-values and asterisks as in a). **c)** Distributions of

750 truncating and silent mutation heteroplasmy (estimated by variant allele frequency) among

751 variants in OXPHOS complex I, III, IV, or V. Difference in heteroplasmy between truncating and

752 silent mutations calculated by two-sided Wilcoxon rank sum test. CI, $P=1\times10^{-6}$, not significant for

753 other complexes. **d)** Percentage of tumors with truncating mtDNA variants per cancer type,

754 among well-covered samples. Right, number of well-covered samples per cancer type. **e)**

755 Percentage of samples per cancer type with truncating variants affecting OXPHOS complex I or

756 III-V. Asterisk indicates cancertypes with enriched truncating variants targeting CI compared to

757 CIII-V, *Q*<0.01, two-sided McNemar's test. **f)** Circular mtDNA genome annotated with 73

758 homopolymer repeat loci ≥5bp in length. Dot height from the circular mtDNA genome indicates

759 the number of affected samples, dot color indicates the identity of the repeated nucleotide (A, C,

760 G, T), dot width indicates the length of the repeat region (5-8bp). Includes putatively somatic

761 truncating variants with tumor-only sequencing coverage. The 6 solid-color homopolymer loci

762 highlighted were found to be statistically enriched hotspots for frame-shift indels in tumors. **g)**

763 The 73 homopolymer repeat loci arranged by gene and repeat size. Dot width indicates –

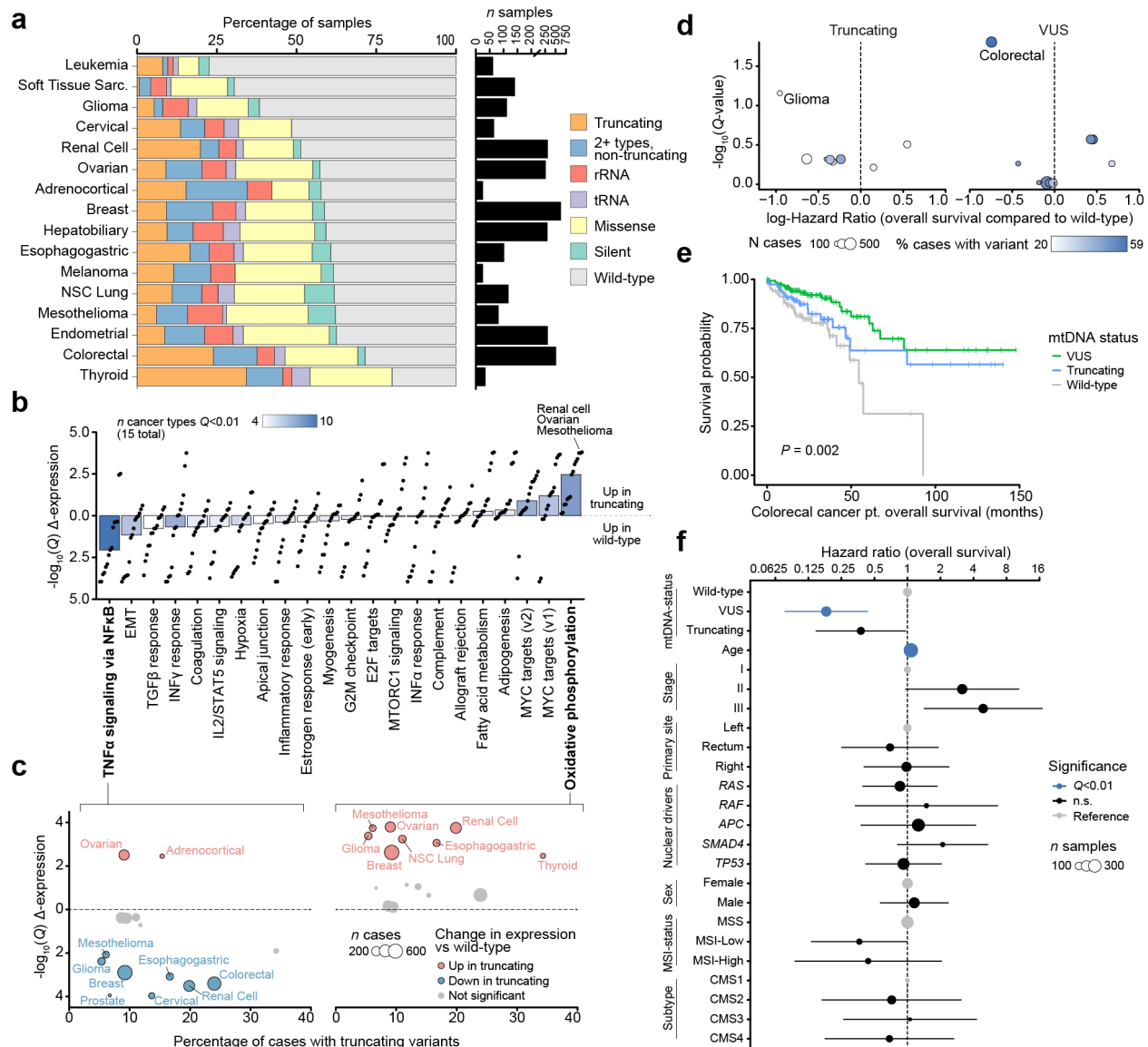764 $\log_{10}(Q$-value) for enriched frame-shift indels in tumors. The 6 hotspot loci are labeled.

765

**Fig. 3: Non-truncating mtDNA mutations arise as rare recurrent alleles in protein-coding and RNA elements. a)** The proportion of truncating, synonymous and VUS somatic mtDNA mutations in this study. VUSs are further classified into missense protein-coding variants, or mutations to rRNA or tRNA genes. **b)** Comparison of the percentage of unique VUSs predicted to be pathogenic by APOGEE [20] between somatic variants which (1) were ever observed to also

772    arise as germline variants among ~200K normal samples from HelixMTdb; (2) were never

773    observed somatically mutated; or (3) were only observed as somatic mutations. All indicated

774    comparisons were statistically significant with $P<10^{-4}$. **c)** Comparison of missense mutation rate

775    (missense variants/Mb) between OXPHOS complexes I, III, IV, V. Synonymous mutation rates

776    shown below for comparison. Missense mutations $n$=1,718; synonymous $n$=475. $P$-values from

777    two-sided Poisson-exact. Single asterisk denotes $P<0.1$;  double asterisk $P<0.01$; n.s., not

778    significant. **d)** Validation of analysis in a) using data from n=1,951 whole-genome sequenced

779    tumors from ICGC/PCAWG after removing samples also in TCGA. Missense mutations $n$=1073;

780    synonymous $n$=263. $P$-values and asterisks as in a). **e)** Rare-recurrent alleles are primarily non-

781    truncating variants. Top portion, number of samples with the given mtDNA mutant allele in

782    decreasing order of prevalence (mutations called in samples with adequate tumor and normal

783    coverage). Bottom portion, tracks indicate consequence of corresponding variant in top portion.

784    **f)** Individual base-pair positions in mtDNA with somatic single-nucleotide variants (SNVs) in ≥5

785    tumors, and their statistical enrichment for mutations. Hotspot positions with $Q<0.01$ are colored

786    by the type of gene in which they arise (protein-coding, rRNA or tRNA). Select hotspots are

787    labeled with their genomic positions (for mutations in tRNAs and rRNAs) or residue (protein

788    coding genes). **g)** Prevalence of SNVs in tRNA genes, aligned to their positions in the folded

789    tRNA cloverleaf structure. Bottom portion, number of samples with SNVs at the given tRNA

790    cloverleaf position across all tRNAs. Top portion, statistical enrichment for the aligned position

791    for mutations. **h)** Mutations at tRNA cloverleaf structural position 31 have greater predicted

792    pathogenicity scores (based on MitoTIP [39]) compared to all possible mutations at other

793    positions. tRNA mutations at position 31 affecting ≥2 samples are highlighted. $P$-value from two-

794    sided Wilcoxon rank sum test. To reduce image size, a random selection of 5% of the mutations

795    not at position 31 are plotted ($P$-value based on the complete set of mutations). **i)** The Mtnd1

796    R25Q mutation lies at a critical region of complex I near the entrance to the ubiquinone binding

797    tunnel (dotted green path), likely affecting its capacity for binding ubiquinone. Larger view: The

798    complete mammalian complex I structure (gray) highlighting Mtnd1 (blue), and the ubiquinone

799    binding tunnel (green) and binding site (large green sphere); black box indicates the region in

800    the closeup view. Closeups, the predicted surface electrostatic potential of Mtnd1 (top) wild-type

801    and R25Q mutant (bottom), proximal to the ubiquinone binding tunnel (green), leading to its

802    binding site at Ndufs2 Y108. **j)** Differentially expressed mSigDB Hallmark genesets between

803    colorectal tumors with *MT-ND1* R25Q and those without non-silent somatic mtDNA variants (*i.e.*

804    wild-type). Normalized enrichment score (NES) and adjusted P-values based on gene set

805    enrichment analysis (GSEA) using the fgsea R package [44].

806



807

808 **Fig. 4: Mitochondrial genotypes associate with transcriptional and clinical phenotypes. a)**

809 Percentage of well-covered tumors with different types of somatic mtDNA variants per cancer

810 type. Right, number of well-covered samples per cancer type. **b)** Differential expression of

811 mSigDB Hallmarks genesets, between samples with truncating mtDNA variants and those with

812 no nonsynonymous somatic mutations (*i.e.* "wild-type" samples). Differential expression is

813 quantified by directional -$\log_{10}$(Q-value): greater than 0 denotes up-regulation in samples with

814 truncating variants, below 0 denotes up-regulation in wild-type samples. Each dot is a single

815 cancer type's level of dysregulation of that geneset. Bars show the median level of

816 dysregulation across 15 cancer types; bar shading shows the number of cancer types with

817 significant dysregulation (*Q*<0.01) in either direction. **c)** Differential expression of TNFα via

818    NFκB Signaling (left) and Oxidative Phosphorylation (right) genesets in individual cancer types.

819    X-axis shows the overall proportion of samples of each cancer type with truncating variants; Y-

820    axis matches the Y-axis in b). Dot width denotes number of well-covered samples for each

821    cancer type. **d)** Effect size and statistical significance of mtDNA truncating variants (left) and

822    VUSs (right) on overall survival among individual cancer types. Effect sizes (quantified as log-

823    hazard ratios) are from univariate Cox proportional-hazards models run for each cancer type

824    independently. *Q*-values are adjusted *P*-values from the model coefficients for each cancer

825    type.

826    **e)** Kaplan–Meier plot showing difference in overall survival time among *n*=344 TCGA colorectal

827    cancer patients with somatic VUSs (*n*=152), truncating variants (*n*=84), or no nonsynonymous

828    mutations (*i.e.* wild-type, *n*=108). **f)** Multivariate analysis of the effect of mtDNA variants on

829    overall survival time among *n*=344 TCGA colorectal cancer patients (stage 1-3). Truncating

830    variants and VUSs are each compared to wild-type samples, while controlling for known

831    prognostic clinical and genomic covariates using a Cox proportional-hazards model. Hazard

832    ratios for each covariate are shown on a log-scale, error-bars are 95% confidence intervals from

833    the Cox proportional hazards regression. Point size indicates the number of samples with the

834    associated covariate value (except for Age, which was coded as a continuous variable, and

835    therefore the size corresponds to the total number of samples). Blue points are statistically

836    significant (*Q*-value < 0.01); black points not significant; gray points are reference categories

837    and were not tested.

838

839

840

841

842                         **Supplementary Materials**

843

844    **Supplementary Tables**

845

846    **Supplementary Table 1:** Table of mtDNA position 1-16,569 annotated with gene symbols,

847    encoding strand, nucleotide, and exclusion criteria.

848    **Supplementary Table 2:** Table of mutation rates in mtDNA and nuclear cancer-associated

849    genes.

850    **Supplementary Table 3:** Table of SNV hotspot positions and associated significance and

851    annotations.

852 **Supplementary Table 4:** Table of homopolymer indel hotspots and associated significance and
853 annotations.

854 **Supplementary Table 5:** Table of tRNA structural alignment hotspots and significance and
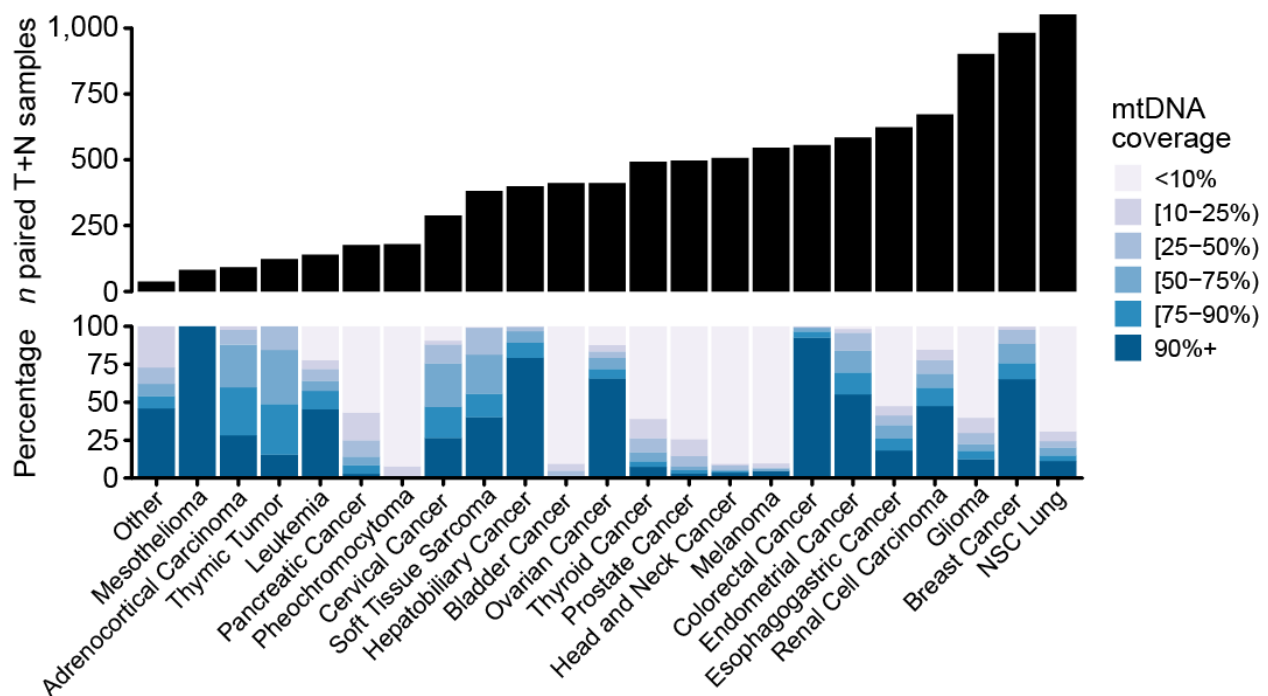855 annotations.

856 **Supplementary Table 6:** Table with mtDNA variants and mtDNA status classifications for all
857 TCGA samples included in this study.

858

859
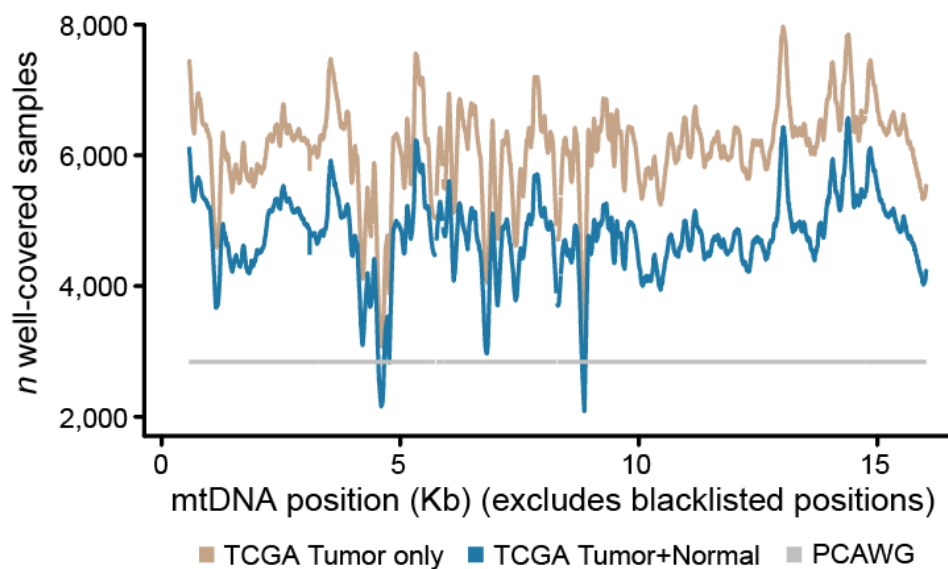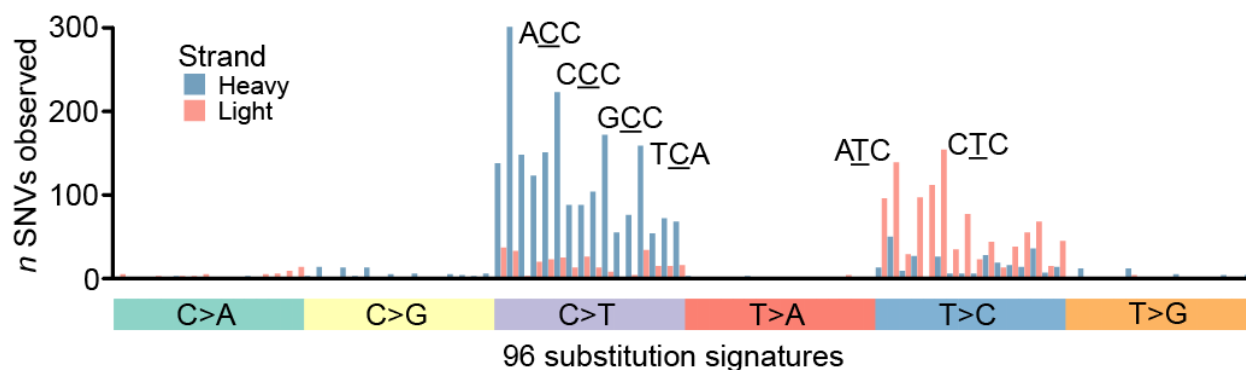
860 **Supplementary Figures**

861



862
863 **Supplementary Fig. 1a: Distribution of cancer types in patient cohort.** Top, number of
864 tumor samples and paired matched normal samples per cancer type in this study. Bottom, the
865 proportion of tumor and normal sample pairs each with ≥5 read coverage in the indicated
866 percentage of genic regions of the mitochondrial genome (*e.g.* darkest blue indicates the
867 percent of well-covered samples of the given cancer type.).

868

869

870 **Supplementary Fig. 1b: mtDNA coverage from off-target reads at each position.** The

871 number of samples for which the given mtDNA position was sequenced to adequate depth to

872 call somatic variants. Brown, the number of samples using unpaired tumor-only data, applicable

873 only for protein-truncating variants which were always assumed to be of somatic origin. Blue,

874 the number using paired tumor and matched-normal data, applicable for all non-truncating

875 variants which required evidence that the variant was absent in the matched normal to be

876 confidently classified as of somatic origin. Gray, the number of whole-genome sequenced

877 samples available from ICGC/PCAWG for comparison.
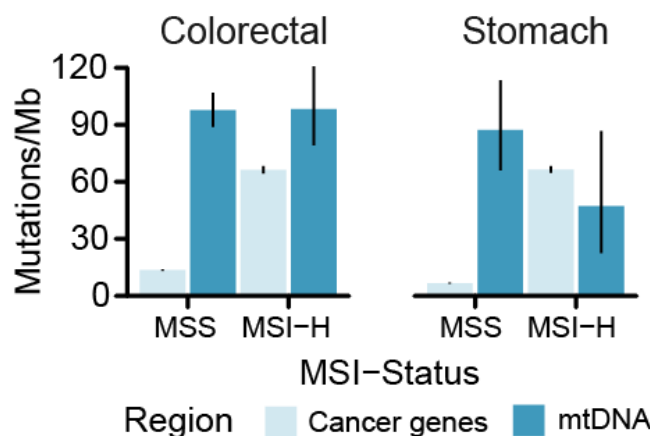
878

879

880



881

882 **Supplementary Fig. 1c: Strand-specific mutational signatures in our dataset.** The

883 frequency of somatic SNVs on the light or heavy mtDNA strand with each of the 96 possible

884 mutational signatures with trinucleotide contexts (among $n$ = 3,872 SNVs). Blue bars indicate

885 the prevalence of mutational signatures for heavy-strand encoded SNVs (substitutions at C or T
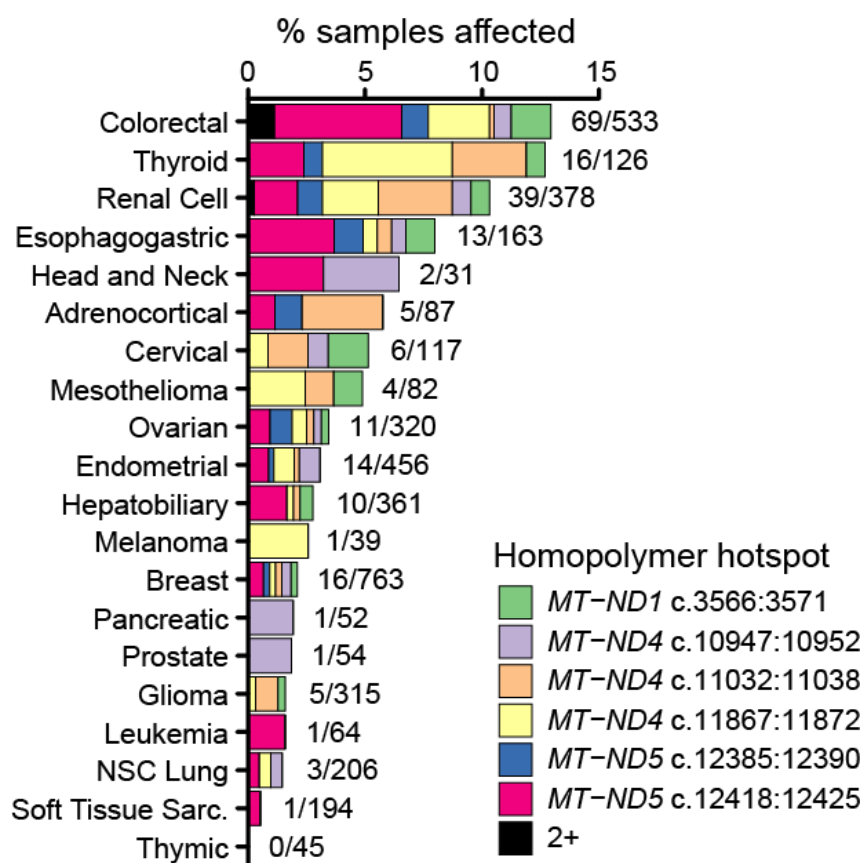
886    central nucleotides); red bars indicate those for light-strand encoded SNVs (substitutions at G or

887    T nucleotides, which were standardized to their C or T complementary nucleotide). The most

888    prevalent mutational signatures are labeled. The underlined central position is mutated with the

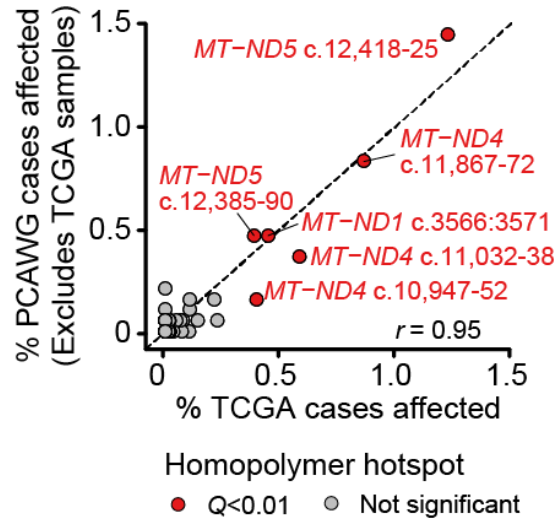889    single nucleotide substitution labeled in the tile below.

890



891

892    **Supplementary Fig. 1d: mtDNA mutation burden does not correlate with nuclear mutation**

893    **burden within cancer types.** Mitochondrial and nuclear tumor mutation burdens (TMB,

894    mutations/Mb) are shown for each well-covered tumor, among cancer types with $n \geq 100$

895    samples. Nuclear TMBs are calculated based on mutations to 468 cancer-associated genes and

896    their total coding-sequence length. Pearson correlation coefficients $r$ indicate no linear

897    correlation between mitochondrial and nuclear TMBs were observed for any cancer type tested.

898

899

900

**Supplementary Fig. 1e: Microsatellite instability does not affect somatic mtDNA mutation rate.** TMBs for somatic mtDNA mutations and mutations to cancer-associated genes are compared between microsatellite stable (MSS) and microsatellite unstable (MSI-High) tumors, for both (*n* colorectal cancer: MSI=65, MSS=318; *n* stomach adenocarcinomas: MSI=75, MSS=256). Although MSI-High tumors have elevated TMB for nuclear cancer genes, there is no effect on mtDNA TMB. Moreover, mtDNA TMB is similar to (or exceeds) that of nuclear cancer associated genes in both cancer types. Error bars are 95% Poisson exact confidence intervals.

908



909

910    **Supplementary Fig. 2a: Prevalence of frame-shift indels at homopolymer hotspots across**

911    **cancer types.** Percentage of cases per cancer type with truncating frame-shift indels at any of 6

912    indel hotspot loci. Plotted cancer types had ≥ 20 well-covered samples (*n*=4,432 paired tumor

913    and matched-normal samples total). Labels indicate the fraction of samples with any indels at

914    homopolymer hotspot out of the total number of well-covered samples for the given cancer type.

915

916



917

918    **Supplementary Fig. 2b: Validation of homopolymeric indel hotspot loci.** The proportion of

919    samples in TCGA (X-axis) or PCAWG (excluding samples also in TCGA, Y-axis) with frame-

920    shift indels at 73 homopolymeric regions. The 6 indel hotspot loci are colored red and labeled.

921    y=x is shown as a dashed line. Pearson correlation coefficient *r* as indicated.
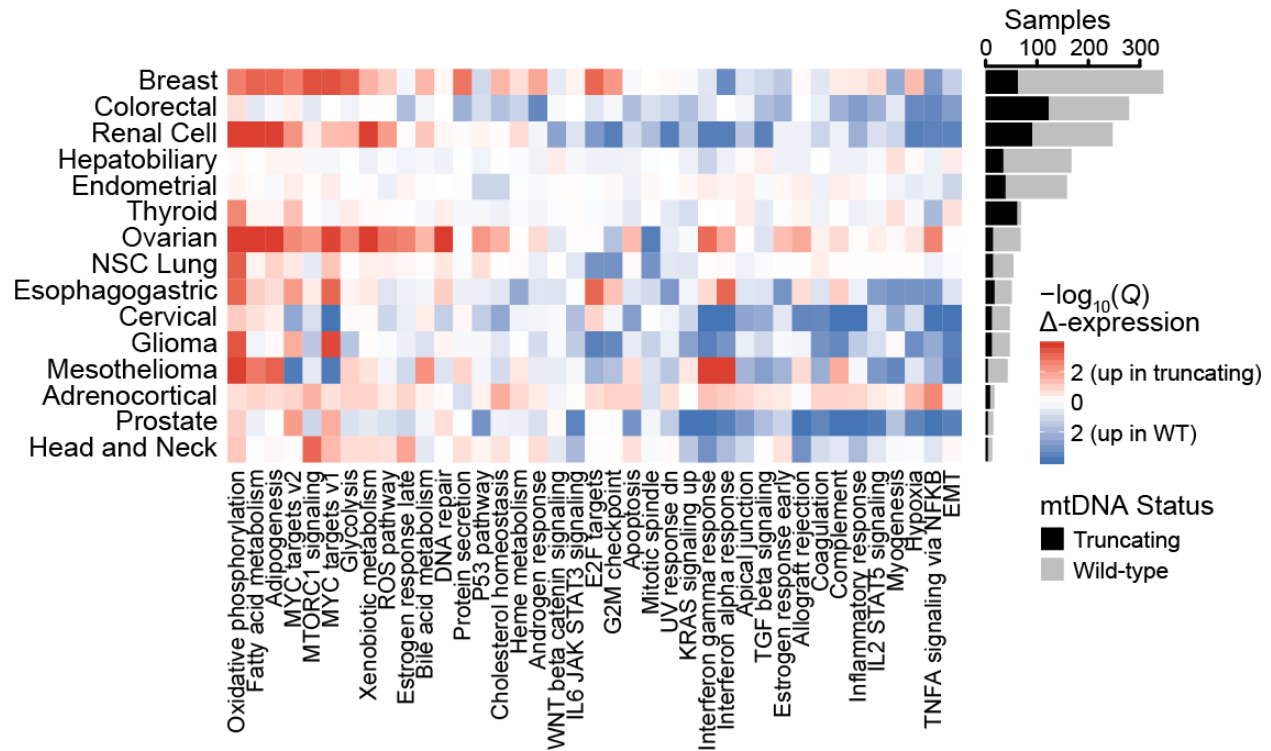
922

**Supplementary Fig. 3c: Validation of tRNA structural hotspots in PCAWG.** The number of samples with SNVs in tRNAs at the indicated cloverleaf structural position, bottom; top, the statistical enriched of the given position for mutations. Position 31 $Q$-value=0.014, $n$=196 tRNA mutations among 1,951 PCAWG samples.



**Supplementary Fig. 4a: Proportion of samples with detectable mutations is not biased by cancer type sequencing coverage.** There is no correlation between the fraction of well-covered samples in a cancer type and the proportion of well-covered samples with a detectable somatic mtDNA mutation. Cancer types with ≥30 well-covered samples shown, $P$ value from linear regression.

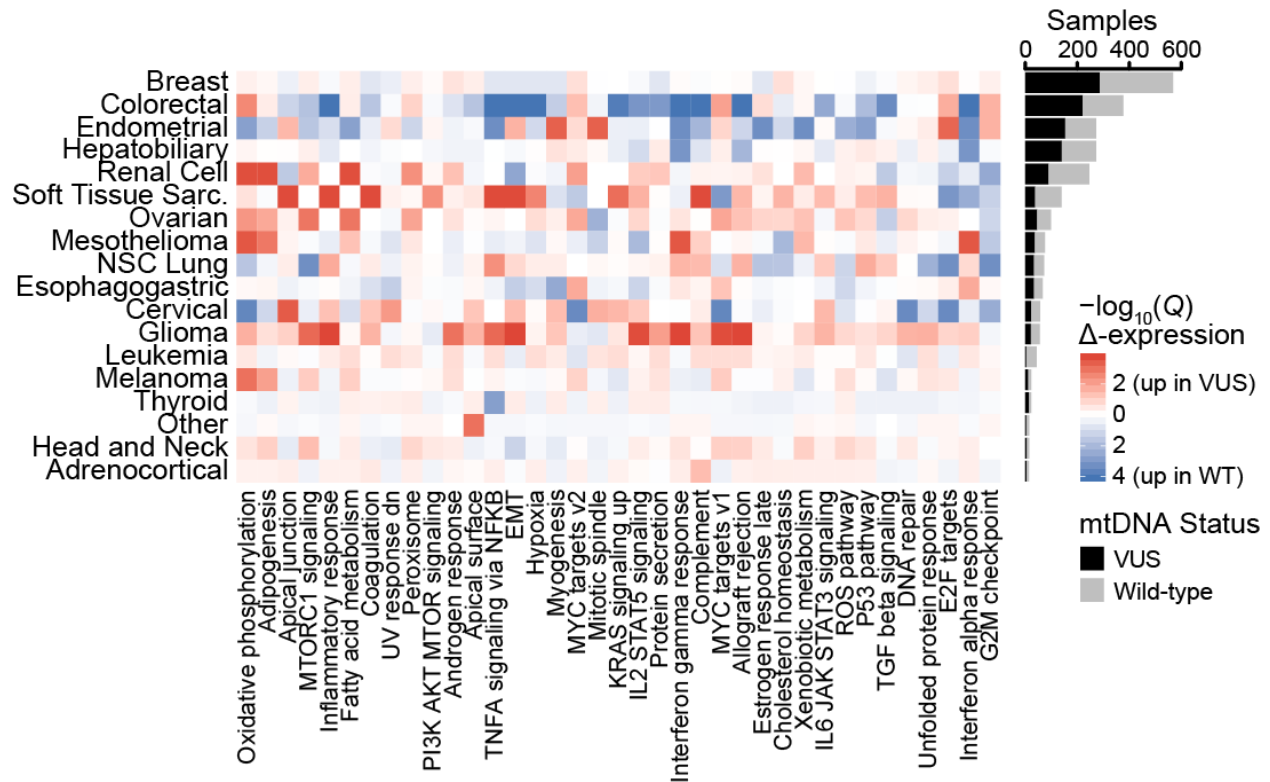**Supplementary Fig. 4b: Transcriptional dysregulation attributed to truncating mtDNA variants.** (Left) Heatmap shows directional significance of dysregulation of a given geneset in tumors with truncating variants among the given cancer type; $-\log_{10}(Q$-value$) > 2$ indicates significant up-regulation, $< -2$ indicates significant down-regulation. (Right) Histogram of wild-type samples and samples with truncating variants used to calculate differentially-expressed genes and dysregulated genesets.
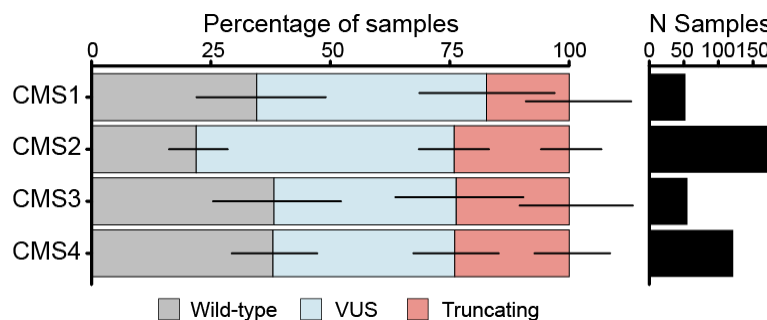
**Supplementary Fig. 4c: Transcriptional dysregulation attributed to mtDNA VUSs.**
Heatmap, differentially expressed mSigDB Hallmarks genesets between tumors with any somatic VUSs or wild-type mtDNA. Genesets ordered from most up-regulated across cancer types to most down-regulated. Barplot, number of cases with VUSs or wild-type mtDNA.



**Supplementary Fig. 4d: Difference in mtDNA mutation status between colorectal cancer consensus molecular subtypes.** Left, the proportion of samples with wild-type mtDNA (*i.e.* no somatic mutations), VUS (any non-truncating) or truncating variants among colorectal tumors with each consensus molecular subtype (CMS) is shown. Right, histogram of the number of well-covered colorectal tumors. There was a statistically significant difference in mtDNA

956　mutation status between different CMS classifications ($P$=0.03, Chi-squared test, $n$=415

957　samples total).

958

959

## Bibliography

961　1.　Rheinbay, E. *et al.* Analyses of non-coding somatic drivers in 2,658 cancer whole genomes.

962　　*Nature* **578,** 102–111 (2020).

963　2.　Ju, Y. S. *et al.* Origins and functional consequences of somatic mitochondrial DNA

964　　mutations in human cancer. *elife* **3,** (2014).

965　3.　Yuan, Y. *et al.* Comprehensive molecular characterization of mitochondrial genomes in

966　　human cancers. *Nat. Genet.* **52,** 342–352 (2020).

967　4.　Stewart, J. B. *et al.* Simultaneous DNA and RNA Mapping of Somatic Mitochondrial

968　　Mutations across Diverse Human Cancers. *PLoS Genet.* **11,** e1005333 (2015).

969　5.　Grandhi, S. *et al.* Heteroplasmic shifts in tumor mitochondrial genomes reveal tissue-

970　　specific signals of relaxed and positive selection. *Hum. Mol. Genet.* **26,** 2912–2922 (2017).

971　6.　Hopkins, J. F. *et al.* Mitochondrial mutations drive prostate cancer aggression. *Nat.*

972　　*Commun.* **8,** 656 (2017).

973　7.　To, T.-L. *et al.* A Compendium of Genetic Modifiers of Mitochondrial Dysfunction Reveals

974　　Intra-organelle Buffering. *Cell* **179,** 1222-1238.e17 (2019).

975　8.　Birsoy, K. *et al.* An essential role of the mitochondrial electron transport chain in cell

976　　proliferation is to enable aspartate synthesis. *Cell* **162,** 540–551 (2015).

977　9.　Samuels, D. C. *et al.* Finding the lost treasures in exome sequencing data. *Trends Genet.*

978　　**29,** 593–599 (2013).

979　10.　Cibulskis, K. *et al.* Sensitive detection of somatic point mutations in impure and

980　　heterogeneous cancer samples. *Nat. Biotechnol.* **31,** 213–219 (2013).

981　11.　Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25,** 2078–

982      2079 (2009).

983   12. Collura, R. V., Auerbach, M. R. & Stewart, C. B. A quick, direct method that can

984      differentiate expressed mitochondrial genes from their nuclear pseudogenes. *Curr. Biol.* **6,**

985      1337–1339 (1996).

986   13. Cheng, D. T. *et al.* Memorial Sloan Kettering-Integrated Mutation Profiling of Actionable

987      Cancer Targets (MSK-IMPACT): A Hybridization Capture-Based Next-Generation

988      Sequencing Clinical Assay for Solid Tumor Molecular Oncology. *J. Mol. Diagn.* **17,** 251–264

989      (2015).

990   14. Bolze, A. *et al.* Selective constraints and pathogenicity of mitochondrial DNA variants

991      inferred from a novel database of 196,554 unrelated individuals. *BioRxiv* (2019).

992      doi:10.1101/798264

993   15. Chang, M. T. *et al.* Identifying recurrent mutations in cancer reveals widespread lineage

994      diversity and mutational specificity. *Nat. Biotechnol.* **34,** 155–163 (2016).

995   16. Chang, M. T. *et al.* Accelerating discovery of functional mutant alleles in cancer. *Cancer*

996      *Discov.* **8,** 174–183 (2018).

997   17. Gopal, R. K. *et al.* Early loss of mitochondrial complex I and rewiring of glutathione

998      metabolism in renal oncocytoma. *Proc Natl Acad Sci USA* **115,** E6283–E6290 (2018).

999   18. Triska, P. *et al.* Landscape of germline and somatic mitochondrial DNA mutations in

1000      pediatric malignancies. *Cancer Res.* **79,** 1318–1330 (2019).

1001   19. Alston, C. L. *et al.* A novel mitochondrial MTND5 frameshift mutation causing isolated

1002      complex I deficiency, renal failure and myopathy. *Neuromuscul. Disord.* **20,** 131–135

1003      (2010).

1004   20. Castellana, S. *et al.* High-confidence assessment of functional impact of human

1005      mitochondrial non-synonymous genome variations by APOGEE. *PLoS Comput. Biol.* **13,**

1006      e1005628 (2017).

1007   21. Martínez-Reyes, I. *et al.* Mitochondrial ubiquinol oxidation is necessary for tumour growth.

1008    *Nature* (2020). doi:10.1038/s41586-020-2475-6

1009    22. El-Hattab, A. W., Adesina, A. M., Jones, J. & Scaglia, F. MELAS syndrome: Clinical

1010    manifestations, pathogenesis, and treatment options. *Mol. Genet. Metab.* **116,** 4–12 (2015).

1011    23. Gorman, G. S. *et al.* Prevalence of nuclear and mitochondrial DNA mutations related to

1012    adult mitochondrial disease. *Ann. Neurol.* **77,** 753–759 (2015).

1013    24. Gopal, R. K. *et al.* Widespread chromosomal losses and mitochondrial DNA alterations as

1014    genetic drivers in hürthle cell carcinoma. *Cancer Cell* **34,** 242-255.e5 (2018).

1015    25. Spagnolo, M. *et al.* A new mutation in the mitochondrial tRNA(Ala) gene in a patient with

1016    ophthalmoplegia and dysphagia. *Neuromuscul. Disord.* **11,** 481–484 (2001).

1017    26. Miller, M. L. *et al.* Pan-Cancer Analysis of Mutation Hotspots in Protein Domains. *Cell Syst.*

1018    **1,** 197–209 (2015).

1019    27. Gao, J. *et al.* 3D clusters of somatic mutations in cancer reveal numerous rare mutations as

1020    functional targets. *Genome Med.* **9,** 4 (2017).

1021    28. Horváth, R., Reilmann, R., Holinski-Feder, E., Ringelstein, E. B. & Klopstock, T. The role of

1022    complex I genes in MELAS: a novel heteroplasmic mutation 3380G>A in ND1 of mtDNA.

1023    *Neuromuscul. Disord.* **18,** 553–556 (2008).

1024    29. Agip, A.-N. A. *et al.* Cryo-EM structures of complex I from mouse heart mitochondria in two

1025    biochemically defined states. *Nat. Struct. Mol. Biol.* **25,** 548–556 (2018).

1026    30. Joshi, S. *et al.* The genomic landscape of renal oncocytoma identifies a metabolic barrier to

1027    tumorigenesis. *Cell Rep.* **13,** 1895–1908 (2015).

1028    31. Ganly, I. *et al.* Integrated genomic analysis of hürthle cell cancer reveals oncogenic drivers,

1029    recurrent mitochondrial mutations, and unique chromosomal landscapes. *Cancer Cell* **34,**

1030    256-270.e5 (2018).

1031    32. Guinney, J. *et al.* The consensus molecular subtypes of colorectal cancer. *Nat. Med.* **21,**

1032    1350–1356 (2015).

1033    33. Yaeger, R. *et al.* Clinical sequencing defines the genomic landscape of metastatic

1034      colorectal cancer. *Cancer Cell* **33,** 125-136.e3 (2018).

1035    34. Priolo, C. *et al.* Impairment of gamma-glutamyl transferase 1 activity in the metabolic

1036      pathogenesis of chromophobe renal cell carcinoma. *Proc Natl Acad Sci USA* **115,** E6274–

1037      E6282 (2018).

1038    35. Cancer Genome Atlas Research Network *et al.* The Cancer Genome Atlas Pan-Cancer

1039      analysis project. *Nat. Genet.* **45,** 1113–1120 (2013).

1040    36. McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing

1041      next-generation DNA sequencing data. *Genome Res.* **20,** 1297–1303 (2010).

1042    37. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping and

1043      population genetical parameter estimation from sequencing data. *Bioinformatics* **27,** 2987–

1044      2993 (2011).

1045    38. Chakravarty, D. *et al.* Oncokb: A precision oncology knowledge base. *JCO Precis. Oncol.*

1046      **2017,** (2017).

1047    39. Sonney, S. *et al.* Predicting the pathogenicity of novel variants in mitochondrial tRNA with

1048      MitoTIP. *PLoS Comput. Biol.* **13,** e1005867 (2017).

1049    40. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and  dispersion

1050      for RNA-seq data with DESeq2. *Genome Biol.* **15,** 550 (2014).

1051    41. Zhu, A., Ibrahim, J. G. & Love, M. I. Heavy-tailed prior distributions for sequence count

1052      data: removing the noise and preserving large differences. *Bioinformatics* **35,** 2084–2092

1053      (2019).

1054    42. Benjamini,  y Y. *et al.* Controlling the False Discovery Rate : A Practical and Powerful

1055      Approach to Multiple Testing Author ( s ): Yoav Benjamini and Yosef Hochberg Source :

1056      Journal of the Royal Statistical Society . Series B ( Methodological ), Vol . 57 , No . 1

1057      Published by : *J R Statist Soc B* **57,** 289–300 (1995).

1058    43. Liberzon, A. *et al.* The Molecular Signatures Database (MSigDB) hallmark gene set

1059      collection. *Cell Syst.* **1,** 417–425 (2015).

1060   44. Sergushichev, A. An algorithm for fast preranked gene set enrichment analysis using

1061      cumulative statistic calculation. *BioRxiv* (2016). doi:10.1101/060012

1062   45. Liu, Y. *et al.* Comparative molecular analysis of gastrointestinal adenocarcinomas. *Cancer*

1063      *Cell* **33,** 721-735.e8 (2018).

1064   46. Ellrott, K. *et al.* Scalable open science approach for mutation calling of tumor exomes using

1065      multiple genomic pipelines. *Cell Syst.* **6,** 271-281.e7 (2018).

1066   47. Pettersen, E. F. *et al.* UCSF Chimera—a visualization system for exploratory research and

1067      analysis. *J. Comput. Chem.* **25,** 1605–1612 (2004).

1068   48. Jurcik, A. *et al.* CAVER Analyst 2.0: analysis and visualization of channels and tunnels in

1069      protein structures and molecular dynamics trajectories. *Bioinformatics* **34,** 3586–3588

1070      (2018).

1071   49. Zhu, J., Vinothkumar, K. R. & Hirst, J. Structure of mammalian respiratory complex I.

1072      *Nature* **536,** 354–358 (2016).

1073   50. Baker, N. A., Sept, D., Joseph, S., Holst, M. J. & McCammon, J. A. Electrostatics of

1074      nanosystems: application to microtubules and the ribosome. *Proc Natl Acad Sci USA* **98,**

1075      10037–10041 (2001).

1076   51. Dolinsky, T. J., Nielsen, J. E., McCammon, J. A. & Baker, N. A. PDB2PQR: an automated

1077      pipeline for the setup of Poisson-Boltzmann electrostatics calculations. *Nucleic Acids Res.*

1078      **32,** W665-7 (2004).