


PRIMARY RESEARCH

Open Access



Mitochondrial genome copy number measured by DNA sequencing in human blood is strongly associated with metabolic traits via cell-type composition differences

Liron Ganel^{1,2} , Lei Chen^{1,2}, Ryan Christ¹, Jagadish Vangipurapu³, Erica Young^{1,4}, Indrani Das¹, Krishna Kanchi¹, David Larson^{1,5}, Allison Regier^{1,2}, Haley Abel^{1,2,5}, Chul Joo Kang¹, Alexandra Scott^{1,2}, Aki Havulinna^{6,7}, Charleston W. K. Chiang^{8,9}, Susan Service¹⁰, Nelson Freimer¹⁰, Aarno Palotie^{6,11,12}, Samuli Ripatti^{6,12,13}, Johanna Kuusisto^{3,14}, Michael Boehnke¹⁵, Markku Laakso^{3,14}, Adam Locke^{1,2}, Nathan O. Stitzel^{1,4,5*} and Ira M. Hall^{1,2,16*}

Abstract

Background: Mitochondrial genome copy number (MT-CN) varies among humans and across tissues and is highly heritable, but its causes and consequences are not well understood. When measured by bulk DNA sequencing in blood, MT-CN may reflect a combination of the number of mitochondria per cell and cell-type composition. Here, we studied MT-CN variation in blood-derived DNA from 19184 Finnish individuals using a combination of genome (N = 4163) and exome sequencing (N = 19034) data as well as imputed genotypes (N = 17718).

Results: We identified two loci significantly associated with MT-CN variation: a common variant at the *MYB-HBS1L* locus ($P = 1.6 \times 10^{-8}$), which has previously been associated with numerous hematological parameters; and a burden of rare variants in the *TMBIM1* gene ($P = 3.0 \times 10^{-8}$), which has been reported to protect against non-alcoholic fatty liver disease. We also found that MT-CN is strongly associated with insulin levels ($P = 2.0 \times 10^{-21}$) and other metabolic syndrome (metS)-related traits. Using a Mendelian randomization framework, we show evidence that MT-CN measured in blood is causally related to insulin levels. We then applied an MT-CN polygenic risk score (PRS) derived from Finnish data to the UK Biobank, where the association between the PRS and metS traits was replicated. Adjusting for cell counts largely eliminated these signals, suggesting that MT-CN affects metS via cell-type composition.

Conclusion: These results suggest that measurements of MT-CN in blood-derived DNA partially reflect differences in cell-type composition and that these differences are causally linked to insulin and related traits.

Keywords: Metabolic syndrome, Mitochondrial content, Human genetics, Human genome sequencing, Genome-wide association studies, Mendelian randomization

* Correspondence: nstitzel@wustl.edu; ira.hall@yale.edu

¹McDonnell Genome Institute, Washington University School of Medicine, St. Louis, MO, USA

Full list of author information is available at the end of the article



© The Author(s). 2021 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

Background

Cardiovascular disease (CVD) is a category comprising numerous diseases of the circulatory system, including coronary heart disease, heart failure, stroke, and hypertension [1]. Collectively, these diseases are the leading cause of mortality both globally and in the USA [1, 2]. Metabolic syndrome (metS), a class of disorders related to CVD with high prevalence in the USA, includes dyslipidemia, obesity, insulin resistance, and prothrombotic and proinflammatory states [3]. Individuals with metS have approximately twofold risk of being diagnosed with CVD over 5 to 10 years and fivefold risk of being diagnosed with type 2 diabetes mellitus [4].

There are many reported links between mitochondrial content and metS-related phenotypes in various tissues, including adipose [5–7], liver [5, 8, 9], skeletal muscle [5, 10–14], and blood [15–21]. Traits associated with mitochondrial (MT) content include CHD, type 2 diabetes, and metabolic syndrome traits such as insulin sensitivity/resistance, obesity, and blood triglycerides. However, these studies have generally been limited by small sample sizes and low statistical power. This, in addition to the use of heterogeneous mitochondrial quantification methods [22], has led to inconsistencies in the literature about the strength and directions of effect between mitochondrial content and metS traits. In one large WGS study of mitochondrial genome copy number (MT-CN) in 2077 Sardinians, Ding et al. estimated the heritability of MT-CN at 54% and detected significant associations between MT-CN and both waist circumference and waist–hip ratio, but found no association with body mass index (BMI) [15]. Another large study (N = 5150) found virtually no evidence of association between qPCR-measured MT-CN and any of several cardiometabolic phenotypes [23]. The only exception was an inverse association with insulin that was identified in one cohort but did not survive meta-analysis across cohorts. However, a study of 21870 individuals from 3 cohorts showed a significant inverse relationship between MT-CN (measured by microarray probe intensities in two cohorts and qPCR in the third) and incident cardiovascular disease [24].

Although variations in MT-CN measured from whole blood can in principle be attributed to either variability of MT copy number within cells or the cell-type composition of the blood (given that different cell types have varying MT content [25–27]), the literature on this subject is inconclusive. Using CpG methylation data, a large (N = 11443), low-coverage (1.7x autosomal; 102x mitochondrial) sequencing study of the link between MT-CN and major depressive disorder using buccal DNA from Chinese women concluded that variability of MT-CN from buccal swabs was not due to differences in cell-type composition [28]. However, this study did not do a

similar experiment in blood. Two small (N = 756 and N = 400) studies identified an association between MT content and CHD that they attributed to variable MT-CN within leukocytes, but they did not directly investigate the possibility of cell-type composition being the true driver of the association [16, 21]. For brevity, we will use the term “MT-CN” to refer to the underlying phenotype reflected by measuring this quantity for the remainder of this work, with these caveats.

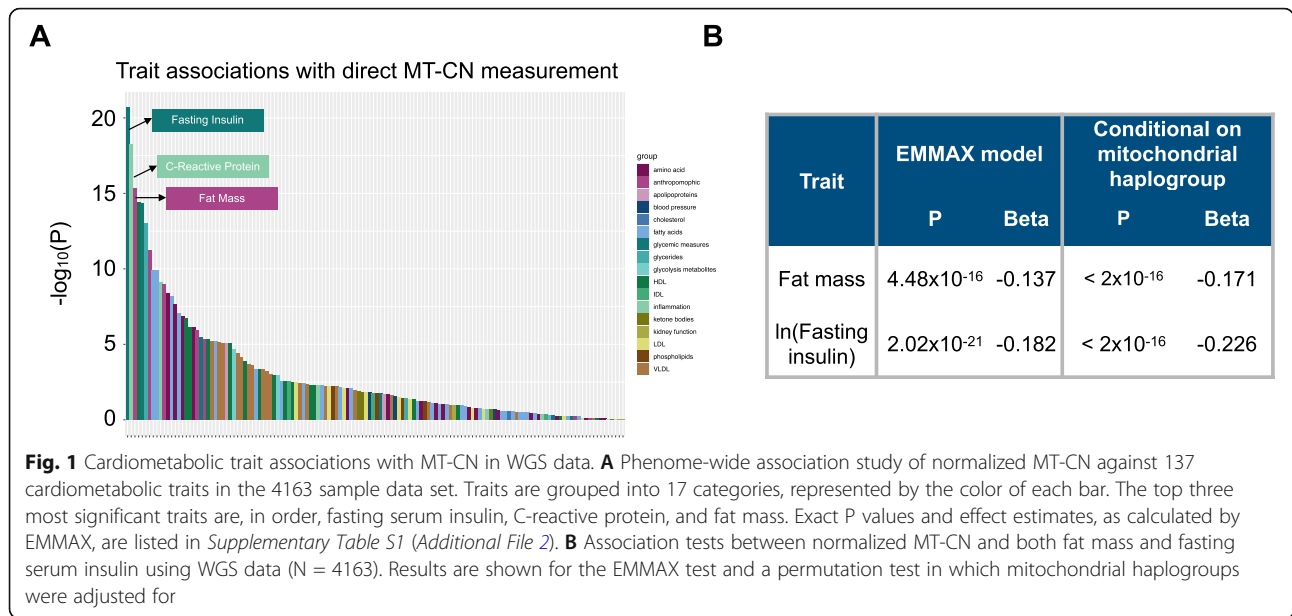
While several studies have found that peripheral blood MT content is heritable, only a small number of MT-CN associated loci have been identified [29–31]. In one of these studies, Curran et al. used linkage analysis in Mexican Americans to find an MT-CN associated locus near a marker previously associated with triglyceride levels [30, 32, 33], providing further indirect evidence for the link between MT-CN and metabolic syndrome.

Here, we take advantage of large-scale genome, exome, and array genotype data to investigate the causes and effects of MT-CN in a large, deeply phenotyped Finnish cohort. Our results reveal novel links with metabolic syndrome and provide evidence supporting a causal role for MT-CN.

Results

Association of MT-CN with metabolic traits

We estimated MT-CN in 4163 individuals from the METSIM and FINRISK studies based on deep (> 20x coverage) WGS data. We did so by measuring the mean coverage depth of reads mapped to the mitochondrial genome in each sample, and normalizing it to the mean autosomal coverage (see Methods). We performed batch normalization separately for METSIM and for two FINRISK batches separated by survey years (see Methods). Each measurement was adjusted for age, age², and sex, then inverse rank normalized separately before combining across batches. We tested the resulting MT-CN estimates for association with 137 quantitative traits that were collected and normalized according to the procedures described previously [34]. MT-CN was strongly associated with fat mass (P = 4.48 × 10⁻¹⁶) and fasting serum insulin (P = 2.02 × 10⁻²¹), as well as numerous additional quantitative traits, many related to metabolic syndrome (Fig. 1a, Table S1 - Additional File 2). Notably, BMI was significantly associated with MT-CN, despite the fact that Ding et al. did not find evidence of this association [15]. Since population structure was a potential confounder in this analysis considering the presence of mtDNA polymorphisms that might adversely affect short-read alignment, we included SNP-inferred mitochondrial haplogroup as a covariate and reran the tests (Fig. 1b). The association signals retained significance even after this adjustment.



To understand the connection between MT-CN and more clinically relevant phenotypes, we tested our MT-CN estimate against Matsuda ISI and disposition index (*Table 1*), which measure insulin sensitivity and secretion, respectively, and were not included in the initial screen. MT-CN was strongly associated with both insulin phenotypes. Notably, the Matsuda ISI signals survived adjustment for fat mass percentage after excluding diabetic individuals, which indicates that the association of peripheral blood MT-CN with insulin sensitivity was independent of fat mass.

To test for this association signal in a larger cohort, we developed a method to estimate mitochondrial genome copy number using 19034 samples with whole exome sequencing (WES) data from the METSIM and FINRISK studies that included most of the WGS samples [34] (see Methods). R² between WGS-based and WES-based estimates was 0.445 (*Figure S6 - Additional File 1*). Consistent with the WGS-based analysis, WES-estimated MT-CN was significantly associated with both fat mass and fasting serum insulin levels, even after removing the samples with WGS data, with identical directions of effect (*Table S2 - Additional File 1*).

Table 1 Associations of normalized MT-CN with disposition index and Matsuda ISI in METSIM. Testing was done by linear regression using disposition index and Matsuda ISI, respectively, as the dependent variable. P* columns represent the P value from linear regression with additional adjustment for fat mass. Follow-up measurements were taken at a later time point.

	Baseline					Follow-up				
	Number	Beta	SE	P	P*	Number	Beta	SE	P	P*
Disposition index										
All subjects	2975	0.094	0.004	3.0 x 10 ⁻⁷	0.0004	2492	0.062	0.004	0.002	0.068
Excludes diabetic subjects at baseline	2842	0.091	0.003	1.3 x 10 ⁻⁶	0.0007	2452	0.067	0.004	0.0009	0.041
Excludes diabetic subjects at baseline and during follow-up	2453	0.069	0.003	0.0007	0.023	2449	0.067	0.004	0.0009	0.042
Matsuda ISI										
All subjects	2975	0.192	0.005	4.3 x 10 ⁻²⁶	7.3 x 10 ⁻¹⁷	2492	0.157	0.006	3.7 x 10 ⁻¹⁵	8.7 x 10 ⁻¹⁰
Excludes diabetic subjects at baseline	2842	0.191	0.005	1.0 x 10 ⁻²⁴	2.4 x 10 ⁻¹⁶	2452	0.161	0.006	1.3 x 10 ⁻¹⁵	3.0 x 10 ⁻¹⁰
Excludes diabetic subjects at baseline and during follow-up	2453	0.173	0.005	7.2 x 10 ⁻¹⁸	5.3 x 10 ⁻¹²	2449	0.16	0.006	1.7 x 10 ⁻¹⁵	3.8 x 10 ⁻¹⁰

Heritability analysis

To assess the extent to which MT-CN is genetically determined, we estimated the heritability of mitochondrial genome copy number using GREML (Table 2). We explored two different approaches available: (1) analysis of the 4149 samples with WGS data that passed quality control measures, where both nuclear genotypes and MT-CN are measured directly from the WGS data, and (2) analysis of the set of 17718 samples with imputed genotype array data, where MT-CN is estimated from WES data. Of these, (1) benefited from more accurate measurement of genotype and phenotype, whereas (2) had noisier measurements but benefited from larger sample size. We focused primarily on the METSIM cohort, both because of the homogeneity of this cohort (see Methods) and because the number of FINRISK samples with WGS data was small.

In the WGS analysis, the GREML-estimated heritability of MT-CN in METSIM was 31%, somewhat less than the 54% value reported in the only prior large-scale study of peripheral blood MT-CN heritability, which was based on low-coverage WGS [15]. For comparison, we used this same approach to estimate heritability of LDL in METSIM WGS data, which yielded an estimate of 34% with a standard error of 7.9% (Table 3). This is broadly consistent with prior work [35, 36], including analysis of the same Finnish sample set using distinct methods [34] (20.2% heritability). These results show that mitochondrial genome copy number is a genetically determined trait with significant heritability, comparable to that of LDL and other quantitative cardiometabolic traits [34].

The analysis of imputed METSIM genotypes using WES-estimated MT-CN yielded an estimated heritability of 11%, which is much lower than the WGS-based estimate (Table 2). To understand this discrepancy, we repeated the GREML analysis with the other two combinations of phenotype source (WGS vs. WES estimation) and genotype source (WGS vs. imputed array). When using the WGS-measured phenotype, the estimated heritability decreased only slightly (31% to 27%)

when switching from the WGS to imputed genotypes. This suggests that the difference in genotyping method was not the main driver of the observed heritability disparity between the WGS and imputed array datasets. Conversely, when analyzing the imputed METSIM genotypes, switching from WGS-measured to WES-measured MT-CN resulted in a large drop (27% to 11%) in estimated heritability. This suggests that the extra noise inherent in WES-based MT-CN estimates was responsible for the reduction in the GREML-estimated heritability despite the increased sample size of the imputed array dataset.

Identification of genetic factors associated with MT-CN

Previous studies have identified three autosomal quantitative trait loci (QTL) reaching genome-wide significance for MT-CN in other populations [29, 30]. Another recent study identified two putative QTLs with suggestive P values [31]. We conducted single-variant GWAS for MT-CN (see Methods). Analysis of WGS (N = 4149) and WES (N = 19034) genotypes yielded no variants exceeding the respective significance thresholds of 5×10^{-8} and 5×10^{-7} (Figure S2 - Additional File 1). However, despite the increased noise in the WES-measured phenotype, GWAS of imputed array genotypes from METSIM (N = 9791) yielded two loci with genome-wide significant associations, identified by lead markers rs2288464 and rs9389268 (Fig. 2, Table 4). Of the previously reported MT-CN QTLs [29–31], we observed an inconclusive signal at rs445 (P = 0.048) and a significant signal at rs709591 (P = 1.61×10^{-4}), a locus associated with neutrophil count [37, 38] (Table S11 - Additional File 1). No significant signal was observed at the other two single-variant QTLs (Table S11 - Additional File 1) or the linkage peak identified by Curran et al. (Figure S3 - Additional File 1).

rs9389268 was the only marker that was strongly associated with MT-CN in the METSIM analyses of both WGS and imputed array data (P = 3.24×10^{-8} and P = 1.26×10^{-10} , respectively). Although this variant was not significantly associated with MT-CN in FINRISK (P =

Table 2. GREML heritability estimates in each cohort separately and in joint analysis. All analyses in this table were limited to sample sets with available imputed genotype data, yielding slightly lower sample sizes than in other tables.

		WGS-measured MT-CN			WES-measured MT-CN		
		Number	h ²	SE	Number	h ²	SE
Joint analysis	WGS genotypes	4149	0.17	0.06	3916	0.11	0.06
	Imputed genotypes	3916	0.16	0.06	17718	0.09	0.02
METSIM	WGS genotypes	3065	0.31	0.07	2974	0.20	0.08
	Imputed genotypes	2974	0.27	0.08	9791	0.11	0.03
FINRISK	WGS genotypes	1084	0.20	0.22	942	0.24	0.27
	Imputed genotypes	942	0.35	0.27	7927	0.08	0.03

Table 3 GREML and GREML-LDMS heritability estimates for normalized MT-CN and low-density lipoprotein (LDL) in METS IM. GREML-LDMS heritability estimates are calculated using PCs 1–10 as fixed-effect covariates. Analyses of imputed array data exclude samples with WGS data.

Trait	Genotype source (phenotype source)	Number	GREML		GREML-LDMS	
			h ²	SE	h ²	SE
Normalized MT-CN	WGS (WGS-measured MT-CN)	3065	0.31	0.07	0.31	0.09
	Imputed array (WES-measured MT-CN)	6789	0.11	0.04	0.14	0.05
LDL	WGS	3062	0.34	0.08	0.38	0.10
	Imputed array	6787	0.25	0.04	0.32	0.05

0.788 and $P = 0.189$ in WGS and imputed array data, respectively) or in a separate random-effects meta-analysis of both cohorts ($P = 0.115$), the lack of signal in FINRISK is likely the product of lower-quality MT-CN measurements in FINRISK, which displayed heterogeneity

across survey years (Figure S4 - Additional File 1). This variant is located in an intergenic region between the *MYB* and *HBS1L* genes, is common across many populations, and is slightly more frequent in Finns compared with non-Finnish Europeans (gnomAD v3 MAF 34.4%

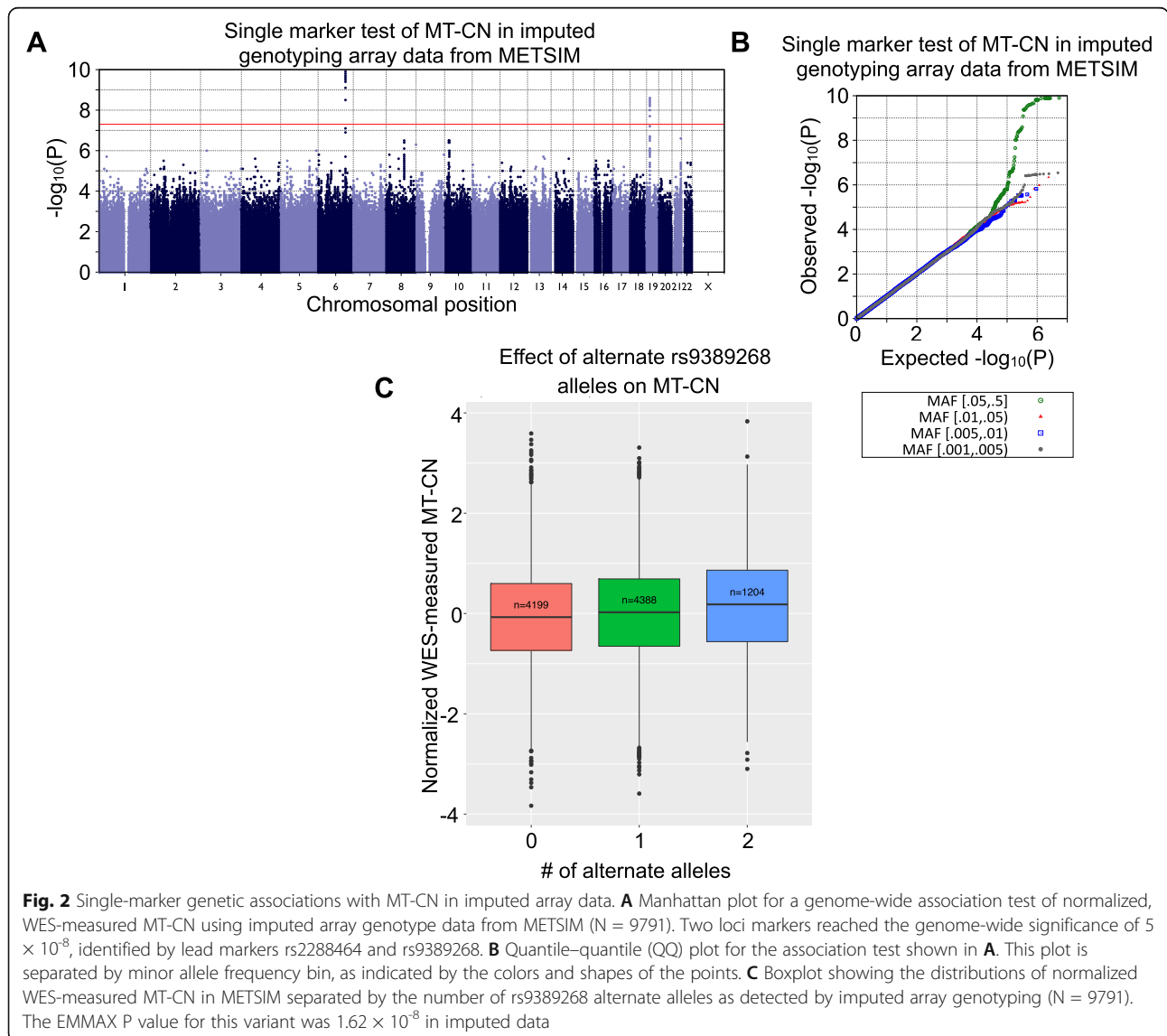
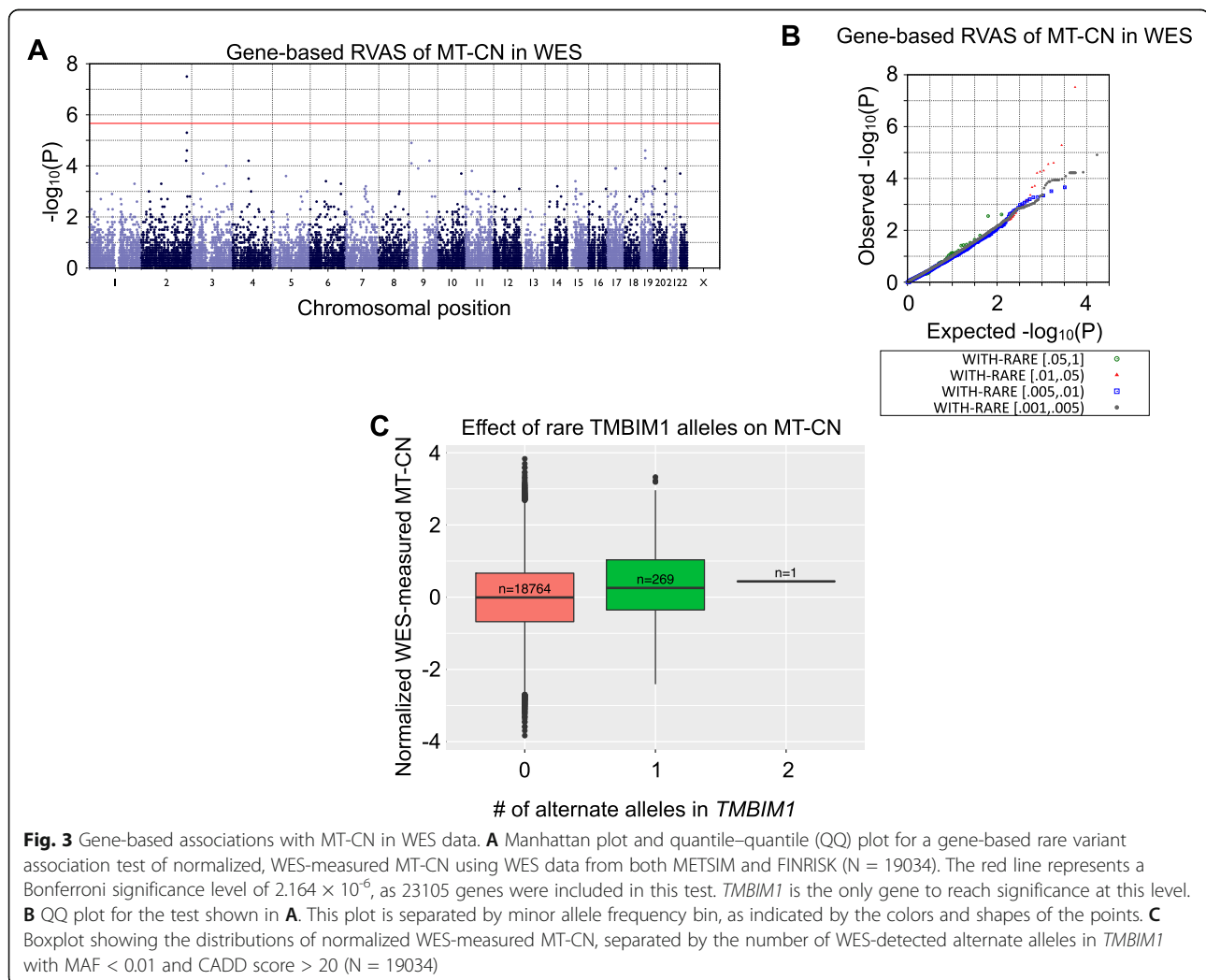


Table 4 Single marker association results for rs2288464 and rs9389268. Analyses of imputed FINRISK array data were performed with covariates for FINRISK genotyping batch. Bolded results are significant at the appropriate threshold for the given test (see Methods).

		FINRISK				METSIM				Joint analysis			
		Number	MAF	P	Beta	Number	MAF	P	Beta	Number	MAF	P	Beta
rs2288464	Imputed	7927	0.148	0.613	0.0113	9791	0.165	2.55×10^{-9}	0.119	17718	0.158	9.77×10^{-7}	0.075
	WES	9221	0.150	0.376	0.0186	9813	0.166	6.75×10^{-9}	0.118	19034	0.158	9.34×10^{-7}	0.0734
	WGS	1084	0.142	0.383	0.0532	3065	0.161	0.113	0.0561	4149	0.156	0.0655	0.0562
rs9389268	Imputed	7927	0.354	0.189	0.0216	9791	0.347	1.26×10^{-10}	0.0973	17718	0.35	1.62×10^{-8}	0.0634
	WGS	1084	0.351	0.788	0.0121	3065	0.347	3.24×10^{-8}	0.150	4149	0.348	7.87×10^{-7}	0.115

vs. 26.0%). *MYB* and *HBS1L* are hematopoietic regulators [39, 40], and the region between them is known to be associated with many hematological parameters including fetal hemoglobin levels, hematocrit, and erythrocyte, platelet, and monocyte counts [41–44]. It has been suggested that these intergenic variants function by disrupting *MYB* transcription factor binding and disrupting enhancer–promoter looping [45]. Conditioning the

METSIM-only imputed array GWAS on rs9399137—a tag SNP shown to be associated with many of these hematological parameters [43]—resulted in elimination of the rs9389268 signal entirely ($P = 0.408$), suggesting that the haplotype responsible for the association of rs9389268 with MT-CN in our data is the same one previously known to be associated with numerous hematological phenotypes.



rs2288464 seemed to be a good candidate due to its location in the 3' untranslated region of *MRPL34*, which codes for a large subunit protein of the mitochondrial ribosome. While the association signal at this marker was not observed in the WGS data ($P = 0.0655$), based on the observed effect size of this variant in WES and imputed data as well as the number of WGS datasets available, there was insufficient power ($\sim 0.5\%$ at $\alpha = 5 \times 10^{-7}$) to robustly detect this association in the WGS data [46].

We next performed rare variant association (RVAS) analyses using a mixed-model version of SKAT-O [47] to test for genes in which the presence of high-impact rare variants might be associated with MT-CN levels (see Methods; Fig. 3, Table 5). Using WES data, the only gene passing the Bonferroni-adjusted P value threshold of 2.16×10^{-6} was *TMBIM1* ($P = 2.96 \times 10^{-8}$), a member of a gene family thought to regulate cell death pathways [48]. *TMBIM1* has been shown to be protective against non-alcoholic fatty liver disease (NAFLD), progression to non-alcoholic steatohepatitis, and insulin resistance in mice and macaques [49]. Interestingly, in our analysis—in which a burden test was determined to be optimal by SKAT-O—rare, putatively high-impact variants in *TMBIM1* were associated with a higher MT-CN (Fig. 3c). Higher MT-CN was, in turn, associated with less severe metabolic syndrome, suggesting that *TMBIM1* is actually a risk gene, not a protective one. Thus, the published function of *TMBIM1* makes it a strong candidate, although the direction of effect in our data disagreed with the direction suggested by prior work in model organisms [49].

Inference of causality in the association between MT-CN and insulin

To further understand the association between MT-CN and fasting serum insulin, we employed a Mendelian randomization (MR) approach with MT-CN as the exposure and insulin as the outcome. Using penalized regression, we leveraged our extensive phenotype data to build a genetic instrument from a large number of genetic variants and adjust for possible confounders via a novel approach (see Methods; Fig. 4). We believe this approach to be more robust to violations of key MR

assumptions than other methods in situations where limited data are available and few robust genotype-exposure associations are known. We restricted our analysis to METSIM samples due to batch effects and inconsistencies in available quantitative trait data observed across FINRISK survey years (Figure S4 - Additional File 1). The effect sizes of the instrument in the causality test for insulin levels are shown in Fig. 4d. We calculated our instrument using either L1 or L2 regularization. In both cases, the MT-CN instrument was not a significant predictor ($\alpha = 0.05$) of insulin when we constructed our instrument from WGS variants, but was significant when the instrument was constructed from imputed array variants. This was likely due to the larger sample size of the imputed array data set. However, the effect estimates were remarkably similar across all four cases. As a result, inverse-variance weighted meta-analysis across datasets yielded highly significant P values for both penalties. In summary, our analysis provided evidence for a significant causal role for MT-CN in determining fasting serum insulin levels that was robust to the choice of regression penalty when building the genetic instrument. We note that this evidence for causality comes with some caveats (see Methods).

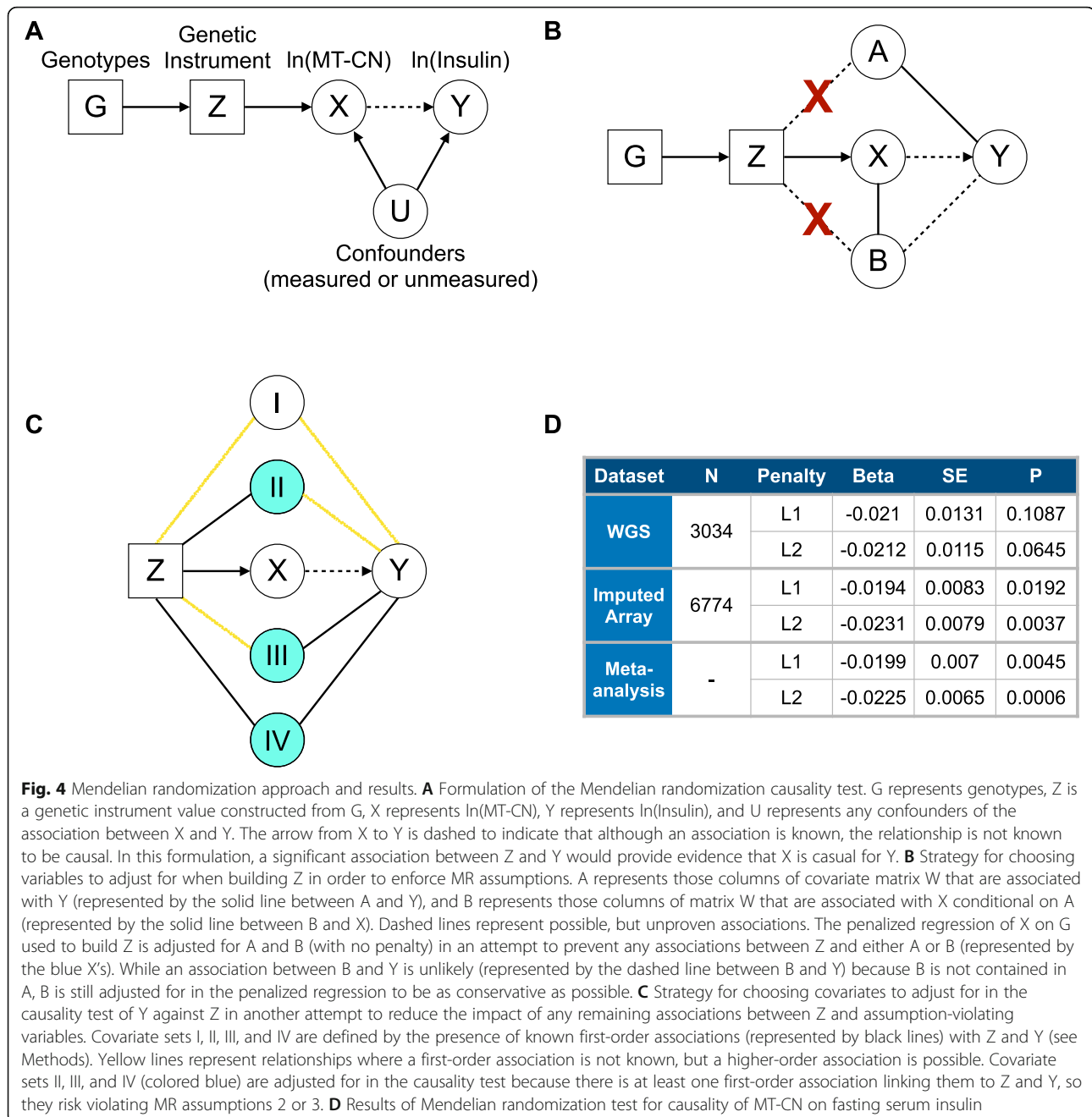
Replication and biological interpretation

In principle, changes in MT-CN can be caused by changes in the number of mitochondrial genome copies within cells or by changes in the blood cell-type composition (see Discussion). Based on the association with rs9389268 and the nuances of the normalization procedure described above, we sought to test the hypothesis that our MT-CN measurement primarily reflects the cell-type composition of the blood rather than the number of mitochondria per cell. Due to the large sample size and rich phenotypic data available in the UK Biobank, we used imputed array genotype and phenotype data from this resource ($N = 357656$) for this purpose [50].

We first tested cell counts from the UK Biobank (UKBB) against a polygenic risk score (PRS) for MT-CN built using the genetic instrument from the Finnish data. Leukocyte, neutrophil, and platelet counts were all significantly associated with MT-CN PRS conditional on

Table 5 Gene-based rare variant association results for *TMBIM1*. *TMBIM1* was the only genome-wide significant gene in the WES rare-variant association tests of METSIM and the whole dataset. Bolded results are significant at the appropriate threshold for the given test (see Methods).

	FINRISK			METSIM			Joint analysis		
	Number	Fraction with rare allele	P	Number	Fraction with rare allele	P	Number	Fraction with rare allele	P
WES	9221	0.016	1.57×10^{-3}	9813	0.013	1.44×10^{-6}	19034	0.014	2.96×10^{-8}
WGS	1084	0.028	0.489	3065	0.014	0.01	4149	0.013	0.01



age, age², and sex (see Methods, Table 6). However, adjusting for neutrophil counts in the leukocyte regression eliminated the signal (PRS regression coefficient $P = 0.839$), suggesting that the leukocyte count signal was driven by the effect of neutrophil count. We removed any high leverage, large residual samples and repeated the neutrophil and platelet count regressions to ensure that this result was robust to outliers and found no appreciable change in significance (Table 6). As a result, we concluded that our MT-CN measurement was significantly associated with neutrophil and platelet counts.

Subsequent analyses were performed both with and without adjustment for these variables, as described below.

We next tested for associations between MT-CN PRS and several cardiometabolic phenotypes from the test in Fig. 1a (see Methods). With the exception of C-reactive protein, which showed no significant association, all tested phenotypes showed nominal association with MT-CN PRS at $\alpha = 0.05$, with total triglycerides and HDL being the only traits surviving Bonferroni correction (Table 7). We interpret this as replication of the

Table 6 Association results between blood cell count traits and MT-CN polygenic risk score in 357656 UK Biobank samples. β refers to the regression coefficient of MT-CN in a linear regression of cell type onto MT-CN PRS and other covariates (see Methods). Bolded results are significant below a nominal $\alpha = 0.05$.

Cell type	All samples			No <i>post hoc</i> high-leverage outliers		
	β	SE	P	β	SE	P
Leukocyte	- 0.00856	0.00170	4.42 × 10⁻⁷	-	-	-
Monocyte	- 0.00119	0.00170	0.482	-	-	-
Lymphocyte	- 0.00250	0.00170	0.142	-	-	-
Neutrophil	- 0.00954	0.00169	1.80 × 10⁻⁸	- 0.00948	0.00169	2.17 × 10⁻⁸
Platelet	0.00548	0.00170	1.24 × 10⁻³	0.00548	0.00170	1.25 × 10⁻³

link between mitochondrial genome copy number and metabolic syndrome in a large, independent data set.

To determine whether there was any association between MT-CN and metabolic syndrome not mediated through cell counts, we repeated the tests of cardiometabolic trait association with MT-CN PRS with adjustment for platelet and neutrophil counts. HDL was the only trait with a nominal ($\alpha = 0.05$) association with PRS under this adjustment, but this signal was not strong enough to survive Bonferroni correction (Table 7). This suggests that the associations we observed between MT-CN and metabolic traits arose simply because MT-CN is a proxy for platelet and neutrophil count. This was supported by the fact that direct testing of platelets and neutrophils against triglycerides, fat mass, and HDL yielded remarkably significant associations, which survived *post hoc* removal of high-leverage, high-residual outlier samples (Table S3 - Additional File 1). This evidence for MT-CN as a proxy for platelet and neutrophil counts strongly suggests that the causal relationship observed in the Mendelian randomization experiment (see above) in fact represents a causative role for neutrophils and platelet counts in setting serum insulin levels.

Given the strong observed associations between blood cell count phenotypes and MT-CN PRS, we used these blood phenotypes to seek replication of the genetic

associations detected in Finnish data. Using imputed UKBB genotype data, we tested the expected alternate allele dosage of both rs2288464 and rs9389268 against the same blood cell traits mentioned above, using linear regression (Table S4 - Additional File 1; expected alternate allele dosage was calculated from genotype call probabilities as $DS = P(0/1) + 2P(1/1)$). As expected given its known associations with multiple hematological parameters (see above), rs9389268 showed strong associations with all tested blood cell phenotypes. rs2288464 was not significantly associated with any of the five phenotypes after correction for multiple testing, although a nominal association was detected with total leukocyte count. This further strengthens our belief that rs9389268 is truly associated with MT-CN through blood cell composition. We also tested *TMBIM1* against the same blood cell traits in UKBB using SKAT-O [47], and found no significant associations (Table S4 - Additional File 1). This may mean that *TMBIM1* affects MT-CN through a mechanism other than altering blood cell-type composition.

As further evidence that MT-CN is a proxy for blood cell composition, we looked up MT-CN association P values in METSIM for the top five neutrophil and platelet count QTLs from the NHGRI-EBI GWAS Catalog [51]. Out of ten variants tested, five had $P < 0.05$ in

Table 7 Association results between metabolic syndrome traits and MT-CN polygenic risk score in 357656 UK Biobank samples. β refers to the regression coefficient of MT-CN in a linear regression of cell type onto MT-CN PRS and other covariates (see Methods). Bolded results are significant below a nominal $\alpha = 0.05$. The weight phenotype tested was that which was measured at the time of impedance measurement.

Trait	Without platelet and neutrophil adjustment			With platelet and neutrophil adjustment		
	β	SE	P	β	SE	P
Type 2 diabetes	- 0.1681	0.0826	0.0419	- 0.1467	0.0844	0.0823
BMI	- 0.0372	0.0175	0.0331	- 0.0233	0.0175	0.1836
Fat mass	- 0.0452	0.0176	0.0100	- 0.0318	0.0177	0.0716
C-reactive protein	- 0.0015	0.0178	0.9305	0.0206	0.0171	0.2291
HDL	0.0573	0.0187	0.0021	0.0416	0.0187	0.0262
Total triglycerides	- 0.0500	0.0176	0.0046	- 0.0330	0.0175	0.0603
Weight	- 0.0359	0.0176	0.0414	- 0.0218	0.0178	0.2189

METSIM (Table S11 - Additional File 1). We note that three of these five were either near or identical to known MT-CN loci (including rs9389268, the marker identified in this study). rs25645, a variant reported to be highly associated with neutrophil count [38], is only 2.5 kb away from rs709591, a SNP with a reported suggestive association with MT-CN [31] and a P value of 1.61×10^{-4} in our METSIM study. Moreover, rs11759553, a platelet-associated variant [38], is 324 kb away from rs9389268, the lead marker for MT-CN in METSIM (rs11759553 P = 2.15×10^{-10} in METSIM). Finally, rs445 was reported as a lead marker for both MT-CN association [29] and platelet count [38]. rs445 has P = 0.048 for association with MT-CN in METSIM. While none of the 10 known cell count-associated markers tested achieved significance beyond a Bonferroni threshold, the overlap between these variants and independently measured MT-CN QTLs was suggestive of a relationship between cell counts and whole blood-derived MT-CN.

Using UKBB data, we further sought to generate hypotheses for other phenotypic associations with MT-CN. To this end, we performed a phenome-wide screen of MT-CN PRS against all of the UKBB phenotypes available to us. To curate and transform these phenotypes, we used a modified version of PHESANT [52, 53], which outputs all continuous variables in both raw and inverse rank normalized form. We chose to interpret the results from the normalized continuous variables (Table S5 - Additional File 2) to be conservative and robust to outliers, although the results of the raw continuous variable analyses were similar (Table S6 - Additional File 2). No metabolic syndrome traits appeared among the tested traits with $q < 0.05$. However, the tests for HDL cholesterol, self-reported heart attack, and doctor-diagnosed heart attack did yield somewhat suggestive results ($q = 0.123, 0.176, \text{ and } 0.176$, respectively). We also repeated this screen with adjustment for neutrophil and platelet counts (Table S7 - Additional File 2 and Table S8 - Additional File 2), resulting again in no metabolic syndrome phenotypes achieving $q < 0.05$. The addition of neutrophil and platelet counts as covariates attenuated the suggestive signals for HDL cholesterol, self-reported heart attack, and doctor-diagnosed heart attack ($q = 0.284, 0.391, \text{ and } 0.402$, respectively).

Discussion

We have described one of the most well-powered studies to date of the genetic relationship between MT-CN measurements in blood and cardiometabolic phenotypes. Our study is one of the very few of which we are aware to utilize WGS data, found to be the most reliable method for estimating MT-CN in a recent study [22], for this purpose. Our data show highly significant associations between blood-derived MT-CN measurements

and several cardiometabolic traits, particularly insulin and fat mass. Anecdotally, it is interesting to note that these MT association signals can also be detected using read-depth analysis of the nuclear genome (Figure S1 - Additional File 1) [54], where reads derived from mtDNA align erroneously to several nuclear loci based on homology between the MT genome and ancient nuclear mitochondrial insertions. This result provides additional evidence for the reported trait associations using an independent MT-CN estimation method, and indicates that these homology-based signals need to be taken into account in future CNV association studies.

We observed strong heritability of MT-CN (31%), on par with other widely studied cardiometabolic traits such as LDL, and identified one single marker association on a haplotype previously associated with several hematological parameters [41–44]. A previous study using qPCR to quantify MT-CN reported two sub-threshold QTLs [31]; of these markers, only rs709591 replicated in our study ($P = 1.61 \times 10^{-4}$). We also report one gene with a rare-variant association with MT-CN, *TMBIM1*, that has a known link to non-alcoholic fatty liver disease [49]. More work is needed to replicate this genetic association.

The association of rs9389268 with MT-CN is not surprising considering our approach for normalizing MT-CN. Because our MT-CN estimate was based on the ratio of mtDNA coverage to nuclear DNA coverage, changes in the cell-type composition of blood could result in changes in our normalized measurement if the underlying cell types have different average numbers of mitochondria. This is especially true of platelets, which can contain mitochondria but not nuclei, and whose counts are known to be associated with rs9399137.

We note that the effect directions of the associations of platelet counts with metS and MT-CN PRS seem inconsistent at first glance, as platelet counts were positively correlated with MT-CN PRS (Table 6) and metS (Table S3 - Additional File 1), while MT-CN and insulin (a proxy for metS) were negatively correlated (Fig. 1b). However, the FinMetSeq regression model in Fig. 1b was not conditional on any other covariates (although age, age², and sex were regressed out of the MT-CN measurement prior to this analysis), while the UKBB models that gave rise to Table 6 and Table S3 (Additional File 1) adjusted for many additional covariates, including 20 PCs and age-sex interaction terms. As a result, the effect directions for the analyses in the two datasets are not directly comparable.

Using a novel multiple-variant instrument-building method, we report evidence from Mendelian randomization supporting a causal role for MT-CN in metabolic syndrome. Further, we used UK Biobank data to show that not only does the link between MT-CN

and metabolic syndrome replicate in an independent data set using a polygenic risk score approach. Contrary to previous claims that variability in the number of mitochondria per cell is responsible for CHD risk [16], this association is mediated by neutrophil and platelet counts.

One important question that our study cannot definitively resolve is the relative contribution of intracellular mitochondrial abundance versus cell-type composition differences in determining the measured MT-CN value. We identified a MT-CN association result at a known QTL for cell-type composition of blood [41–44] (*HBSIL-MYB*), and we further replicated a prior sub-threshold association at a different neutrophil-associated locus [37, 38] (rs709591). Together, these results argue that cell-type composition is an important component of this measurement. On the other hand, two other significant associations from the Finnish dataset (rs2288464, *TMBIM1*) showed no effect on cell-type composition in the UK Biobank. Future work in large cohorts with both WGS and cell count data—which were not simultaneously measured in any samples in this study—will be required to rigorously determine what blood-derived MT-CN primarily measures. However, the results of our MR and UK Biobank analyses together suggest that MT-CN is causally related to metabolic syndrome traits, and that this relationship is mediated by cell-type composition differences.

There is prior evidence to support the role of inflammation—specifically via innate immune cells such as neutrophils—in the etiology of type 2 diabetes (T2D) and insulin resistance [55–57], which suggests a plausible model by which peripheral blood neutrophil count could influence metabolic syndrome. Nutrient excess and high-fat diets are known to recruit neutrophils into tissues, which then cause insulin resistance both by releasing TNF- α and IL-6 and by upregulating cyclooxygenase [55]. This leads to increased LTB4 and subsequent upregulation of NF- κ B, a central regulator of inflammation. Moreover, free fatty acids also cause neutrophils to stay in tissues longer, resulting in persistent inflammation and leading to insulin resistance [55]. While it is known that inflammation, and particularly neutrophils, play a role in metabolic syndrome, our results strongly suggest that peripheral blood neutrophil count causally contributes to this process and is associated with heritable genetic variation in the human population. Overall, our work provides further insight into the role that inflammation plays in metabolic syndrome and supports the idea that targeting inflammation may be a fruitful avenue of investigation in developing future therapeutics.

Conclusions

In summary, we have shown that peripheral blood MT-CN as measured by sequencing is significantly associated with metS and have identified two loci significantly associated with this phenotype, one of which is a novel gene. We have used a Mendelian randomization framework to provide evidence that MT-CN is causal for fasting insulin levels. Finally, we replicated the association between MT-CN and metS in a separate dataset and showed that this association is largely eliminated by adjusting for neutrophil and platelet counts, providing further insight into the role that inflammation plays in metabolic syndrome. Our work uses a large cohort and improved methods to add clarity to a field with often contradictory reports.

Methods

Genotype and phenotype data

Whole genome sequencing (WGS) was performed on a cohort of 4163 samples comprising 3074 male samples from the METSIM study [58] and 1089 male and female samples from the FINRISK study [59]. For details, see Supplementary Methods - *Additional File 3*. Separately, whole exome sequencing (WES) data (N = 19034), genotyping array data (N = 17718) imputed using the Haplotype Reference Consortium panel [60] v1.1, and transformed, normalized quantitative cardiometabolic trait data were obtained from an earlier study [34]. FINRISK array data came in nine genotyping batches, two of which were excluded from the present study due to small sample size. The traits, normalization and transformation procedures, and sample sizes are described in a previous publication [34]. The WES and imputed sample sets contained 4013 and 3929 of the 4163 WGS samples included in the present study.

Mitochondrial genome copy number estimation

We estimated mitochondrial genome copy number (MT-CN) from both WGS and WES data. In WGS data, we used BEDTools [61] to calculate per-base coverage on the mitochondrial genome from the latest available 4163 WGS CRAM files. MT-CN was then calculated by normalizing the mean coverage of the mitochondrial genome to the “haploid coverage” of the autosomes as calculated by Picard [62]. The result was then doubled to account for the diploidy of the autosomal genome. This normalization is summarized by the following equation:

$$\text{MT_CN}_{\text{WGS}} = 2 \times \frac{\text{Mean mtDNA coverage}}{\text{Haploid autosomal coverage}}$$

The output from the above equation served as the raw measurement of per-sample MT-CN. To reduce batch effects, we separated the 4163 samples into three groups:

METSIM, FINRISK collected in 1992 or 1997, and FINRISK collected in 2002 or 2007 (the FINRISK batching decisions were made based on the means shown in *Figure S4 - Additional File 1*). Within each cohort, the raw estimates were regressed on age, age², and sex (FINRISK only), and the residuals were inverse-normal transformed. We combined the three batches of normalized MT-CN values and inverse-normal transformed the combined values for downstream analysis.

We used a similar procedure to estimate MT-CN from WES data, with mean autosomal coverage estimates taken from XHMM [63]. However, as mitochondrial genomic coverage was nonuniform due to the use of hybrid capture probes, mean mtDNA coverage was not an obvious choice of metric for MT-CN estimation (*Figure S5 - Additional File 1*). To summarize this nonuniform mitochondrial genomic coverage into a single number, we tried taking the mean and the maximum depth of reads that aligned to the mitochondrial chromosome; the resulting values were then processed in the same way as the WGS-estimated values. We evaluated the approaches by measuring the R² between WGS-estimated and WES-estimated MT-CN in the 4013 samples for which both data types were available (*Figure S6 - Additional File 1*). While R² was fairly high using both approaches, the maximum coverage method was ultimately selected for use as it yielded a higher R² (0.445 vs 0.380). As a result, the WES MT-CN estimate was calculated as follows:

$$\text{MT-CN}_{\text{WES}} = 2 \times \frac{\text{Maximum mtDNA coverage}}{\text{Mean haploid autosomal coverage}}$$

Mitochondrial haplogroup estimation

We assigned mitochondrial haplogroups using HaploGrep [64] v1.0. Mitochondrial SNP/indel variants were genotyped using GATK GenotypeGVCFs, and a customized filter based on allele balance was applied to the combined callset. HaploGrep was then used to call mitochondrial haplotypes for each individual. We adjusted for major haplogroups in the same linear regressions of metabolic traits onto MT-CN (see Results) and calculated the summary statistics from a permutation test as implemented in the R package lmPerm [65].

Heritability analysis

To estimate heritability of MT-CN, a genomic relatedness-based restricted maximum-likelihood (GREML) method was used as implemented in GCTA [66]. The original GREML [67] method was used first, followed by GREML-LDMS [68] to account for biases arising from differences in minor allele frequency (MAF)

spectrum or linkage disequilibrium (LD) properties between the genotyped variants and the true causal variants [69]. For both analyses, MT-CN values were normalized and residualized for sex, age, and age² as described above. Heritability estimation was performed jointly and separately for METSIM and FINRISK samples using WGS and imputed array genotypes. In all cases, a minimum MAF threshold of 1% was applied. Beyond those covariates already adjusted for in the normalization process, sensitivity analyses were performed on imputed array data to determine whether heritability estimates were sensitive to inclusion of covariates. In these experiments, either cohort or FINRISK genotyping array batch were included as fixed-effect covariates in joint analyses of imputed array data; in neither case was the final heritability estimate significantly affected ($h^2 = 0.09$, SE = 0.02 in both cases). In GREML-LDMS, genotypes were split into four SNP-based LD score quartiles and two MAF bins (1% > MAF > 5% and MAF > 5%), and genetic relatedness matrices (GRMs) were estimated separately for each of the eight combinations. The GREML algorithm was then run on all eight GRMs simultaneously using the first ten principal components (PCs) of the genotype matrix (as calculated by smartPCA [70] v13050) as fixed covariates [68]. The use of GREML-LDMS over GREML also did not affect estimated heritability values (*Table 3*), suggesting that the properties of the causal variants for this trait do not lead to significant biases when using the standard GREML approach.

We observed that WGS heritability estimates are lower when analyzing FINRISK and METSIM data together as compared with the analysis of METSIM alone (*Table 2*) (note that FINRISK-only heritability estimates are not reliable as they have large standard errors resulting from the small number of FINRISK samples sequenced). One potential explanation for this is that there exists substantial heterogeneity across FINRISK survey years (*Figure S4 - Additional File 1*), and between the FINRISK and METSIM cohorts, with respect to the reliability with which mtDNA was captured (likely due to different DNA preparation protocols).

Genome-wide association analyses

Genome-wide association studies (GWAS) were performed using the same normalized phenotype used in heritability analyses. Single-variant GWAS were conducted using EMMAX [71] as implemented in EPACTS [72]. Kinship matrices required by EMMAX were generated by EPACTS; kinship matrices for WGS GWAS were generated from WGS data, while those for WES and imputed array-based GWAS were generated from WES data. A P value threshold of 5×10^{-8} was used for the WGS and imputed array GWAS, while 5×10^{-7} was

used for significance in the WES GWAS. Single-variant association analyses of WGS and WES data did not include any covariates in the EMMAX model, although all association analyses were performed using MT-CN values that adjusted for age, age², sex, and cohort (see “Mitochondrial genome copy number estimation” section). All association tests labeled “joint” were performed on METSIM and FINRISK cohorts together; in one case, a random-effects meta-analysis was performed using individual-cohort summary statistics and the R package meta [73].

Gene-based variant aggregation studies (RVAS) were done using a mixed-model version of SKAT-O [47] as implemented in EPACTS. Variants with CADD [74] score greater than 20 and minor allele frequency less than 1% were grouped into genes as annotated by VEP [75] (which by default annotates a variant with a gene name if the gene falls within 5 kb of that gene by default). For gene-based RVAS, Bonferroni-corrected genome-wide significance thresholds varied slightly due to differing the number of genes with at least two variants meeting the above criteria in each test, but were approximately 2×10^{-6} in all cases.

Mendelian randomization

To assess the evidence for a causal relationship between mitochondrial genome copy number and fasting serum insulin levels, the METSIM cohort alone was used due to its homogeneity of sex, collection procedures, and location. A penalized regression-based, multiple-variant Mendelian randomization (MR) approach was employed to enforce the necessary assumptions of MR methods. While some MR studies have tested one or more assumptions *post hoc*, to our knowledge, there is no published method that tries to enforce these assumptions during the process of building a genetic instrument from multiple variants in the absence of a large set of known genotype-exposure associations. In our formulation (Fig. 4), X , the natural log of MT-CN (adjusted for nuclear genomic coverage but not for age, age², or sex), and a genotype matrix G were used to build a genetic instrument Z , which was then tested against Y , the natural log of fasting serum insulin. The goal of the MR approach was to use a large number of common variants to build a genetic instrument Z that satisfies the three assumptions of MR [76] (see Supplementary Methods - Additional File 3).

To account for missing phenotype data, missing values were multiply imputed using regression trees as implemented in the R package mice [77] v3.4.0 (`maxit = 25`). This imputation was repeated 1000 times in parallel, with each set of imputed values being carried through the entire procedure described above. The resulting 1000 computed instrument effect sizes and standard

errors were combined using Rubin’s method as implemented in the R package Amelia [78] v1.7.5. The combined effect size and standard error were then tested for significance using a t test with 998 degrees of freedom.

The above procedure was performed separately for METSIM samples with WGS data ($N = 3034$) and METSIM samples with only imputed array data ($N = 6774$) using an L1 penalty in the instrument-building regression, and again using an L2 penalty. Both sample sets were limited to those for which relevant quantitative traits were available. An inverse-variance weighted meta-analysis was performed across data sets for L1- and L2-penalized regression separately. The resulting effect size and standard error were tested for significance using a Z test.

To ensure that our results were not driven by outlier samples, we removed outliers in two stages. Before the MR analyses, we used principal components analysis (PCA), Mahalanobis distance, and multi-trait extreme outlier identification to remove 5 WGS samples and 15 imputed array samples based on quantitative trait data. We also removed high-leverage, high-residual outliers from the causality test regression (see below) *post hoc* and recomputed the instrument effect sizes to ensure that there was no significant change in the results. In each of the 1000 multiple imputation runs, among the samples with standardized residual greater than 1, the top 10 samples by leverage were recorded. Any sample that was recorded in this way in at least one run was then excluded from the re-analysis as a *post hoc* outlier. The results of this additional analysis showed only very small differences in effect estimates, and their interpretation remained the same (Table S9 - Additional File 1). Thus, we concluded that our causal inference results were not driven by outlier samples.

One caveat of this method is that, as mentioned above, exclusion of sets A and B from the regression penalty (see Supplementary Methods - Additional File 3) did not perfectly orthogonalize the resulting instrument from these variables in practice (Figure S7 - Additional File 1). Reasons for this may include relatively low levels of shrinkage in the instrument-building regression or higher order associations between MT-CN and the confounding variables. However, our method still represents an improvement over the current standard, which is not to adjust for these covariates at all. Another caveat is that it is impossible to determine the perfect set of covariates for which adjustment is appropriate. Lack of adjustment for truly confounding variables can result in an instrument which does not satisfy MR assumptions 2 and/or 3 (see Supplementary Methods - Additional File 3), yielding a biased effect estimate. Conversely, unnecessary adjustment for certain variables can also result in biases. For example, adjusting for an intermediate

phenotype that truly lies along the path from Z to X to Y can cause a false-negative signal, making the causality test overly conservative. Alternatively, adjusting for some variables can result in collider biases [79]. That is, if both Z and Y are causal for a confounder U, then adjusting for U can induce a dependency between Z and Y (*Figure S8 - Additional File 1*) that did not previously exist.

We note that a known source of bias in MR studies is the selection of samples based on case–control status for a related disease [80]. While METSIM is a population-based study, samples were selected for WGS based on cardiovascular disease case–control status so as to enrich the sequenced samples for cases. This has the potential to bias a MR experiment if both the exposure and the outcome are associated with the disease, which is certainly possible. However, in our design, all of the METSIM samples not chosen for WGS were tested in the imputed array experiment. The consistency of effect estimates between the WGS and imputed array samples both in the L1 and L2 penalty cases (*Fig. 4d*) suggests that there is little to no bias arising from sample selection in this experiment.

Calculation and testing of polygenic risk score in the UK Biobank

To search for associations between MT-CN and other phenotypes, the genetic instrument calculated in Finnish imputed array data was computed and treated as a polygenic risk score (PRS) in a relatively homogenous subset of 357656 UK Biobank samples identified by a previous study [53]. We calculated $\overline{\beta_G}$, the average of the five values of $\beta_G^{(i)}$ across all 1000 multiple imputation runs using an L2 penalty and imputed array data—the L2 penalty was chosen because it performed better than the L1 on both METSIM data types, and the imputed array data set was chosen due to its larger sample size than the WGS set (*Fig. 4d*). Next, to keep the procedure as consistent as possible with the imputation protocol used for METSIM—which used haploid dosage values to call imputed genotypes [34]—we called imputed genotypes using the expected alternate allele dosage from the UK Biobank by setting thresholds of 0.5 and 1.5. Using the resulting imputed variant calls, we calculated our PRS as $\tilde{Z} = \overline{\beta_G} \times \tilde{G}$, where \tilde{G} is the UK Biobank genotype dosage matrix constructed in the same way as G in METSIM.

To test for associations with MT-CN PRS in the UK Biobank, we employed two approaches: a hypothesis-driven analysis targeted to the phenotypes associated with MT-CN in the Finnish data as well as a hypothesis-free screen of all the phenotypes available to us.

In the targeted analysis, we used our genetic instrument from the MR experiment as a PRS for MT-CN in our chosen subset of the UK Biobank and tested for

associations with several blood cell count and metabolic syndrome traits. Given the association of MT-CN with rs9389268 (see Results), we selected as cell count traits total leukocyte count as well as lymphocyte, neutrophil, monocyte, and platelet counts for testing (because lymphocyte count was not readily available, it was calculated as the product of leukocyte count and lymphocyte percentage). We did not include basophils and eosinophils in this analysis considering that they comprise a small minority of white blood cells and are unlikely to affect MT-CN measured from whole blood. All cell count traits were log-transformed and standardized separately by sex.

We took several steps to eliminate outlier samples in the dataset. Through three iterations of PCA on the cell count matrix and subsequent outlier removal, we removed 1637 outlier samples. We then fit null linear models of the form $Cell\ count \sim Age + Age^2 + Sex + Age : Sex + Age^2 : Sex + PCs$ (the first 20 PCs were included) for each cell count trait and subsequently removed samples with either large residuals or high leverage and moderate residuals in at least one model (following the example of [53]). Through two iterations of null model fitting and outlier removal, we removed 7 additional samples based on null model fit. Tests of association between cell counts and MT-CN PRS were based on the PRS regression coefficient in linear models of the form $Cell\ count \sim PRS + Age + Age^2 + Sex + Age : Sex + Age^2 : Sex + PCs$.

We repeated this process for those cardiometabolic traits found to be suggestively associated with MT-CN in the Finnish dataset ($P < 10^{-6}$) that were also readily available in the UK Biobank (*Fig. 1a, Table S1 - Additional File 2*); these phenotypes were body mass index (BMI), fat mass, C-reactive protein, high-density lipoprotein, total triglycerides, and weight. We also chose to include T2D status because of the lack of insulin measurement in the UK Biobank. Except for T2D, a binary trait, all traits were log-transformed before further analysis (after removing 817 samples with negative values for T2D, representing missing information). The above outlier removal steps were repeated for the cardiometabolic traits after excluding the outliers already identified from the cell count data, with the only major modification being the use of logistic regression for the T2D models. This process resulted in the removal of 42 and 53 samples from PCA and null model fitting, respectively.

SKAT-O tests of association between *TMBIM1* and the cell count traits identified above were also performed. Similarly to the RVAS in Finnish data, variants within 5 kb of *TMBIM1* with MAF < 1% and CADD v1.6 score > 20 were selected for inclusion in this analysis. Rather than the mixed-model version of SKAT-O

used in the Finnish data, standard SKAT-O was used due to the lower expected level of cryptic relatedness in the UK Biobank population.

We also performed a hypothesis-free, phenome-wide screen of UK Biobank traits to which we had access (*Table S10 - Additional File 2*), to search for other associations with MT-CN PRS. The statistical models used in this screen were of the same form as those described above, both with and without adjustment for neutrophil and platelet counts. To curate and transform phenotypes, we used an adapted version of PHESANT [52, 53]. A few further modifications were made to the pipeline, the most significant being the direct use of logistic regression for testing categorical unordered variables, the inclusion of cancer phenotypes, and the exclusion of sex-specific (or nearly sex-specific) categorical traits. The PHESANT pipeline we used [53] outputs continuous variables both in their raw form and after applying an inverse rank normal transformation. For the sake of being conservative and robust to outliers, we chose to interpret the results from the normalized continuous variables. To control false discovery rate, we performed a Benjamini-Hochberg procedure with Storey correction as implemented in the R package *q* value [81] v2.18.0 on the categorical and normalized continuous variables together. As a secondary analysis, this same correction was applied to the categorical and raw continuous variables together.

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s40246-021-00335-2>.

Additional file 1: PDF format. Supplementary Figures and Small Supplementary Tables. Contains Figures S1-S8 and Tables S2-S4, S9, and S11.

Additional file 2: Excel format (.xlsx). Large Supplementary Tables. Contains Tables S1, S5-S8, and S10.

Additional file 3: PDF format. Supplementary Methods. Describes methods in further detail than the main text.

Acknowledgements

We thank Hyeim Jung for her help in identifying outlier individuals as well as the WashU data production team, in particular Robert Fulton, Lucinda Fulton, Catrina Fronick, Aye Wollam, and Susan K. Dutcher. This research has been conducted using the UK Biobank Resource under Application Number 56546. The FINRISK samples used for the research were obtained from the FINRISK Study and from THL Biobank. We thank all study participants for their generous participation at THL Biobank and the FINRISK Study.

Authors' contributions

LG led trait association, heritability, and MR analyses. LC discovered the original links between the nuclear genome, MT-CN, and metabolic traits and developed the MT-CN estimation method. RC and LG conceived the MR covariate correction approach, and RC provided critical guidance on statistical methods. HJA, DL, ID, AR, and LG led variant callset creation and QC. JV and ML performed the analysis of insulin sensitivity and resistance. AP, SR, ML, JK, and AH contributed samples and phenotypic data. All authors edited the manuscript and/or provided intellectual contributions. IMH and NOS conceived the study and supervised the overall project. LG and IMH wrote the

paper, with contributions from LC. The author(s) read and approved the final manuscript.

Funding

This work was funded by the NHGRI Centers for Common Disease Genetics (grant number UM1 HG008853 to IMH and NOS), the NHGRI large-scale sequencing grant (grant number 5U54HG003079), the Sigrid Jusélius Foundation (to SR), the University of Helsinki HiLIFE Fellow grants 2017–2020 (to SR), the Academy of Finland Center of Excellence in Complex Disease Genetics (grant number 312062 to SR), the Academy of Finland (grant number 285380 to SR), the National Heart, Lung and Blood Institute (grant number T32HL007081 to EY), and the National Center for Advancing Translational Sciences (grant number UL1TR002345 to EY). The funders had no role in the design of the study or the collection, analysis, or interpretation of the data.

Availability of data and materials

METSIM WGS, METSIM WES, and FINRISK WES sequence data are available through dbGaP (accession numbers phs001579, phs000752, and phs000756). METSIM callsets from WGS and imputed array data as well as MT-CN phenotype values will soon be available through AnVIL. Imputed array GWAS summary statistics from METSIM and WES SKAT-O summary statistics from the joint dataset are freely available at <https://wustl.box.com/s/7xfbmqxq2r4kg8p8bfc7vpqlmqvhm0lx>. Genomic and phenotypic data for the FINRISK cohort are obtainable through THL Biobank, the Finnish Institute for Health and Welfare, Finland (<https://thl.fi/en/web/thl-biobank>).

Declarations

Ethics approval and consent to participate

All participants in both the METSIM and FINRISK studies provided informed consent, and study protocols were approved by the Ethics Committees at participating institutions (National Public Health Institute of Finland, Hospital District of Helsinki and Uusimaa, and Hospital District of Northern Savo). All relevant ethics committees approved this study.

Consent for publication

Not applicable

Competing interests

NOS has received grant funding from Regeneron Pharmaceuticals for unrelated work. The rest of the authors declare no competing interests.

Author details

¹McDonnell Genome Institute, Washington University School of Medicine, St. Louis, MO, USA. ²Department of Medicine, Washington University School of Medicine, St. Louis, MO, USA. ³Institute of Clinical Medicine, Internal Medicine, University of Eastern Finland, Kuopio, Finland. ⁴Department of Medicine, Cardiovascular Division, Washington University School of Medicine, St. Louis, MO, USA. ⁵Department of Genetics, Washington University School of Medicine, St. Louis, MO, USA. ⁶Institute for Molecular Medicine Finland (FIMM), HiLIFE, University of Helsinki, Helsinki, Finland. ⁷Finnish Institute for Health and Welfare (THL), Helsinki, Finland. ⁸Center for Genetic Epidemiology, Department of Preventive Medicine, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA. ⁹Quantitative and Computational Biology Section, Department of Biological Sciences, University of Southern California, Los Angeles, CA, USA. ¹⁰Center for Neurobehavioral Genetics, Jane and Terry Semel Institute for Neuroscience and Human Behavior, University of California Los Angeles, Los Angeles, CA, USA. ¹¹Analytical and Translational Genetics Unit (ATGU), Psychiatric & Neurodevelopmental Genetics Unit, Departments of Psychiatry and Neurology, Massachusetts General Hospital, Boston, MA, USA. ¹²Broad Institute of MIT and Harvard, Cambridge, MA, USA. ¹³Faculty of Medicine, University of Helsinki, Helsinki, Finland. ¹⁴Department of Medicine, Kuopio University Hospital, Kuopio, Finland. ¹⁵Department of Biostatistics and Center for Statistical Genetics, University of Michigan School of Public Health, Ann Arbor, MI, USA. ¹⁶Department of Genetics, Yale University School of Medicine, New Haven, CT, USA.

Received: 23 March 2021 Accepted: 26 May 2021

Published online: 07 June 2021

References

- Virani SS, Alonso A, Benjamin EJ, Bittencourt MS, Callaway CW, Carson AP, et al. Heart disease and stroke statistics—2020 update: a report from the American Heart Association. *Circulation*. American Heart Association. 2020; 141:e139–596.
- University of Washington Institute for Health Metrics and Evaluation. GBD results tool. Global Health Data Exchange. [cited 2021 Feb 10]. Available from: <http://ghdx.healthdata.org/gbd-results-tool>
- Grundy SM, Brewer HB Jr, Cleeman JI, Smith SC Jr, Lenfant C. American Heart Association, et al. Definition of metabolic syndrome: report of the National Heart, Lung, and Blood Institute/American Heart Association conference on scientific issues related to definition. *Circulation*. 2004;109:433–8.
- Samson SL, Garber AJ. Metabolic syndrome. *Endocrinol Metab Clin North Am*. 2014;43(1):1–23. <https://doi.org/10.1016/j.ecl.2013.09.009>.
- Koliaki C, Roden M. Alterations of mitochondrial function and insulin sensitivity in human obesity and diabetes mellitus | Annual Review of Nutrition. *Annu Rev Nutr*. 2016;36(1):337–67. <https://doi.org/10.1146/annurev-nutr-071715-050656>.
- Weisberg SP, McCann D, Desai M, Rosenbaum M, Leibel RL, Ferrante AW. Obesity is associated with macrophage accumulation in adipose tissue. *J Clin Invest*. 2003;112(12):1796–808. <https://doi.org/10.1172/JCI200319246>.
- Kim J-A, Wei Y, Sowers JR. Role of Mitochondrial dysfunction in insulin resistance. *Circ Res*. 2008;102(4):401–14. <https://doi.org/10.1161/CIRCRESA.HA.107.165472>.
- Burgueño AL, Cabrerizo R, Gonzales Mansilla N, Sookoian S, Pirola CJ. Maternal high-fat intake during pregnancy programs metabolic-syndrome-related phenotypes through liver mitochondrial DNA copy number and transcriptional activity of liver PPARGC1A. *J Nutr Biochem*. 2013;24(1):6–13. <https://doi.org/10.1016/j.jnutbio.2011.12.008>.
- Sookoian S, Rosselli MS, Gemma C, Burgueño AL, Gianotti TF, Castaño GO, et al. Epigenetic regulation of insulin resistance in nonalcoholic fatty liver disease: impact of liver methylation of the peroxisome proliferator-activated receptor γ coactivator 1 α promoter. *Hepatology*. 2010;52(6):1992–2000. <https://doi.org/10.1002/hep.23927>.
- Begrache K, Igoudjil A, Pessayre D, Fromenty B. Mitochondrial dysfunction in NASH: causes, consequences and possible means to prevent it. *Mitochondrion*. 2006;6(1):1–28. <https://doi.org/10.1016/j.mito.2005.10.004>.
- Zhou X, Li R, Liu X, Wang L, Hui P, Chan L, et al. ROCK1 reduces mitochondrial content and irisin production in muscle suppressing adipocyte browning and impairing insulin sensitivity. *Sci Rep*. 2016;6(1):29669. <https://doi.org/10.1038/srep29669>.
- Ren J, Pulakat L, Whaley-Connell A, Sowers JR. Mitochondrial biogenesis in the metabolic syndrome and cardiovascular disease. *J Mol Med*. 2010;88(10):993–1001. <https://doi.org/10.1007/s00109-010-0663-9>.
- Stephenson EJ, Hawley JA. Mitochondrial function in metabolic health: a genetic and environmental tug of war. *Biochimica et Biophysica Acta (BBA) - General Subjects*. 2014;1840(4):1285–94. <https://doi.org/10.1016/j.bbagen.2013.12.004>.
- Szendroedi J, Phielix E, Roden M. The role of mitochondria in insulin resistance and type 2 diabetes mellitus. *Nat Rev Endocrinol*. 2012;8(2):92–103. <https://doi.org/10.1038/nrendo.2011.138>.
- Ding J, Sidore C, Butler TJ, Wing MK, Qian Y, Meirelles O, et al. Assessing mitochondrial DNA variation and copy number in lymphocytes of ~2,000 Sardinians using tailored sequencing analysis tools. *PLoS Genet*. 2015;11(7):e1005306. <https://doi.org/10.1371/journal.pgen.1005306>.
- Chen S, Xie X, Wang Y, Gao Y, Xie X, Yang J, et al. Association between leukocyte mitochondrial DNA content and risk of coronary heart disease: a case-control study. *Atherosclerosis*. 2014;237(1):220–6. <https://doi.org/10.1016/j.atherosclerosis.2014.08.051>.
- Lee HK, Song JH, Shin CS, Park DJ, Park KS, Lee KU, et al. Decreased mitochondrial DNA content in peripheral blood precedes the development of non-insulin-dependent diabetes mellitus. *Diabetes Res Clin Pract*. 1998; 42(3):161–7. [https://doi.org/10.1016/S0168-8227\(98\)00110-7](https://doi.org/10.1016/S0168-8227(98)00110-7).
- Shoar Z, Goldenthal MJ, De Luca F, Suarez E. Mitochondrial DNA content and function, childhood obesity, and insulin resistance. *Endocr Res*. 2016; 41(1):49–56. <https://doi.org/10.3109/07435800.2015.1068797>.
- Song J, Oh JY, Sung Y-A, Pak YK, Park KS, Lee HK. Peripheral blood mitochondrial DNA content is related to insulin sensitivity in offspring of type 2 diabetic patients. *Diabetes Care*. 2001;24(5):865–9. <https://doi.org/10.2337/diacare.24.5.865>.
- Weng S-W, Lin T-K, Liou C-W, Chen S-D, Wei Y-H, Lee H-C, et al. Peripheral blood mitochondrial DNA content and dysregulation of glucose metabolism. *Diabetes Res Clin Pract*. 2009;83(1):94–9. <https://doi.org/10.1016/j.diabres.2008.10.002>.
- Liu L-P, Cheng K, Ning M-A, Li H-H, Wang H-C, Li F, et al. Association between peripheral blood cells mitochondrial DNA content and severity of coronary heart disease. *Atherosclerosis*. 2017;261:105–10. <https://doi.org/10.1016/j.atherosclerosis.2017.02.013>.
- Longchamps RJ, Castellani CA, Yang SY, Newcomb CE, Sumpter JA, Lane J, et al. Evaluation of mitochondrial DNA copy number estimation techniques. *PLoS One*. 2020;15(1):e0228166. <https://doi.org/10.1371/journal.pone.0228166>.
- Guyatt AL, Burrows K, Guthrie PAI, Ring S, McArdle W, Day INM, et al. Cardiometabolic phenotypes and mitochondrial DNA copy number in two cohorts of UK women. *Mitochondrion*. 2018;39:9–19. <https://doi.org/10.1016/j.mito.2017.08.007>.
- Ashar FN, Zhang Y, Longchamps RJ, Lane J, Moes A, Grove ML, et al. Association of mitochondrial DNA copy number with cardiovascular disease. *JAMA Cardiol*. 2017;2(11):1247–55. <https://doi.org/10.1001/jamacardio.2017.3683>.
- Robin ED, Wong R. Mitochondrial DNA molecules and virtual number of mitochondria per cell in mammalian cells. *J Cell Physiol*. 1988;136(3):507–13. <https://doi.org/10.1002/jcp.1041360316>.
- Maianski NA, Geissler J, Srinivasula SM, Alnemri ES, Roos D, Kuipers TW. Functional characterization of mitochondria in neutrophils: a role restricted to apoptosis. *Cell Death Differ*. 2004;11(2):143–53. <https://doi.org/10.1038/sj.cdd.4401320>.
- Zharikov S, Shiva S. Platelet mitochondrial function: from regulation of thrombosis to biomarker of disease. *Biochem Soc Trans*. 2013;41(1):118–23. <https://doi.org/10.1042/BST20120327>.
- Cai N, Chang S, Li Y, Li Q, Hu J, Liang J, et al. Molecular signatures of major depression. *Curr Biol*. 2015;25(9):1146–56. <https://doi.org/10.1016/j.cub.2015.03.008>.
- Cai N, Li Y, Chang S, Liang J, Lin C, Zhang X, et al. Genetic control over mtDNA and its relationship to major depressive disorder. *Curr Biol*. 2015; 25(24):3170–7. <https://doi.org/10.1016/j.cub.2015.10.065>.
- Curran JE, Johnson MP, Dyer TD, Göring HHH, Kent JW, Charlesworth JC, et al. Genetic determinants of mitochondrial content. *Hum Mol Genet*. 2007;16(12):1504–14. <https://doi.org/10.1093/hmg/ddm101>.
- Guyatt AL, Brennan RR, Burrows K, Guthrie PAI, Ascione R, Ring SM, et al. A genome-wide association study of mitochondrial DNA copy number in two population-based cohorts. *Hum Genomics*. 2019;13(1):6. <https://doi.org/10.1186/s40246-018-0190-2>.
- Pajukanta P, Terwilliger JD, Perola M, Hiekkalinna T, Nuotio I, Ellonen P, et al. Genomewide scan for familial combined hyperlipidemia genes in Finnish families, suggesting multiple susceptibility loci influencing triglyceride, cholesterol, and apolipoprotein B levels. *Am J Hum Genet*. 1999;64(5):1453–63. <https://doi.org/10.1086/302365>.
- Pajukanta P, Allayee H, Krass KL, Kuraishy A, Soro A, Lilja HE, et al. Combined analysis of genome scans of dutch and finnish families reveals a susceptibility locus for high-density lipoprotein cholesterol on chromosome 16q. *Am J Hum Genet*. 2003;72(4):903–17. <https://doi.org/10.1086/374177>.
- Locke AE, Steinberg KM, Chiang CWK, Service SK, Havulinna AS, Stell L, et al. Exome sequencing of Finnish isolates enhances rare-variant association power. *Nature*. 2019;572(7769):323–8. <https://doi.org/10.1038/s41586-019-1457-z>.
- Kaess B, Fischer M, Baessler A, Stark K, Huber F, Kremer W, et al. The lipoprotein subfraction profile: heritability and identification of quantitative trait loci. *J Lipid Res*. 2008;49(4):715–23. <https://doi.org/10.1194/jlr.M700338-JLR200>.
- Weiss LA, Pan L, Abney M, Ober C. The sex-specific genetic architecture of quantitative traits in humans. *Nat Genet*. 2006;38(2):218–22. <https://doi.org/10.1038/ng1726>.
- Nalls MA, Couper DJ, Tanaka T, van Rooij FJA, Chen M-H, Smith AV, et al. Multiple loci are associated with white blood cell phenotypes. *PLoS Genet*. 2011;7(6):e1002113. <https://doi.org/10.1371/journal.pgen.1002113>.
- Chen M-H, Raffield LM, Mousas A, Sakae S, Huffman JE, Moscati A, et al. Trans-ethnic and ancestry-specific blood-cell genetics in 746,667 individuals from 5 global populations. *Cell*. 2020;182:1198–1213.e14.

39. Wang X, Angelis N, Thein SL. MYB - a regulatory factor in hematopoiesis. *Gene*. 2018;665:6–17. <https://doi.org/10.1016/j.gene.2018.04.065>.
40. Pandit RA, Svasti S, Sripichai O, Munkongdee T, Triwitayakorn K, Winichagoon P, et al. Association of SNP in exon 1 of HBS1L with hemoglobin F level in beta0-thalassemia/hemoglobin E. *Int J Hematol*. 2008; 88(4):357–61. <https://doi.org/10.1007/s12185-008-0167-3>.
41. Thein SL, Menzel S, Peng X, Best S, Jiang J, Close J, et al. Intergenic variants of HBS1L-MYB are responsible for a major quantitative trait locus on chromosome 6q23 influencing fetal hemoglobin levels in adults. *Proceedings of the National Academy of Sciences*. 2007;104(27):11346–51. <https://doi.org/10.1073/pnas.0611393104>.
42. Ganesh SK, Zakai NA, van Rooij FJA, Soranzo N, Smith AV, Nalls MA, et al. Multiple loci influence erythrocyte phenotypes in the CHARGE Consortium. *Nat Genet*. 2009;41(11):1191–8. <https://doi.org/10.1038/ng.466>.
43. Menzel S, Jiang J, Silver N, Gallagher J, Cunningham J, Surdulescu G, et al. The HBS1L-MYB intergenic region on chromosome 6q23.3 influences erythrocyte, platelet, and monocyte counts in humans. *Blood*. 2007;110(10): 3624–6. <https://doi.org/10.1182/blood-2007-05-093419>.
44. Lin BD, Carnero-Montoro E, Bell JT, Boomsma DI, de Geus EJ, Jansen R, et al. 2SNP heritability and effects of genetic variants for neutrophil-to-lymphocyte and platelet-to-lymphocyte ratio. *J Hum Genet*. 2017;62(11): 979–88. <https://doi.org/10.1038/jhg.2017.76>.
45. Stadhouders R, Aktuna S, Thongjuea S, Aghajani-refah A, Pourfarzad F, van Ijcken W, et al. HBS1L-MYB intergenic variants modulate fetal hemoglobin via long-range MYB enhancers. *J Clin Invest*. 2014;124(4):1699–710. <https://doi.org/10.1172/JCI171520>.
46. Purcell S, Cherny SS, Sham PC. Genetic Power Calculator: design of linkage and association genetic mapping studies of complex traits. *Bioinformatics*. 2003;19(1):149–50. <https://doi.org/10.1093/bioinformatics/19.1.149>.
47. Lee S, Emond MJ, Bamshad MJ, Barnes KC, Rieder MJ, Nickerson DA, et al. Optimal unified approach for rare-variant association testing with application to small-sample case-control whole-exome sequencing studies. *Am J Hum Genet*. 2012;91(2):224–37. <https://doi.org/10.1016/j.ajhg.2012.06.007>.
48. Zhou J, Zhu T, Hu C, Li H, Chen G, Xu G, et al. Comparative genomics and function analysis on B11 family. *Comput Biol Chem*. 2008;32(3):159–62. <https://doi.org/10.1016/j.compbiolchem.2008.01.002>.
49. Zhao G-N, Zhang P, Gong J, Zhang X-J, Wang P-X, Yin M, et al. Tmbim1 is a multivesicular body regulator that protects against non-alcoholic fatty liver disease in mice and monkeys by targeting the lysosomal degradation of Tlr4. *Nat Med*. 2017;23(6):742–52. <https://doi.org/10.1038/nm.4334>.
50. Sudlow C, Gallacher J, Allen N, Beral V, Burton P, Danesh J, et al. UK biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLoS Med*. 2015;12(3):e1001779. <https://doi.org/10.1371/journal.pmed.1001779>.
51. Buniello A, MacArthur JAL, Cerezo M, Harris LW, Hayhurst J, Malangan C, et al. The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res*. Oxford University Press (OUP). 2019;47:D1005–12.
52. Millard LAC, Davies NM, Gaunt TR, Davey Smith G, Tilling K. Software Application Profile: PHESANT: a tool for performing automated phenome scans in UK Biobank. *Int J Epidemiol*. 2018;47(1):29–35. <https://doi.org/10.1093/ije/dyx204>.
53. Neale B. Neale Lab UK Biobank analysis [Internet]. Available from: <http://www.nealelab.is/uk-biobank/>
54. Chen, L. et al. Association of structural variation with cardiometabolic traits in Finns. *Am J Hum Genet*. 2021;108:583–96.
55. Welty FK, Alfaddagh A, Elajami TK. Targeting inflammation in metabolic syndrome. *Transl Res*. 2016;167(1):257–80. <https://doi.org/10.1016/j.trsl.2015.06.017>.
56. Creely SJ, McTernan PG, Kusminski CM, Fisher ff M, Da Silva NF, Khanolkar M, et al. Lipopolysaccharide activates an innate immune system response in human adipose tissue in obesity and type 2 diabetes. *Am J Physiol Endocrinol Metab*. 2007;292(3):E740–7. <https://doi.org/10.1152/ajpendo.003.02.2006>.
57. Hotamisligil GS. Inflammation, metaflammation and immunometabolic disorders. *Nature*. 2017;542(7640):177–85. <https://doi.org/10.1038/nature21363>.
58. Laakso M, Kuusisto J, Stančáková A, Kuulasmaa T, Pajukanta P, Lusin AJ, et al. The metabolic syndrome in men study: a resource for studies of metabolic and cardiovascular diseases. *J Lipid Res*. 2017;58(3):481–93. <https://doi.org/10.1194/jlr.O072629>.
59. Borodulin K, Tolonen H, Jousilahti P, Jula A, Juolevi A, Koskinen S, et al. Cohort profile: the national FINRISK study. *Int J Epidemiol*. 2018;47:696–696i.
60. McCarthy S, Das S, Kretzschmar W, Delaneau O, Wood AR, Teumer A, et al. A reference panel of 64,976 haplotypes for genotype imputation. *Nat Genet*. 2016;48(10):1279–83. <https://doi.org/10.1038/ng.3643>.
61. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*. 2010;26(6):841–2. <https://doi.org/10.1093/bioinformatics/btq033>.
62. Picard Toolkit. Broad Institute, GitHub repository; 2019. Available from: <http://broadinstitute.github.io/picard/>
63. Fromer M, Moran JL, Chambert K, Banks E, Bergen SE, Ruderfer DM, et al. Discovery and statistical genotyping of copy-number variation from whole-exome sequencing depth. *Am J Hum Genet*. 2012;91(4):597–607. <https://doi.org/10.1016/j.ajhg.2012.08.005>.
64. Kloss-Brandstätter A, Pacher D, Schönherr S, Weissensteiner H, Binna R, Specht G, et al. HaploGrep: a fast and reliable algorithm for automatic classification of mitochondrial DNA haplogroups. *Hum Mutat*. 2011;32(1):25–32. <https://doi.org/10.1002/humu.21382>.
65. Wheeler B, Torchiano M. ImPerm [Internet]. Available from: <https://github.com/mtorchiano/ImPerm>
66. Yang J, Lee SH, Goddard ME, Visscher PM. GCTA: a tool for genome-wide complex trait analysis. *Am J Hum Genet*. 2011;88(1):76–82. <https://doi.org/10.1016/j.ajhg.2010.11.011>.
67. Lee SH, Wray NR, Goddard ME, Visscher PM. Estimating missing heritability for disease from genome-wide association studies. *Am J Hum Genet*. 2011; 88(3):294–305. <https://doi.org/10.1016/j.ajhg.2011.02.002>.
68. Yang J, Bakshi A, Zhu Z, Hemani G, Vinkhuyzen AAE, Lee SH, et al. Genetic variance estimation with imputed variants finds negligible missing heritability for human height and body mass index. *Nat Genet*. 2015;47(10): 1114–20. <https://doi.org/10.1038/ng.3390>.
69. Evans LM, Tahmasbi R, Vrieze SJ, Abecasis GR, Das S, Gazal S, et al. Comparison of methods that use whole genome data to estimate the heritability and genetic architecture of complex traits. *Nat Genet*. 2018;50(5): 737–45. <https://doi.org/10.1038/s41588-018-0108-x>.
70. Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. Principal components analysis corrects for stratification in genome-wide association studies. *Nat Genet*. Springer Nature. 2006;38(8):904–9. <https://doi.org/10.1038/ng1847>.
71. Kang HM, Sul JH, Service SK, Zaitlen NA, Kong S-Y, Freimer NB, et al. Variance component model to account for sample structure in genome-wide association studies. *Nat Genet*. 2010;42(4):348–54. <https://doi.org/10.1038/ng.548>.
72. Kang HM. EPCATS [Internet]. Available from: <https://genome.sph.umich.edu/wiki/EPCATS>
73. Balduzzi S, Rücker G, Schwarzer G. How to perform a meta-analysis with R: a practical tutorial. *Evid Based Ment Health*. 2019;22(4):153–60. <https://doi.org/10.1136/ebmental-2019-300117>.
74. Kircher M, Witten DM, Jain P, O’Roak BJ, Cooper GM, Shendure J. A general framework for estimating the relative pathogenicity of human genetic variants. *Nat Genet*. 2014;46(3):310–5. <https://doi.org/10.1038/ng.2892>.
75. McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GRS, Thormann A, et al. The ensembl variant effect predictor. *Genome Biol*. 2016;17(1):122. <https://doi.org/10.1186/s13059-016-0974-4>.
76. Palmer TM, Lawlor DA, Harbord RM, Sheehan NA, Tobias JH, Timpson NJ, et al. Using multiple genetic variants as instrumental variables for modifiable risk factors. *Stat Methods Med Res*. 2012;21(3):223–42. <https://doi.org/10.1177/0962280210394459>.
77. van Buuren S, Groothuis-Oudshoorn K. Mice: multivariate imputation by chained equations in R. *Journal of Statistical Software, Articles*. 2011;45:1–67.
78. Honaker J, King G, Blackwell M. Amelia II: a program for missing data. *Journal of Statistical Software, Articles*. 2011;45:1–47.
79. Cole SR, Platt RW, Schisterman EF, Chu H, Westreich D, Richardson D, et al. Illustrating bias due to conditioning on a collider. *Int J Epidemiol*. 2010; 39(2):417–20. <https://doi.org/10.1093/ije/dyp334>.
80. VanderWeele TJ, Tchetgen Tchetgen EJ, Cornelis M, Kraft P. Methodological challenges in mendelian randomization. *Epidemiology*. 2014;25(3):427–35. <https://doi.org/10.1097/EDE.0000000000000081>.
81. Storey JD, Bass AJ, Dabney A, Robinson D. qvalue: Q-value estimation for false discovery rate control [Internet]. 2019. Available from: <http://github.com/jdstorey/qvalue>

Publisher’s Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.