

# Gene-level analysis of rare variants in 363,977 whole exome sequences reveals an association of *GIGYF1* loss of function with diabetes

Aimee M. Deaton<sup>1</sup>, Margaret M. Parker<sup>1</sup>, Lucas D. Ward<sup>1</sup>, Alexander O. Flynn-Carroll<sup>1</sup>, Lucas BonDurant<sup>1</sup>, Gregory Hinkle<sup>1</sup>, Paul Nioi<sup>1</sup>

<sup>1</sup>Alnylam Pharmaceuticals, Cambridge, MA 02142

## Abstract

Sequencing of large cohorts offers an unprecedented opportunity to identify rare genetic variants and to find novel contributors to human disease. We used gene-based collapsing tests to identify genes associated with glucose, HbA1c and T2D diagnosis in 363,977 exome-sequenced participants in the UK Biobank. We identified known associations with diabetes including variants in *GCK*, *HNF1A* and *PDX1*, genes involved in Mendelian forms of diabetes. Novel associations were identified for *GIGYF1* predicted loss of function (pLOF), *TNRC6B* pLOF and *PFAS* predicted damaging missense variants. Multiple rare variants contributed to these associations, including singleton variants. The most significant novel associations were seen for *GIGYF1* pLOF which associated with increased levels of glucose (0.77 mmol/L increase,  $p = 4.42 \times 10^{-12}$ ) and HbA1c (4.33 mmol/mol,  $p = 1.28 \times 10^{-14}$ ) as well as T2D diagnosis (OR = 4.15,  $p = 6.14 \times 10^{-11}$ ). *GIGYF1* pLOF also associated with decreased cholesterol levels as well as an increased risk of hypothyroidism. An independent common variant association for glucose and T2D was identified at *GIGYF1* which replicated in additional datasets. Our results highlight the role of *GIGYF1* in regulating insulin signaling and protecting from diabetes.

## Introduction

Human genetics provides powerful methods for understanding the roles of genes and proteins in disease and can lead to new therapeutic hypotheses and drug targets. Genetic evidence based on sequence variants within coding regions of the genome is better at predicting the efficacy and safety of novel therapeutics than evidence from genome-wide association studies (GWAS), which tends to involve common noncoding variants [1-3]. Among coding variants, predicted loss of function (pLOF) variants are particularly informative in association studies because they establish a direct causal link between reduction in gene function and biological outcomes. Additionally, rare missense variants predicted to be deleterious can provide valuable biological insights [4, 5]. However, interrogation of the effects of such variants is hampered by the rarity of these variants and the cohort sizes needed to identify associations [6]. Exome or whole-genome sequencing of large biobanks coupled with gene-level aggregation of rare high impact variants can help to circumvent these challenges [4]. Biobanks offer a considerable advantage over case-control cohorts as they contain richer phenotyping data which often includes biomarker measurements as well as disease diagnoses. This allows a more complete understanding of the biological consequences of damaging variants in particular genes [7, 8].

Diabetes is a disease that has been extensively studied in traditional array-based GWAS with hundreds of associations identified to date [9-12]. Although these studies have given insight into some of the

biological mechanisms contributing to diabetes, most of the reported associations are with variants in non-coding regions, making identification of the causal gene challenging. More recently, exome sequencing has been applied to discover protein-coding variants that alter the risk of developing type 2 diabetes (T2D). Sequencing of 20,791 T2D cases followed by the use of gene-based collapsing tests (to aggregate predicted damaging variants) identified associations of *SLC30A8*, *MC4R* and *PAM* with T2D diagnosis [5].

Using 363,977 whole exome sequences from the UK Biobank (UKBB) we performed gene-level collapsing tests to examine the association of pLOF and damaging missense variants in ~17,000 genes with biomarkers of glycemic control, glucose and glycated hemoglobin (HbA1c), as well as T2D diagnosis.

## Results

### Association of rare variants with glucose, HbA1c and T2D

We used 454,787 whole exome sequences from the UK Biobank (UKBB) to identify rare variants with a minor allele frequency (MAF)  $\leq 1\%$  likely to have functional impact; pLOF variants (i.e. frameshift, stop gain, splice donor or splice acceptor variants) called as high confidence by LOFTEE [13] or missense variants predicted to be damaging (Combined Annotation Dependent Depletion [CADD] score  $\geq 25$ ). We identified 726,422 rare pLOF variants affecting 16,477 genes, 58.5% of which were singletons (carried by a single individual), and 2.14 million damaging missense variants in 17,312 genes, 49.6% of which were singletons (Supplementary Table 1).

Given the large proportion of variants present in just a single individual, we used gene-based collapsing tests to look for associations with biomarkers of glycemic control and T2D diagnosis. We used two variant aggregation strategies; 1) pLOF variants with MAF  $\leq 1\%$  and 2) damaging missense variants with MAF  $\leq 1\%$  and performed burden testing in the unrelated White population ( $n=363,977$ ) adjusting for age, sex and genetic ancestry via 12 principal components.

First, we tested genes for association with glucose and HbA1c levels. We required at least 10 variant carriers per gene to have measurements based on an examination of genomic inflation at different carrier thresholds (Supplementary Figure 1). Using a p-value threshold adjusted for the number of variant sets and phenotypes tested ( $p \leq 7.82 \times 10^{-7}$ ), four genes significantly associated with glucose levels: *GCK* pLOF ( $p = 1.56 \times 10^{-9}$ , 1.24 mmol/L increase), *GCK* damaging missense ( $p = 6.15 \times 10^{-11}$ , 0.61 mmol/L increase), *GIGYF1* pLOF ( $p = 4.42 \times 10^{-12}$ , 0.77 mmol/L increase) and *G6PC2* damaging missense variants ( $p = 4.62 \times 10^{-83}$ , 0.33 mmol/L decrease) (Figure 1, Table 1). The same variant sets also associated with HbA1c levels along with 27 additional sets including *HNF1A* pLOF ( $p = 2.14 \times 10^{-7}$ , 4.01 mmol/mol increase), *TNRC6B* pLOF ( $p = 2.36 \times 10^{-7}$ , 3.94 mmol/mol increase) and *PDX1* damaging missense variants ( $p = 2.54 \times 10^{-7}$ , 0.41 mmol/mol increase) (Figure 1, Table 1).

We then tested aggregated pLOF and damaging missense variants for association with T2D diagnosis (Supplementary Figure 1). Using a p-value threshold adjusted for the number of variant sets tested ( $p \leq 1.46 \times 10^{-6}$ ), 6 variant sets significantly associated with T2D; pLOF variants in *GIGYF1*, *GCK*, *HNF1A* and *TNRC6B* and damaging missense variants in *GCK* and *PAM* (Figure 2, Table 2). As the time of available follow-up differs between England, Scotland and Wales, we controlled for country of recruitment in the regression (see Methods). In addition, we confirmed that significant hits did not associate with country of recruitment (all  $p > 0.035$ ) and that hits remained significant when only data from England were considered (Supplementary Table 2).

## Identification of genes with a biological role in diabetes

Variants in two genes, *GCK* and *GIGYF1*, significantly associated with glucose, HbA1c and T2D diagnosis, strongly suggesting a biological role in diabetes; *GCK* is involved in Mendelian forms of diabetes while *GIGYF1* has not previously been implicated by genetics in the disease. Both *GCK* and *GIGYF1* are located on chromosome 7 but are 56Mb apart, strongly suggesting that these signals are independent; this independence was confirmed by conditional analysis (Supplementary Table 3). Two additional variant sets, *HNF1A* pLOF and *TNRC6B* pLOF, had genome-wide associations with both T2D diagnosis and HbA1c levels while *G6PC2* damaging missense associated with decreased levels of both glucose and HbA1c but not T2D (Table 3).

To see which other significant genes were likely to have a role in diabetes we looked at all variant sets with a significant glucose, HbA1c, or T2D association and examined whether they had associations with additional diabetes traits using a more permissive p-value threshold correcting for the number of variant sets tested ( $p \leq 0.0016$ , 32 sets tested). Damaging missense variants in *PDX1* and *PFAS*, which had significant associations with HbA1c levels in our primary analysis, associated with T2D diagnosis using this threshold (Table 3 and Supplementary Table 4).

Many HbA1c associations appeared to be secondary to effects on red blood cells. 22 out of 31 variant sets associated with HbA1c did not show effects on glucose levels or T2D diagnosis (Supplementary Table 4) and were not implicated in Mendelian forms of diabetes. Out of these 22 variant sets, 12 were in genes implicated in Mendelian disorders affecting red blood cells (for example *EPB42* and *TFR2*; see Supplementary Table 5) and an additional five had highly significant associations with red blood cell traits in our data ( $p \leq 7.82 \times 10^{-7}$ ; Supplementary Table 6).

We focused on the variant sets associated with multiple diabetes traits as these are strong candidates for regulating glucose homeostasis. The genes fall into three main groups; known MODY (maturity-onset diabetes of the young) genes (*GCK*, *HNF1A* and *PDX1*) [14], known genes reported in previous exome-wide analyses of glucose levels or T2D (*G6PC2* and *PAM*) [5, 15], and novel genes not previously implicated by genetics in diabetes (*GIGYF1*, *TNRC6B* and *PFAS*).

Because obesity is linked to the development of T2D, we adjusted for body mass index (BMI) in the burden tests and found that the association of variants in these genes with diabetes-related traits remained significant (Supplementary Tables 7 and 8).

Associations for rare variants can be susceptible to confounders such as population stratification and sample relatedness leading to false positives. Therefore, we used the generalized linear mixed model implemented by SAIGE-Gene which accounts for relatedness and adjusts for unbalanced case-control ratios [16] to verify association of our variant sets of interest with glucose, HbA1c, and T2D diagnosis. SAIGE-Gene was run in the White population including related individuals ( $n=398,574$ ). Using the p-value thresholds previously employed, all associations were statistically significant using this method apart from the association of *TNRC6B* pLOF with HbA1c ( $p = 6.85 \times 10^{-6}$ ) and T2D diagnosis ( $p = 4.77 \times 10^{-5}$ ) (Supplementary Table 9). Therefore, additional data or replication is needed in order to determine if *TNRC6B* pLOF truly associates with HbA1c levels and T2D.

We sought to use independent measurements of glucose and HbA1c to validate the associations of interest seen in our primary analysis which used measurements taken as part of the UKBB assessment. To do this we extracted lab test values for glucose and HbA1c from primary care data, which is available

for approximately half of the cohort, taking the mean measurement per individual. In gene-based burden tests all variants showed a direction of effect consistent with that seen in the primary analysis (Supplementary Figure 2). Nine out of the 12 associations with glucose and HbA1c were significant when correcting for the number of tests performed ( $p \leq 0.004$ ), including the association of *GIGYF1* pLOF with glucose ( $p=1.21 \times 10^{-5}$ , effect = 0.63 SD) and HbA1c ( $p=1.38 \times 10^{-5}$ , effect = 0.73 SD) levels (Supplementary Table 10).

## Multiple variants contribute to associations with diabetes diagnosis and biomarkers

To examine whether specific variants were driving the associations with diabetes traits we conducted “leave-one-out” burden tests. The association of *PAM* missense variants with T2D diagnosis was driven entirely by a previously reported variant Ser539Trp (rs78408340;  $p = 0.43$  when Ser539Trp is excluded). For all other variant sets, multiple variants contributed to the associations observed (Supplementary Figure 3). Notably, when singleton variants were excluded, half of the associations no longer reached significance including those for *GCK* pLOF and glucose ( $p = 0.0015$  without singletons versus  $p = 1.56 \times 10^{-9}$ ) and *GIGYF1* pLOF and T2D ( $p = 2.9 \times 10^{-5}$  without singletons versus  $p = 6.14 \times 10^{-11}$ ) (Supplementary Table 11), demonstrating the power of including singletons in gene-based tests.

For the variants contributing to our novel discovered associations, *GIGYF1* pLOF, *TNRC6B* pLOF and *PFAS* damaging missense variants, we examined the quality scores, sequencing depth, transcripts affected and presence of contributing variants in gnomAD. We found that for *GIGYF1* and *PFAS* the variants contributing most to the associations had good quality scores and depth and were present in the non-Finnish European population in gnomAD. In contrast, *TNRC6B* is a highly constrained gene and the most common pLOF variant is not present in gnomAD. This further underscores the need for replication of the *TNRC6B* associations as pLOF variants for constrained genes may not result in true loss of function (see Supplementary Note and Supplementary Figure 4).

## PheWAS of *GIGYF1* pLOF reveals associations with cholesterol levels, hypothyroidism and complications of diabetes

The most significant novel associations were seen for *GIGYF1* pLOF which associates with increased glucose, HbA1c and increased incidence of T2D diagnosis. *GIGYF1* encodes a protein named for its binding to GRB10 (GRB10 interacting GYF protein 1), an adapter protein that has been shown to bind both the insulin and IGF-1 receptors. The association between *GIGYF1* pLOF and increased diabetes risk indicates that *GIGYF1* has a role in regulating insulin signaling and in protecting from diabetes. To give additional insight into the biological roles of *GIGYF1* we performed a phenome-wide association study (PheWAS) testing *GIGYF1* pLOF for association with 142 quantitative traits and 262 ICD10-coded diagnoses. Based on the number of tests performed, the threshold for significance was  $p \leq 1.22 \times 10^{-4}$  (Figure 4).

*GIGYF1* pLOF strongly associated with decreased levels of total cholesterol ( $p=2.44 \times 10^{-12}$ , effect = -0.61 SD) which was, in large part, driven by LDL cholesterol ( $p = 2.40 \times 10^{-10}$ , effect = -0.56 SD) although an effect on HDL cholesterol was also observed (Table 4). We also detected associations with decreased grip strength and decreased peak expiratory flow which may reflect changes in body size, muscle mass or general health in carriers [17, 18]. Notably, *GIGYF1* pLOF also associated with increased levels of the kidney injury biomarker cystatin c ( $p= 6.65 \times 10^{-6}$ , effect = 0.36 SD) and increased diagnosis of urinary

system disorders ( $p = 7.32 \times 10^{-5}$ , OR = 2.71) which might suggest renal complications of diabetes in carriers (Table 4 and Table 5).

After diabetes, the next most significant disease association of *GIGYF1* pLOF was with increased risk of hypothyroidism ( $p = 1.25 \times 10^{-9}$ , OR = 4.53). 21 out of the 131 *GIGYF1* pLOF carriers had a diagnosis of unspecified hypothyroidism and 7 of these also had a diagnosis of T2D. We also saw a significant association with increased risk of syncope and collapse ( $p = 1.92 \times 10^{-6}$ , OR = 3.75), possibly reflecting complications of diabetes or thyroid disorders (Table 5).

Other phenome-wide significant associations with quantitative traits included waist circumference, total protein and mean corpuscular hemoglobin as well increased time to complete a cognitive test (Table 4). To ensure that the association of *GIGYF1* pLOF with HbA1c was independent of effects on hemoglobin we adjusted for mean corpuscular hemoglobin level and verified that the association remained highly significant ( $p = 4.10 \times 10^{-12}$ ). *GIGYF1* pLOF also associated with increased diagnosis of emphysema and anemia (Table 5).

We also performed PheWAS on *TNRC6B* pLOF and *PFAS* damaging missense variants. PheWAS of *TNRC6B* pLOF revealed associations with decreased grip strength, increased blood urea levels as well as decreased measures of lung function, forced vital capacity and forced expiratory volume (Supplementary Table 12). We also saw associations with increased time to complete a cognitive test and increased speech recognition threshold (SRT) suggesting hearing impairment in carriers. This was further supported by an association with increased diagnosis of hearing loss in *TNRC6B* pLOF carriers ( $p = 1.63 \times 10^{-4}$ , OR = 5.18). *PFAS* damaging missense variants did not have any significant associations beyond those with HbA1c levels and T2D.

### Common variants at *GIGYF1* associate with glucose, T2D and *GIGYF1* expression

As *GIGYF1* pLOF variants are very rare it is extremely challenging to replicate the associations seen in UKBB. Therefore, we looked for more common variants that could further shed light on the role of the locus and gene in diabetes. We tested array genotyped and imputed variants at the *GIGYF1* locus for association with glucose levels in 294,042 unrelated White individuals with measurements available. We found a cluster of variants in a linkage disequilibrium block covering *GIGYF1* and *EPO* significantly associating with glucose levels (Figure 4). This signal is represented by rs221783, an intergenic variant whose minor T allele associated with decreased glucose ( $p = 1.8 \times 10^{-11}$ , effect = -0.03 SD,) and HbA1c ( $p = 3.6 \times 10^{-7}$ , effect = -0.02 SD,) levels as well as increased cholesterol ( $p = 7.0 \times 10^{-12}$ , effect = 0.03 SD,). This variant also associated with a decreased risk of T2D ( $p = 0.005$ , OR = 0.96) and hypothyroidism ( $p = 6.95 \times 10^{-7}$ , OR=0.92) (Table 6). rs221783 is the best eQTL ( $R^2 > 0.8$ ) for *GIGYF1* in several tissues including pancreas, adipose and thyroid [19] (Supplementary Table 13). In all tissues, the T allele associating with decreased glucose and decreased T2D risk associated with increased *GIGYF1* expression. Conditional analysis showed that the glucose and HbA1c associations of *GIGYF1* pLOF and rs221783 are independent of each other (Supplementary Table 14).

The association of rs221783 with glucose levels replicated in Biobank Japan ( $p = 1.7 \times 10^{-4}$ , effect = -0.05 SD for T allele) [20] whilst in FinnGen, rs221783 showed a nominal association with T2D risk ( $p = 0.02$ , OR = 0.96 for T allele) (Supplementary Table 15). The association with thyroid disease has been replicated elsewhere [21].

The independent glucose and T2D associations at the *GIGYF1* locus and their replication in other datasets further support the hypothesis that decreasing *GIGYF1* predisposes to diabetes while increasing *GIGYF1* levels may protect from diabetes.

## Discussion

Our results highlight the power of whole exome sequencing to make novel discoveries relevant to human disease and to detect known associations of Mendelian disease genes. Gene-level aggregation and burden testing of rare pLOF and predicted damaging missense variants identified genes associating with diabetes and biomarkers of glycemic control. These included several genes not previously implicated in diabetes, *GIGYF1*, *TNRC6B* and *PFAS*, as well as *GCK*, *HNF1A* and *PDX1*, known MODY genes [14, 22-24]. We also identified *PAM* and *G6PC2*, genes identified by other rare-variant studies of T2D and glucose levels [5, 15]. Gene-level tests were needed to detect the majority of these associations owing to the rarity of damaging variants. For example, out of 363,977 individuals, just 40 carried a pLOF variant in *GCK* and 131 carried a pLOF variant in *GIGYF1*. In general, singleton variants contributed a large part of the signal arguing strongly, as others have done [4], for including such variants in gene-based collapsing tests.

Test statistic inflation can be a challenge when testing rare variants as statistical assumptions break down when the number of carriers expected to have the disease of interest is low [4, 25]. To avoid false positives in our analysis of diabetes, we initially examined associations with glucose and HbA1c because quantitative traits are less susceptible to inflation. All of the variant sets that associated with T2D also affected HbA1c and/or glucose levels giving us confidence in these associations. In addition, T2D associations for all genes, apart from *TNRC6B*, were significant ( $p \leq 1.46 \times 10^{-6}$ ) using the linear mixed model implemented by SAIGE-Gene which can be more robust when dealing with low numbers of variant carriers [16]. Additional confidence in our results comes from the fact that we identified genes known to be involved in Mendelian forms of diabetes and previously reported genes.

We uncovered novel associations with T2D and biomarkers of glycemic control for aggregated variants in *GIGYF1*, *TNRC6B* and *PFAS*. *TNRC6B* encodes trinucleotide repeat-containing gene 6B protein which is involved siRNA- and miRNA-mediated gene silencing [26-28]. As *TNRC6B* is highly constrained, one must view pLOF variants in this gene with suspicion as they may be sequencing errors or not result in a true loss of function. Additional data will be needed to verify the association of *TNRC6B* pLOF with T2D. Interestingly, heterozygous protein-truncating variants in *TNRC6B* have been implicated in developmental delay and autism spectrum disorders [29]. Consistent with this, we see an association of *TNRC6B* pLOF with worse performance in a cognitive test as well as with hearing impairment. Damaging missense variants in *PFAS* associate with decreased HbA1c levels and decreased incidence of T2D diagnosis with a large part of the signal driven by two missense variants Glu434Lys and Ala466Thr. *PFAS* encodes phosphoribosylformylglycinamide synthase which is involved in de novo synthesis of the purine inosine monophosphate (IMP), a process which is required by growing and proliferating cells [30]. A mechanistic link to glucose homeostasis or diabetes is not apparent but *PFAS* may have roles beyond purine biosynthesis as it interacts with and deamidates other proteins [31]. *PFAS* is ubiquitously expressed [19, 32].

We focused our analysis on understanding the consequences of *GIGYF1* pLOF as it strongly associated with glucose, HbA1c and T2D. *GIGYF1* encodes a protein that was initially identified for its binding to the

adapter protein GRB10 which negatively regulates both the insulin and IGF-1 receptors [33]. Transfection of cells with GRB10-binding fragments of GIGYF1 lead to greater activation of both the insulin and IGF-1 receptors [34]. This supports a hypothesis whereby GIGYF1 enhances insulin signaling by reducing the negative regulation of the insulin receptor by GRB10. When GIGYF1 is reduced, as is the case in individuals carrying pLOF variants, GRB10 presumably inhibits insulin signaling to a greater degree thereby reducing the action of insulin in its target tissues and leading to increased risk of T2D. However, the exact mechanistic details of these interactions remain to be determined. *GRB10* variants have also been reported to associate with glycemic traits although interpretation of these results is complicated by imprinting [35]. *GIGYF1* is broadly expressed with high levels observed in endocrine tissues, pancreas and brain [19, 32]. GIGYF1 and the related protein GIGYF2 have also been implicated in translational repression [36] and translation-coupled mRNA decay [37] suggesting biological roles beyond regulation of insulin and IGF-1 receptor signaling.

PheWAS of *GIGYF1* pLOF revealed a strong association with decreased cholesterol levels reflecting altered energy homeostasis in carriers. An inverse relationship between glucose and cholesterol levels has been observed for variants in other genes [38]. We also observed several associations that could reflect complications of diabetes in *GIGYF1* pLOF carriers including increased cystatin c levels and increased diagnosis of urinary disorders suggesting renal complications as well as syncope and collapse which may be a side-effect of hyperglycemia and/or hypoglycemia in diabetics. Other associations may reflect poor health in carriers including decreased grip strength and decreased peak expiratory flow. *GIGYF1* pLOF also associated with decreased mean corpuscular hemoglobin levels and increased diagnosis of anemia as well as increased emphysema diagnosis. The biological basis for these associations is not clear. *GIGYF1* is highly expressed in lung [19, 32] although the emphysema association is driven by small numbers of individuals, so replication is required.

*GIGYF1* pLOF associated with a 4.5-fold increased risk of hypothyroidism and *GIGYF1* is highly expressed in thyroid [19, 32] consistent with a biological function in this tissue. IGF-1 and insulin have been implicated in the proliferation of thyroid cells which may, in part, explain the association with thyroid dysfunction [39-41]. An alternative explanation is that GIGYF1 contributes to thyroid function by affecting secretion of thyroid stimulating hormone in the anterior pituitary gland. Damaging variants in *GIGYF1* have recently been implicated in conferring risk for developmental delay and autism spectrum disorders [42]. Consistent with this, we see an association of *GIGYF1* pLOF with increased time to complete a cognitive test. It may be that metabolic aberrations in carriers affect cognitive performance, that brain development is altered due to perturbation of IGF-1 signaling, or that other functions of GIGYF1 such as regulation of mRNA expression and decay are responsible for cognitive phenotypes.

Replication of associations is a challenge in rare variant analysis. The UKBB dataset used consists of exomes and baseline biomarker measurements for 363,977 individuals which is, to our knowledge, the largest dataset of its kind available. We replicated the majority of our associations with glucose and HbA1c levels, including those for *GIGYF1* pLOF, using independent measurements from primary care data. Additionally, we turned to common genetic variants to further investigate the role of the *GIGYF1* locus in diabetes. Non-coding variants at the *GIGYF1* locus associated with glucose levels and T2D, and we replicated these findings in independent datasets. These variants associated with increased *GIGYF1* expression but a lower risk of T2D. This direction of effect is consistent with what we see for the pLOF variants – reduced levels of *GIGYF1* increases diabetes risk but increased levels of *GIGYF1* are protective.

We assessed the impact of pLOF and damaging missense variants in approximately 17,000 genes on glycemic traits and uncovered a hitherto unappreciated role for *GIGYF1* in regulating blood sugar and protecting from T2D. By highlighting the importance of GIGYF1 and GRB adapter proteins in modulating insulin signaling this finding may lead to new therapeutic approaches for the treatment of diabetes. Discoveries such as this are only possible by combining health-related data with the sequencing of rare variants on a biobank scale.

## Methods

### The UK Biobank resource and data access

The UK Biobank (UKBB) recruited ~500,000 participants in England, Wales, and Scotland between 2006 and 2010 [43]. Written informed consent was obtained from all participants. Phenotypic data available includes age, sex, biomarker data and self-reported diseases collected at the time of baseline assessment as well as disease diagnoses from inpatient hospital stays, the cancer registry and death records obtained through the NHS. Approximately half of the participants also have diagnoses from primary care available. Array genotypes are available for nearly all participants and exome sequencing data is available for 454,787 participants. The data used in this study were obtained from the UKBB through application 26041. All phenotypic data and array genotypes are accessible through application to UKBB. Currently, exome sequencing data for ~200,000 participants is available [44]; the remainder of the exome data is scheduled for public release in 2021.

### Population definition and PC calculation for subjects with exome data

Subject quality control was performed by Regeneron Genetics Center (RGC) and removed subjects with evidence of contamination, unresolved duplications, sex discrepancies and discordance between exome sequencing and genotyping data. Genetic relationships between participants were determined by RGC using the PRIMUS program [45]. For the unrelated subset all first- and second-degrees relatives and some third-degree relatives were excluded.

Populations were defined through a combination of self-reported ethnicity and genetic principal components. We selected the unrelated individuals who identify as White (Field 21000) and ran an initial principal component analysis (PCA) on high quality common variants using eigenstrat [46]. SNPs were filtered for missingness across individuals < 2%, MAF > 1%, regions of known long range LD [47], and pruned to independent markers with pairwise LD < 0.1. We then projected the principal components (PCs) onto related individuals and removed all individuals +/- 3 standard deviations from the mean of PCs 1-6. A final PC estimation was performed in eigenstrat [46] using unrelated subjects. We then projected related individuals onto the PCs.

### Exome sequencing and variant calling

DNA was extracted from whole blood and was sequenced by the RGC as described elsewhere [48]. Briefly, the xGen exome capture was used and reads were sequenced using the Illumina NovaSeq 6000 platform. Reads were aligned to the GRCh38 reference genome using BWA-mem [49]. Duplicate reads were identified and excluded using the Picard MarkDuplicates tool (Broad Institute). Variant calling of SNVs and indels was done using the WeCall variant caller (Genomics Plc.) to produce a GVCF for each subject. GVCFs were combined to using the GLnexus joint calling tool [50]. Post-variant calling filtering was applied using the Goldilocks pipeline [48]. Variants were annotated using the Ensembl Variant Effect Predictor v95 [51] which includes a LOFTEE plug-in to identify high confidence (HC) pLOF variants [13].



Combined Annotation Dependent Depletion (CADD) scores were generated using the Whole Genome Sequence Annotator (WGSA) AML version 0.8.

## Phenotype definitions

Blood biochemistry values were obtained for glucose (Field 30740) and HbA1c (Field 30750) and inverse rank normalized using the RNOmni R package [52], resulting in an approximately normal distribution.

For disease diagnoses, ICD10 codes were obtained from inpatient hospital diagnoses (Field 41270), causes of death (Field 40001 and 40002) and the cancer registry (Field 40006). Diagnoses also included additional hospital episode statistics (HESIN) and death registry data made available by UKBB in July 2020. T2D was defined as ICD10 E11.

For phenome-wide analyses, a selection of quantitative traits was obtained from other fields, encompassing anthropometric measurements, blood counts, as well as blood and urine biochemistry. Beyond these measurements, we selected additional quantitative traits found to be heritable ( $h^2$  significance flagged as at least “nominal” with a confidence level flagged as “medium” or “high”) by the Neale lab [25], using PHESANT to transform values to quantitative traits when necessary as they describe. These included the results of cognitive tests. All quantitative traits were inverse rank normalized using the RNOmni R package. [52]. For burden testing, we required at least 10 carriers to have measurements. We also tested associations with ICD10-coded diagnoses (using 3 character codes) that had more than 500 cases in the White subset of participants with exome data and at least one expected case carrier based on variant frequency and disease prevalence.

Glucose and HbA1c values were also extracted from primary care data available for about half of the cohort using the following read codes. Glucose: read 2 codes 44U., 44g., 44g1., 44TJ., 44f., 44TK., 44f1., 44g0., 44f0. and read 3 codes XM0ly, X772z, XE2mq; HbA1c: read 2 codes 42W5., 44TB., 66Ae0, 44TC., 42W4. and read 3 codes XaPbt, X772q, XaWP9, XaBLm, XaERp.

Values were converted to IFCC units where necessary. Aberrantly high (> 50 mmol/L for glucose, > 200 mmol/mol for HbA1c) and values of zero were excluded. The mean measurement per individual was then taken and inverse rank normalized prior to association testing.

## Gene-based association testing

For gene-based tests, autosomal rare pLOF variants were identified as follows; LOFTEE high confidence LOFs, MAF  $\leq$  1%, missingness across individuals  $\leq$  2%, HWE p-value  $\geq 10^{-10}$ . Predicted damaging missense variants were defined as missense variants with a CADD PHRED-scaled score  $\geq$  25, MAF  $\leq$  1%, missingness across individuals  $\leq$  2%, HWE p-value  $\geq 10^{-10}$ . Only genes with more than one pLOF variant or damaging missense variant were tested.

Burden testing was performed unrelated White subset using glm in R, using a gaussian model for quantitative traits and a binomial model for case-control analyses. Genotype was coded as 0 (no variant) or 1 (any number of variants). We adjusted for age, sex and the first 12 PCs of genetic ancestry in the regression. Additionally, when testing for association with disease diagnoses, we included country of recruitment as a covariate as the time of available follow-up differs between England, Scotland and Wales. Recruitment country was defined using the location of the relevant UKBB recruitment center (Field 54). Associations were later confirmed using just participants recruited in England. For case-control analyses we only ran tests where there was at least one expected case carrier based on variant

frequency and disease prevalence. For quantitative traits we required at least 10 carriers to have measurements.

For glucose and HbA1c, to convert effect sizes from normalized values back to measured units, the estimates from the regression were multiplied by the standard deviation of these traits in the entire cohort.

SAIGE-Gene was run using the SAIGE R package (v0.36.5) [53] using settings recommend by the developers and related individuals were included.

Manhattan plots were created using the R Package CMplot (<https://github.com/YinLiLin/R-CMplot>).

## Array association testing

Genotypes were obtained through array typing and imputation as described previously [54]. Population definition and PC estimation for individuals with array data was performed as previously described [55]. We tested all variants with imputation quality score (info)  $\geq 0.8$  and minor allele frequency (MAF)  $\geq 0.1\%$  in a 200Mb region around *GIGYF1* for association with glucose, HbA1c, T2D and hypothyroidism. Association analyses were performed using an additive model in PLINK adjusting for age at recruitment to UKBB, sex and the first 12 PCs of genetic ancestry. We also adjusted for country of recruitment where appropriate. The most significant variant with info  $> 0.95$  was selected as the lead variant at the locus.

We replicated the association of rs221783 with glucose using available summary statistics for Biobank Japan for the trait “blood sugar” (<http://jenger.riken.jp/en/result>) [20]. We replicated the association of this variant with T2D diagnosis using summary statistics from FinnGen release 3 for the phenotype “E4\_DM2” ([https://www.finnngen.fi/en/access\\_results](https://www.finnngen.fi/en/access_results)). The effect allele in these datasets was the alternate allele “C”. For consistency with the UKBB associations we have shown the effect for the “T” allele.

Meta-analysis of the UKBB and replication dataset association results was performed with the METAL software package using the classical method [56].

Region plots were created using LocusZoom [57]. LD calculations were performed in the White population for array variants in a 500kb sliding window as follows; we extracted genotypes with info  $> 0.9$ , rounded them to whole numbers, mean-imputed missing genotypes and used the R “cor” function to compute R which was then squared to get an  $R^2$  value.

## Gene expression and eQTL analysis

The expression of *GIGYF1* in various tissues was assessed using the GTEx portal (accessed 08/04/2020) [19] and Human Protein Atlas (<http://www.proteinatlas.org>) [32]. eQTL data for rs221783 was obtained from GTEx v8. For each tissue of interest, the best eQTL for *GIGYF1* was identified (GTEx v8 “eGene”).  $R^2$  for rs221783 and the best *GIGYF1* eQTL was calculated as described above.

## Author Contributions

A.D, L.W., M.P., A.F.C. and P.N. performed computational analyses; A.D. wrote the manuscript; all authors interpreted results and edited the manuscript.

## Competing Interests

The authors are all employees of Alnylam Pharmaceuticals.

## Acknowledgements

We would like to thank the participants and researchers of UK Biobank for creating an open-access resource. We thank Regeneron Genetics Center, the UK Biobank Exome Sequencing Consortium and UK Biobank for facilitating exome sequencing of participants. We also thank the participants and investigators of the FinnGen study and Biobank Japan. Data management and analytics were performed using the REVEAL/SciDB translational analytics platform from Paradigm4.

## References

1. King EA, Davis JW, Degner JF. Are drug targets with genetic support twice as likely to be approved? Revised estimates of the impact of genetic support for drug mechanisms on the probability of drug approval. *PLOS Genetics*. 2019;15(12):e1008489. doi: 10.1371/journal.pgen.1008489.
2. Nelson MR, Tipney H, Painter JL, Shen J, Nicoletti P, Shen Y, et al. The support of human genetic evidence for approved drug indications. *Nature Genetics*. 2015;47(8):856-60. doi: 10.1038/ng.3314.
3. Nguyen PA, Born DA, Deaton AM, Nioi P, Ward LD. Phenotypes associated with genes encoding drug targets are predictive of clinical trial side effects. *Nature Communications*. 2019;10(1):1579. doi: 10.1038/s41467-019-09407-3.
4. Cirulli ET, White S, Read RW, Elhanan G, Metcalf WJ, Tanudjaja F, et al. Genome-wide rare variant analysis for thousands of phenotypes in over 70,000 exomes from two cohorts. *Nat Commun*. 2020;11(1):542. Epub 2020/01/30. doi: 10.1038/s41467-020-14288-y. PubMed PMID: 31992710; PubMed Central PMCID: PMC6987107.
5. Flannick J, Mercader JM, Fuchsberger C, Udler MS, Mahajan A, Wessel J, et al. Exome sequencing of 20,791 cases of type 2 diabetes and 24,440 controls. *Nature*. 2019;570(7759):71-6. Epub 2019/05/24. doi: 10.1038/s41586-019-1231-2. PubMed PMID: 31118516; PubMed Central PMCID: PMC6699738.
6. Moutsianas L, Agarwala V, Fuchsberger C, Flannick J, Rivas MA, Gaulton KJ, et al. The Power of Gene-Based Rare Variant Methods to Detect Disease-Associated Variation and Test Hypotheses About Complex Disease. *PLOS Genetics*. 2015;11(4):e1005165. doi: 10.1371/journal.pgen.1005165.
7. Bush WS, Oetjens MT, Crawford DC. Unravelling the human genome-phenome relationship using phenome-wide association studies. *Nat Rev Genet*. 2016;17(3):129-45. Epub 2016/02/16. doi: 10.1038/nrg.2015.36. PubMed PMID: 26875678.
8. Diogo D, Tian C, Franklin CS, Alanne-Kinnunen M, March M, Spencer CCA, et al. Phenome-wide association studies across large population cohorts support drug target validation. *Nat Commun*. 2018;9(1):4285. Epub 2018/10/18. doi: 10.1038/s41467-018-06540-3. PubMed PMID: 30327483; PubMed Central PMCID: PMC6191429.
9. Vujkovic M, Keaton JM, Lynch JA, Miller DR, Zhou J, Tcheandjieu C, et al. Discovery of 318 new risk loci for type 2 diabetes and related vascular outcomes among 1.4 million participants in a multi-ancestry meta-analysis. *Nat Genet*. 2020;52(7):680-91. Epub 2020/06/17. doi: 10.1038/s41588-020-0637-y. PubMed PMID: 32541925; PubMed Central PMCID: PMC67343592.
10. Xue A, Wu Y, Zhu Z, Zhang F, Kemper KE, Zheng Z, et al. Genome-wide association analyses identify 143 risk variants and putative regulatory mechanisms for type 2 diabetes. *Nature Communications*. 2018;9(1):2941. doi: 10.1038/s41467-018-04951-w.

11. Huang J, Ellinghaus D, Franke A, Howie B, Li Y. 1000 Genomes-based imputation identifies novel and refined associations for the Wellcome Trust Case Control Consortium phase 1 Data. *Eur J Hum Genet.* 2012;20(7):801-5. Epub 2012/02/02. doi: 10.1038/ejhg.2012.3. PubMed PMID: 22293688; PubMed Central PMCID: PMC3376268.
12. Scott RA, Scott LJ, Magi R, Marullo L, Gaulton KJ, Kaakinen M, et al. An Expanded Genome-Wide Association Study of Type 2 Diabetes in Europeans. *Diabetes.* 2017;66(11):2888-902. Epub 2017/06/02. doi: 10.2337/db16-1253. PubMed PMID: 28566273; PubMed Central PMCID: PMC5652602.
13. Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alfoldi J, Wang Q, et al. The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature.* 2020;581(7809):434-43. Epub 2020/05/29. doi: 10.1038/s41586-020-2308-7. PubMed PMID: 32461654; PubMed Central PMCID: PMC7334197.
14. Fajans SS, Bell GI, Polonsky KS. Molecular mechanisms and clinical pathophysiology of maturity-onset diabetes of the young. *N Engl J Med.* 2001;345(13):971-80. Epub 2001/09/29. doi: 10.1056/NEJMra002168. PubMed PMID: 11575290.
15. Wessel J, Chu AY, Willems SM, Wang S, Yaghootkar H, Brody JA, et al. Low-frequency and rare exome chip variants associate with fasting glucose and type 2 diabetes susceptibility. *Nature Communications.* 2015;6(1):5897. doi: 10.1038/ncomms6897.
16. Zhou W, Zhao Z, Nielsen JB, Fritsche LG, LeFaive J, Gagliano Taliun SA, et al. Scalable generalized linear mixed model for region-based association tests in large biobanks and cohorts. *Nature Genetics.* 2020;52(6):634-9. doi: 10.1038/s41588-020-0621-6.
17. Tikkanen E, Gustafsson S, Amar D, Shcherbina A, Waggott D, Ashley EA, et al. Biological Insights Into Muscular Strength: Genetic Findings in the UK Biobank. *Scientific Reports.* 2018;8(1):6451. doi: 10.1038/s41598-018-24735-y.
18. Willems SM, Wright DJ, Day FR, Trajanoska K, Joshi PK, Morris JA, et al. Large-scale GWAS identifies multiple loci for hand grip strength providing biological insights into muscular fitness. *Nature Communications.* 2017;8(1):16015. doi: 10.1038/ncomms16015.
19. Aguet F, Brown AA, Castel SE, Davis JR, He Y, Jo B, et al. Genetic effects on gene expression across human tissues. *Nature.* 2017;550(7675):204-13. doi: 10.1038/nature24277.
20. Kanai M, Akiyama M, Takahashi A, Matoba N, Momozawa Y, Ikeda M, et al. Genetic analysis of quantitative traits in the Japanese population links cell types to complex human diseases. *Nat Genet.* 2018;50(3):390-400. Epub 2018/02/07. doi: 10.1038/s41588-018-0047-6. PubMed PMID: 29403010.
21. Saevarsdottir S, Olafsdottir TA, Ivarsdottir EV, Halldorsson GH, Gunnarsdottir K, Sigurdsson A, et al. FLT3 stop mutation increases FLT3 ligand level and risk of autoimmune thyroid disease. *Nature.* 2020;584(7822):619-23. doi: 10.1038/s41586-020-2436-0.
22. Froguel P, Vaxillaire M, Sun F, Velho G, Zouali H, Butel MO, et al. Close linkage of glucokinase locus on chromosome 7p to early-onset non-insulin-dependent diabetes mellitus. *Nature.* 1992;356(6365):162-4. Epub 1992/03/12. doi: 10.1038/356162a0. PubMed PMID: 1545870.
23. Ellard S. Hepatocyte nuclear factor 1 alpha (HNF-1 alpha) mutations in maturity-onset diabetes of the young. *Hum Mutat.* 2000;16(5):377-85. Epub 2000/11/03. doi: 10.1002/1098-1004(200011)16:5<377::AID-HUMU1>3.0.CO;2-2. PubMed PMID: 11058894.
24. Stoffers DA, Ferrer J, Clarke WL, Habener JF. Early-onset type-II diabetes mellitus (MODY4) linked to IPF1. *Nat Genet.* 1997;17(2):138-9. Epub 1997/11/05. doi: 10.1038/ng1097-138. PubMed PMID: 9326926.
25. Churchhouse C. Details and Considerations of the UK Biobank GWAS. Neale lab. 2017; <http://www.nealelab.is/blog/2017/9/11/details-and-considerations-of-theuk-biobank-gwas>.
26. Rasch F, Weber R, Izaurrealde E, Igreja C. 4E-T-bound mRNAs are stored in a silenced and deadenylated form. *Genes Dev.* 2020;34(11-12):847-60. Epub 2020/05/02. doi: 10.1101/gad.336073.119. PubMed PMID: 32354837; PubMed Central PMCID: PMC7263148.

27. Meister G, Landthaler M, Peters L, Chen PY, Urlaub H, Luhrmann R, et al. Identification of novel argonaute-associated proteins. *Curr Biol*. 2005;15(23):2149-55. Epub 2005/11/18. doi: 10.1016/j.cub.2005.10.048. PubMed PMID: 16289642.
28. Weinmann L, Hock J, Ivacevic T, Ohrt T, Mutze J, Schwille P, et al. Importin 8 is a gene silencing factor that targets argonaute proteins to distinct mRNAs. *Cell*. 2009;136(3):496-507. Epub 2009/01/27. doi: 10.1016/j.cell.2008.12.023. PubMed PMID: 19167051.
29. Granadillo JL, A PAS, Guo H, Xia K, Angle B, Bontempo K, et al. Pathogenic variants in TNRC6B cause a genetic disorder characterised by developmental delay/intellectual disability and a spectrum of neurobehavioural phenotypes including autism and ADHD. *J Med Genet*. 2020;57(10):717-24. Epub 2020/03/11. doi: 10.1136/jmedgenet-2019-106470. PubMed PMID: 32152250.
30. Ali ES, Sahu U, Villa E, O'Hara BP, Gao P, Beaudet C, et al. ERK2 Phosphorylates PFAS to Mediate Posttranslational Control of De Novo Purine Synthesis. *Mol Cell*. 2020;78(6):1178-91 e6. Epub 2020/06/03. doi: 10.1016/j.molcel.2020.05.001. PubMed PMID: 32485148; PubMed Central PMCID: PMC7306006.
31. Lu A, Disoma C, Zhou Y, Chen Z, Zhang L, Shen Y, et al. Protein interactome of the deamidase phosphoribosylformylglycinamide synthetase (PFAS) by LC-MS/MS. *Biochem Biophys Res Commun*. 2019;513(3):746-52. Epub 2019/04/17. doi: 10.1016/j.bbrc.2019.04.039. PubMed PMID: 30987822.
32. Uhlen M, Fagerberg L, Hallstrom BM, Lindskog C, Oksvold P, Mardinoglu A, et al. Proteomics. Tissue-based map of the human proteome. *Science*. 2015;347(6220):1260419. Epub 2015/01/24. doi: 10.1126/science.1260419. PubMed PMID: 25613900.
33. Holt LJ, Siddle K. Grb10 and Grb14: enigmatic regulators of insulin action--and more? *Biochem J*. 2005;388(Pt 2):393-406. Epub 2005/05/20. doi: 10.1042/BJ20050216. PubMed PMID: 15901248; PubMed Central PMCID: PMC1138946.
34. Giovannone B, Lee E, Laviola L, Giorgino F, Cleveland KA, Smith RJ. Two novel proteins that are linked to insulin-like growth factor (IGF-I) receptors by the Grb10 adapter and modulate IGF-I signaling. *J Biol Chem*. 2003;278(34):31564-73. Epub 2003/05/29. doi: 10.1074/jbc.M211572200. PubMed PMID: 12771153.
35. Prokopenko I, Poon W, Magi R, Prasad BR, Salehi SA, Almgren P, et al. A central role for GRB10 in regulation of islet function in man. *PLoS Genet*. 2014;10(4):e1004235. Epub 2014/04/05. doi: 10.1371/journal.pgen.1004235. PubMed PMID: 24699409; PubMed Central PMCID: PMC3974640.
36. Peter D, Weber R, Sandmeir F, Wohlbold L, Helms S, Bawankar P, et al. GIGYF1/2 proteins use auxiliary sequences to selectively bind to 4EHP and repress target mRNA expression. *Genes Dev*. 2017;31(11):1147-61. Epub 2017/07/13. doi: 10.1101/gad.299420.117. PubMed PMID: 28698298; PubMed Central PMCID: PMC5538437.
37. Weber R, Chung MY, Keskeny C, Zinnall U, Landthaler M, Valkov E, et al. 4EHP and GIGYF1/2 Mediate Translation-Coupled Messenger RNA Decay. *Cell Rep*. 2020;33(2):108262. Epub 2020/10/15. doi: 10.1016/j.celrep.2020.108262. PubMed PMID: 33053355.
38. Klimentidis YC, Arora A, Newell M, Zhou J, Ordovas JM, Renquist BJ, et al. Type-2 diabetes with low LDL-C: genetic insights into a unique phenotype. *bioRxiv*. 2019;10.1101/837013 %J bioRxiv:837013. doi: 10.1101/837013 %J bioRxiv.
39. Clement S, Refetoff S, Robaye B, Dumont JE, Schurmans S. Low TSH requirement and goiter in transgenic mice overexpressing IGF-I and IGF-Ir receptor in the thyroid gland. *Endocrinology*. 2001;142(12):5131-9. Epub 2001/11/20. doi: 10.1210/endo.142.12.8534. PubMed PMID: 11713206.
40. Kimura T, Van Keymeulen A, Golstein J, Fusco A, Dumont JE, Roger PP. Regulation of thyroid cell proliferation by TSH and other factors: a critical evaluation of in vitro models. *Endocr Rev*. 2001;22(5):631-56. Epub 2001/10/06. doi: 10.1210/edrv.22.5.0444. PubMed PMID: 11588145.

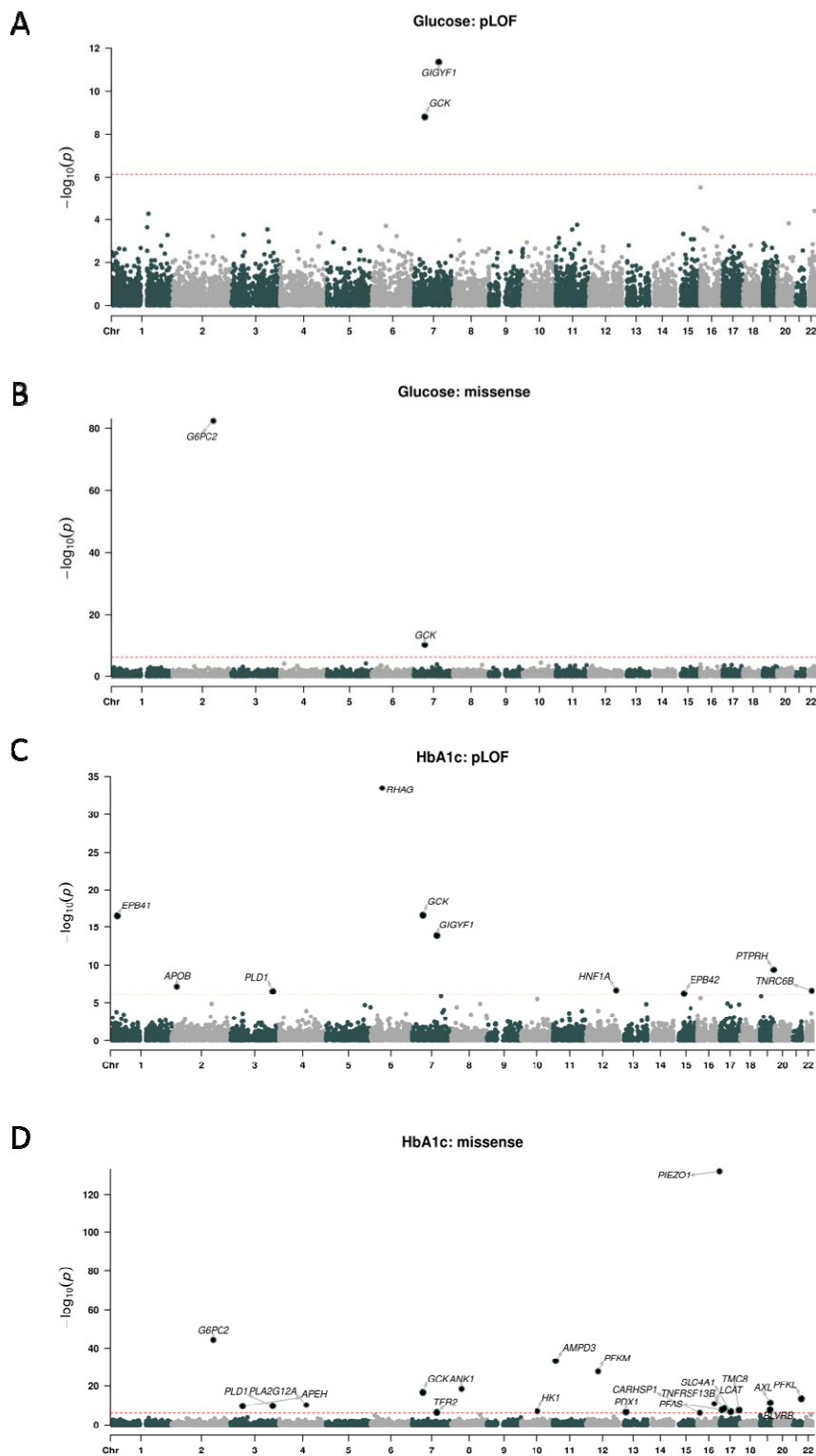
41. Zaballos MA, Santisteban P. FOXO1 controls thyroid cell proliferation in response to TSH and IGF-I and is involved in thyroid tumorigenesis. *Mol Endocrinol.* 2013;27(1):50-62. Epub 2012/11/20. doi: 10.1210/me.2012-1032. PubMed PMID: 23160481; PubMed Central PMCID: PMCPMC5416949.
42. Satterstrom FK, Kosmicki JA, Wang J, Breen MS, De Rubeis S, An JY, et al. Large-Scale Exome Sequencing Study Implicates Both Developmental and Functional Changes in the Neurobiology of Autism. *Cell.* 2020;180(3):568-84 e23. Epub 2020/01/26. doi: 10.1016/j.cell.2019.12.036. PubMed PMID: 31981491; PubMed Central PMCID: PMCPMC7250485.
43. Allen N, Sudlow C, Downey D, Peakman T, Danesh J, Elliott P, et al. UK Biobank: current status and what it means for epidemiology. *Health Policy Technol.* 2012;1(3):123-6.
44. Szustakowski JD, Balasubramanian S, Sasson A, Khalid S, Bronson PG, Kvikstad E, et al. Advancing Human Genetics Research and Drug Discovery through Exome Sequencing of the UK Biobank. 2020:2020.11.02.20222232. doi: 10.1101/2020.11.02.20222232 %J medRxiv.
45. Staples J, Qiao D, Cho MH, Silverman EK, University of Washington Center for Mendelian G, Nickerson DA, et al. PRIMUS: rapid reconstruction of pedigrees from genome-wide estimates of identity by descent. *Am J Hum Genet.* 2014;95(5):553-64. Epub 2014/12/03. doi: 10.1016/j.ajhg.2014.10.005. PubMed PMID: 25439724; PubMed Central PMCID: PMCPMC4225580.
46. Wang L, Zhang W, Li Q. AssocTests: An R Package for Genetic Association Studies. *Journal of Statistical Software; Vol 1, Issue 5 (2020).* 2020.
47. Price AL, Weale ME, Patterson N, Myers SR, Need AC, Shianna KV, et al. Long-range LD can confound genome scans in admixed populations. *Am J Hum Genet.* 2008;83(1):132-5; author reply 5-9. Epub 2008/07/09. doi: 10.1016/j.ajhg.2008.06.005. PubMed PMID: 18606306; PubMed Central PMCID: PMCPMC2443852.
48. Van Hout CV, Tachmazidou I, Backman JD, Hoffman JD, Liu D, Pandey AK, et al. Exome sequencing and characterization of 49,960 individuals in the UK Biobank. *Nature.* 2020;586(7831):749-56. Epub 2020/10/23. doi: 10.1038/s41586-020-2853-0. PubMed PMID: 33087929.
49. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics.* 2009;25(14):1754-60. Epub 2009/05/20. doi: 10.1093/bioinformatics/btp324. PubMed PMID: 19451168; PubMed Central PMCID: PMCPMC2705234.
50. Lin MF, Rodeh O, Penn J, Bai X, Reid JG, Krasheninina O, et al. GLnexus: joint variant calling for large cohort sequencing. *bioRxiv.* 2018:343970. doi: 10.1101/343970.
51. McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GR, Thormann A, et al. The Ensembl Variant Effect Predictor. *Genome Biol.* 2016;17(1):122. Epub 2016/06/09. doi: 10.1186/s13059-016-0974-4. PubMed PMID: 27268795; PubMed Central PMCID: PMCPMC4893825.
52. McCaw ZR, Lane JM, Saxena R, Redline S, Lin X. Operating Characteristics of the Rank-Based Inverse Normal Transformation for Quantitative Trait Analysis in Genome-Wide Association Studies. *bioRxiv.* 2019:635706. doi: 10.1101/635706.
53. Zhou W, Nielsen JB, Fritsche LG, Dey R, Gabrielsen ME, Wolford BN, et al. Efficiently controlling for case-control imbalance and sample relatedness in large-scale genetic association studies. *Nat Genet.* 2018;50(9):1335-41. Epub 2018/08/15. doi: 10.1038/s41588-018-0184-y. PubMed PMID: 30104761; PubMed Central PMCID: PMCPMC6119127.
54. Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, et al. Genome-wide genetic data on ~500,000 UK Biobank participants. *bioRxiv.* 2017:166298. doi: 10.1101/166298.
55. Ward LD, Tu H-C, Quenneville C, Flynn-Carroll AO, Parker MM, Deaton AM, et al. Genome-wide association study of circulating liver enzymes reveals an expanded role for manganese transporter SLC30A10 in liver health. 2020:2020.05.19.104570. doi: 10.1101/2020.05.19.104570 %J bioRxiv.
56. Willer CJ, Li Y, Abecasis GR. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics.* 2010;26(17):2190-1. Epub 2010/07/10. doi: 10.1093/bioinformatics/btq340. PubMed PMID: 20616382; PubMed Central PMCID: PMCPMC2922887.

57. Pruim RJ, Welch RP, Sanna S, Teslovich TM, Chines PS, Gliedt TP, et al. LocusZoom: regional visualization of genome-wide association scan results. *Bioinformatics*. 2010;26(18):2336-7. Epub 2010/07/17. doi: 10.1093/bioinformatics/btq419. PubMed PMID: 20634204; PubMed Central PMCID: PMC2935401.

## Figures

### Figure 1: Gene-level associations with glucose and HbA1c levels

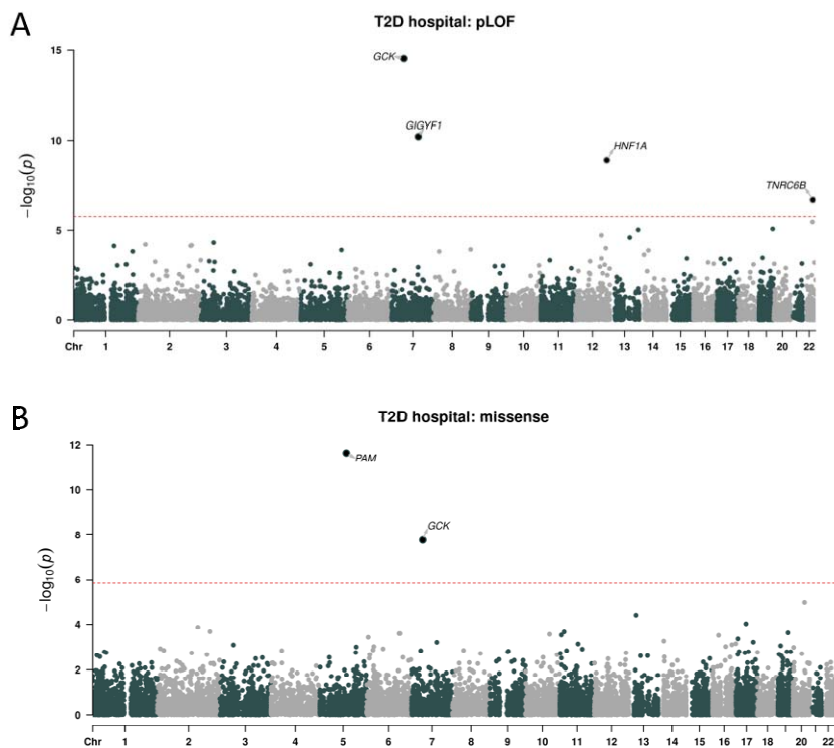
A) pLOF associations with glucose levels. B) Damaging missense variant associations with glucose levels. C) pLOF associations with HbA1c. D) Damaging missense variant associations with HbA1c levels. The red line indicates the threshold for significance, genes with significant associations are labeled.





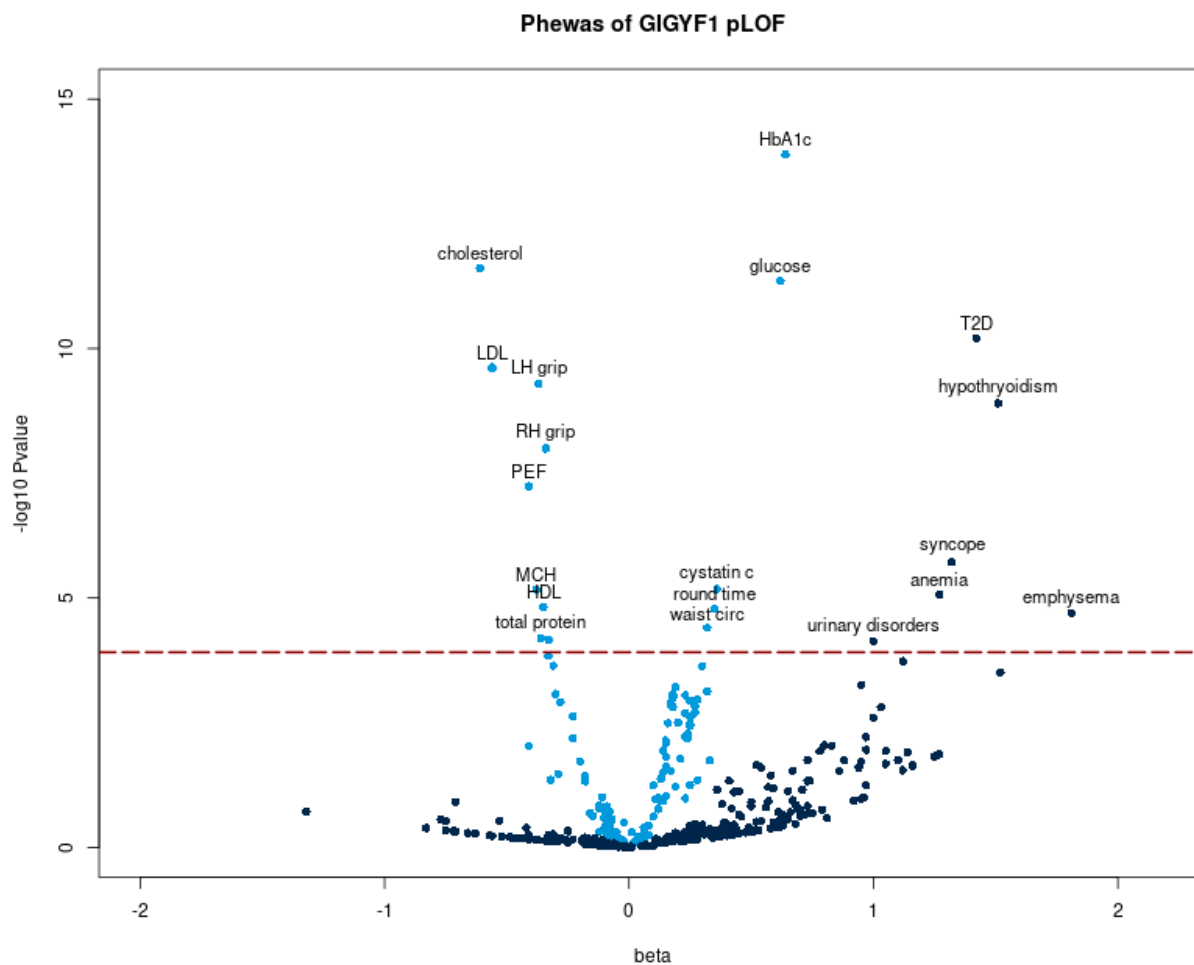
## Figure 2: Gene-level associations with T2D

A) pLOF associations with T2D diagnosis. B) Damaging missense variant associations with T2D diagnosis. The red line indicates the threshold for significance, genes with significant associations are labeled.



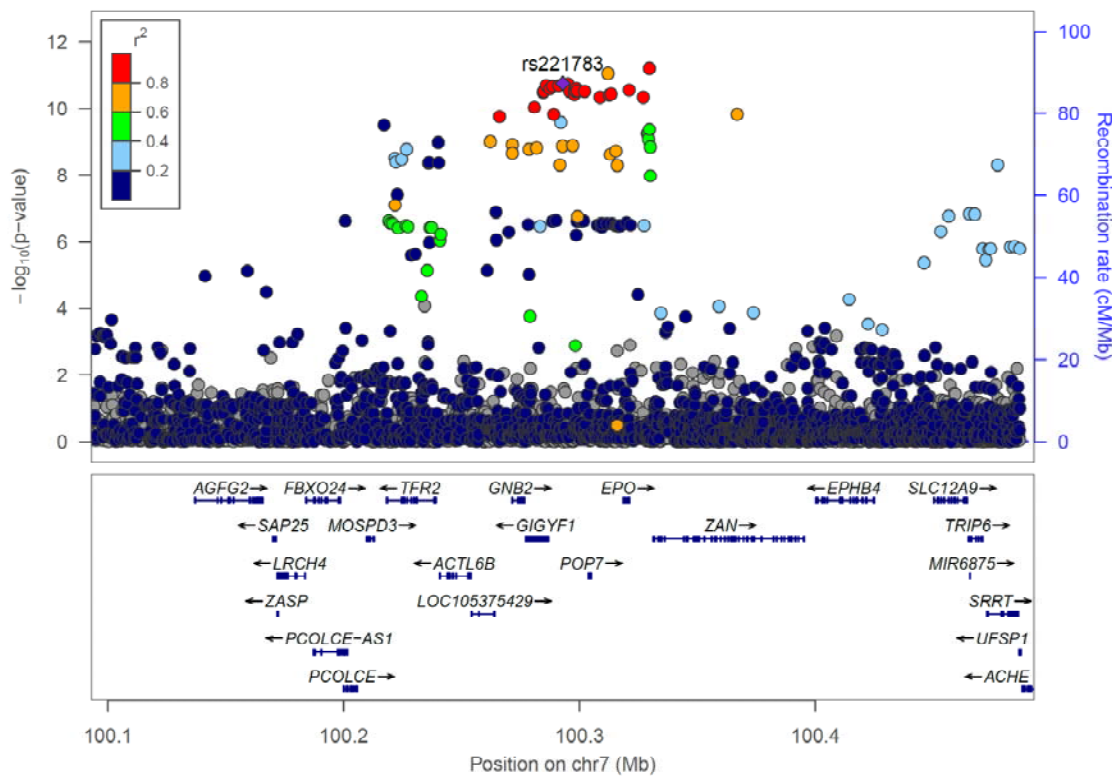
### Figure 3: PheWAS of *GIGYF1* pLOF

The x-axis is the beta (effect size in standard deviations) for the association and the y-axis is  $-\log_{10}(\text{p-value})$ . Quantitative traits are colored light blue and ICD10 diagnoses colored dark blue. Phenome-wide significant associations are labeled. The dashed line indicates the p-value threshold for phenome-wide significance. Protein; total protein, RH grip; right hand grip strength, round time: time to complete round (cognitive test), LH grip; left hand grip strength, PEF; peak expiratory flow.



## Figure 4: Locus plot of glucose associations at the *GIGYF1* locus

Association results for array genotyped and imputed variants are shown. The purple diamond represents the lead variant rs221783. Other variants are colored according to correlation ( $R^2$ ) with this marker (legend at top-left). The region displayed is chr7: 100092914-100492914. Genomic coordinates are for hg19.



## Tables

**Table 1: Gene-level associations with glucose and HbA1c levels**

Association of pLOF or damaging missense variants (CADD score  $\geq 25$ ) aggregated per gene with glucose and HbA1c levels. The effect is shown in standard deviations (SD) of transformed values as well as in International Federation of Clinical Chemistry (IFCC) units. CI; confidence interval.

gene	Variant set	title	pvalue	Effect (SD)	effect (units)	95% CI -	95% CI +	units	n carrier measured
GCK	pLOF	glucose	$1.56 \times 10^{-9}$	1.00	1.24	0.84	1.65	mmol/L	35
GIGYF1	pLOF	glucose	$4.42 \times 10^{-12}$	0.62	0.77	0.55	0.98	mmol/L	121
GCK	missense CADD $\geq 25$	glucose	$6.15 \times 10^{-11}$	0.49	0.61	0.42	0.79	mmol/L	173
G6PC2	missense CADD $\geq 25$	glucose	$4.62 \times 10^{-83}$	-0.27	-0.33	-0.36	-0.3	mmol/L	5128
GCK	pLOF	HbA1c	$2.64 \times 10^{-17}$	1.29	8.75	6.73	10.78	mmol/mol	38
GIGYF1	pLOF	HbA1c	$1.28 \times 10^{-14}$	0.64	4.33	3.23	5.43	mmol/mol	129
HNF1A	pLOF	HbA1c	$2.14 \times 10^{-7}$	0.59	4.01	2.50	5.53	mmol/mol	68
TNRC6B	pLOF	HbA1c	$2.36 \times 10^{-7}$	0.58	3.94	2.45	5.43	mmol/mol	70
RHAG	pLOF	HbA1c	$3.31 \times 10^{-34}$	-0.86	-5.81	-6.75	-4.88	mmol/mol	179
EPB41	pLOF	HbA1c	$3.14 \times 10^{-17}$	-0.53	-3.58	-4.41	-2.75	mmol/mol	226
PTPRH	pLOF	HbA1c	$4.39 \times 10^{-10}$	0.11	0.74	0.51	0.97	mmol/mol	2924
APOB	pLOF	HbA1c	$6.94 \times 10^{-8}$	0.23	1.57	1.00	2.15	mmol/mol	478
PLD1	pLOF	HbA1c	$2.99 \times 10^{-7}$	0.23	1.56	0.96	2.16	mmol/mol	438
EPB42	pLOF	HbA1c	$6.11 \times 10^{-7}$	-0.31	-2.08	-2.90	-1.26	mmol/mol	234
GCK	missense CADD $\geq 25$	HbA1c	$1.86 \times 10^{-17}$	0.56	3.83	2.94	4.71	mmol/mol	201
G6PC2	missense CADD $\geq 25$	HbA1c	$6.71 \times 10^{-45}$	-0.18	-1.21	-1.38	-1.04	mmol/mol	5574
PFAS	missense CADD $\geq 25$	HbA1c	$2.09 \times 10^{-8}$	-0.05	-0.32	-0.44	-0.21	mmol/mol	12621
PDX1	missense CADD $\geq 25$	HbA1c	$2.54 \times 10^{-7}$	0.06	0.41	0.25	0.56	mmol/mol	6694
PIEZO1	missense CADD $\geq 25$	HbA1c	$1E \times 10^{-132}$	-0.15	-1.00	-1.07	-0.92	mmol/mol	26726
AMPD3	missense CADD $\geq 25$	HbA1c	$7.28 \times 10^{-34}$	0.13	0.86	0.72	1.00	mmol/mol	8258
PFKM	missense CADD $\geq 25$	HbA1c	$2.16 \times 10^{-28}$	-0.28	-1.92	-2.26	-1.58	mmol/mol	1353
ANK1	missense CADD $\geq 25$	HbA1c	$3.12 \times 10^{-19}$	-0.13	-0.87	-1.06	-0.68	mmol/mol	4342
PFKL	missense CADD $\geq 25$	HbA1c	$2.69 \times 10^{-14}$	0.10	0.68	0.50	0.85	mmol/mol	5245
AXL	missense CADD $\geq 25$	HbA1c	$4.11 \times 10^{-12}$	-0.08	-0.54	-0.69	-0.39	mmol/mol	6827
LCAT	missense CADD $\geq 25$	HbA1c	$1.29 \times 10^{-11}$	-0.27	-1.82	-2.34	-1.29	mmol/mol	565
PLA2G12A	missense CADD $\geq 25$	HbA1c	$5.74 \times 10^{-11}$	-0.08	-0.51	-0.67	-0.36	mmol/mol	6761
PLD1	missense	HbA1c	$1.51 \times 10^{-10}$	0.08	0.57	0.40	0.75	mmol/mol	5180

	CADD $\geq$ 25								
APEH	missense CADD $\geq$ 25	HbA1c	$1.86 \times 10^{-10}$	0.29	1.96	1.36	2.57	mmol/mol	429
TNFRSF13B	missense CADD $\geq$ 25	HbA1c	$4.63 \times 10^{-9}$	-0.08	-0.51	-0.68	-0.34	mmol/mol	5470
BLVRB	missense CADD $\geq$ 25	HbA1c	$2.53 \times 10^{-8}$	-0.07	-0.48	-0.65	-0.31	mmol/mol	5533
TMC8	missense CADD $\geq$ 25	HbA1c	$3.01 \times 10^{-8}$	0.11	0.77	0.50	1.05	mmol/mol	2095
HK1	missense CADD $\geq$ 25	HbA1c	$1.08 \times 10^{-7}$	-0.16	-1.08	-1.48	-0.68	mmol/mol	988
SLC4A1	missense CADD $\geq$ 25	HbA1c	$1.76 \times 10^{-7}$	-0.15	-1.04	-1.43	-0.65	mmol/mol	1025
TFR2	missense CADD $\geq$ 25	HbA1c	$4.24 \times 10^{-7}$	-0.12	-0.84	-1.16	-0.51	mmol/mol	1491
CARHSP1	missense CADD $\geq$ 25	HbA1c	$6.78 \times 10^{-7}$	-0.13	-0.86	-1.20	-0.52	mmol/mol	1360

**Table 2: Gene-level associations with T2D diagnosis**

Association of pLOF or damaging missense variants (CADD score  $\geq$  25) aggregated per gene with T2D diagnosis. OR; odds ratio, CI; confidence interval.

gene	Variant set	title	pvalue	OR	95% CI -	95% CI +	N cases	N carrier	N carrier cases	N expected
GCK	pLOF	T2D	$2.96 \times 10^{-15}$	14.16	7.33	27.34	24695	40	19	2.71
GIGYF1	pLOF	T2D	$6.14 \times 10^{-11}$	4.15	2.71	6.37	24695	131	29	8.89
HNF1A	pLOF	T2D	$1.23 \times 10^{-9}$	5.27	3.08	9	24695	73	20	4.95
TNRC6B	pLOF	T2D	$2.00 \times 10^{-7}$	4.44	2.53	7.79	24695	71	17	4.82
PAM	missense CADD $\geq$ 25	T2D	$2.26 \times 10^{-12}$	1.31	1.21	1.41	24695	9357	801	634.85
GCK	missense CADD $\geq$ 25	T2D	$1.70 \times 10^{-8}$	2.96	2.03	4.32	24695	202	34	13.71

**Table 3: Genes and variant sets associated with multiple diabetes-related traits**

Variant sets significant for at least one trait in our primary analysis that are also associated with additional diabetes traits ( $p \leq 0.0016$ , 32 sets tested). Effect is shown in SD of transformed values or as an odds ratio (OR).

gene	Variant set	Pvalue glucose	Effect glucose	Pvalue HbA1c	Effect Hba1c	Pvalue T2D	OR T2D
GCK	pLOF	$1.56 \times 10^{-9}$	0.999	$2.64 \times 10^{-17}$	1.292	$2.96 \times 10^{-15}$	14.16
HNF1A	pLOF	0.01	0.317	$2.14 \times 10^{-7}$	0.592	$1.23 \times 10^{-9}$	5.27
GIGYF1	pLOF	$4.42 \times 10^{-12}$	0.616	$1.28 \times 10^{-14}$	0.639	$6.14 \times 10^{-11}$	4.15
GCK	missense CADD $\geq 25$	$6.15 \times 10^{-11}$	0.487	$1.86 \times 10^{-17}$	0.565	$1.70 \times 10^{-8}$	2.96
PAM	missense CADD $\geq 25$	0.92	0.001	0.009	0.026	$2.26 \times 10^{-12}$	1.31
TNRC6B	pLOF	$4.01 \times 10^{-5}$	0.507	$2.36 \times 10^{-7}$	0.582	$2.00 \times 10^{-7}$	4.44
PDX1	missense CADD $\geq 25$	0.02	0.029	$2.54 \times 10^{-7}$	0.060	$3.99 \times 10^{-5}$	1.21
PFAS	missense CADD $\geq 25$	0.32	0.009	$2.09 \times 10^{-8}$	-0.048	$4.43 \times 10^{-4}$	0.88
G6PC2	missense CADD $\geq 25$	$4.62 \times 10^{-83}$	-0.266	$6.71 \times 10^{-45}$	-0.179	0.97	1.00

**Table 4: PheWAS of GIGYF1 pLOF – quantitative traits**

Showing significant results for burden tests on quantitative traits ( $p \leq 1.22 \times 10^{-4}$ ). Effect is shown in standard deviations (SD) of transformed values. RH; right hand, LH; left hand.

gene	variant set	title	pvalue	Effect (SD)	95% CI -	95% CI +	n carrier measured
GIGYF1	pLOF	HbA1c	$1.28 \times 10^{-14}$	0.64	0.48	0.80	129
GIGYF1	pLOF	cholesterol	$2.44 \times 10^{-12}$	-0.61	-0.78	-0.44	128
GIGYF1	pLOF	glucose	$4.42 \times 10^{-12}$	0.62	0.44	0.79	121
GIGYF1	pLOF	LDL cholesterol	$2.4 \times 10^{-10}$	-0.56	-0.73	-0.38	128
GIGYF1	pLOF	apolipoprotein b	$2.52 \times 10^{-10}$	-0.56	-0.73	-0.39	127
GIGYF1	pLOF	LH grip strength	$5.11 \times 10^{-10}$	-0.37	-0.49	-0.25	131
GIGYF1	pLOF	RH grip strength	$1.01 \times 10^{-8}$	-0.34	-0.46	-0.23	131
GIGYF1	pLOF	peak expiratory flow	$5.73 \times 10^{-8}$	-0.41	-0.56	-0.26	114
GIGYF1	pLOF	cystatin c	$6.65 \times 10^{-6}$	0.36	0.20	0.51	128
GIGYF1	pLOF	mean corpuscular hemoglobin	$6.80 \times 10^{-6}$	-0.38	-0.55	-0.22	128
GIGYF1	pLOF	HDL cholesterol	$1.53 \times 10^{-5}$	-0.35	-0.52	-0.19	121
GIGYF1	pLOF	time to complete round (cognitive test)	$1.67 \times 10^{-5}$	0.35	0.19	0.51	129
GIGYF1	pLOF	waist circumference	$3.98 \times 10^{-5}$	0.32	0.16	0.47	130
GIGYF1	pLOF	total protein	$6.45 \times 10^{-5}$	-0.36	-0.53	-0.18	121
GIGYF1	pLOF	apolipoprotein a	$6.88 \times 10^{-5}$	-0.33	-0.49	-0.17	121

### Table 5: PheWAS of *GIGYF1* pLOF – ICD10-coded diagnoses

Showing significant results for burden tests on ICD10 coded diagnoses with  $\geq 500$  cases and  $\geq 1$  expected case carrier ( $p \leq 1.22 \times 10^{-4}$ ). OR; odds ratio.

gene	Variant set	title	pvalue	OR	95% CI -	95% CI +	N cases	N carrier cases	N expected
GIGYF1	pLOF	E11 T2D	$6.14 \times 10^{-11}$	4.15	2.71	6.37	24695	29	8.89
GIGYF1	pLOF	E03 other hypothyroidism	$1.25 \times 10^{-9}$	4.53	2.78	7.38	19417	21	6.99
GIGYF1	pLOF	R55 syncope and collapse	$1.90 \times 10^{-6}$	3.75	2.18	6.47	12706	15	4.57
GIGYF1	pLOF	D50 iron deficiency anemia	$8.52 \times 10^{-6}$	3.56	2.04	6.23	12886	14	4.64
GIGYF1	pLOF	J43 emphysema	$1.99 \times 10^{-5}$	6.13	2.67	14.10	3015	6	1.09
GIGYF1	pLOF	N39 other disorders of urinary system	$7.32 \times 10^{-5}$	2.71	1.66	4.45	24581	19	8.85

### Table 6: Common variant associations at the *GIGYF1* locus

Associations for the array-typed variant rs221783. For quantitative traits the effect is shown in standard deviations (beta) and for diagnoses as an odds ratio (OR). MAF; minor allele frequency.

phenotype	chrom	Pos (hg19/hg38)	Ref (effect allele)	Alt	rsid	MAF	pvalue	Effect (beta/OR)	95% CI -	95% CI +
glucose	7	100292914/100695291	T	C	rs221783	11%	$1.82 \times 10^{-11}$	-0.03	-0.03	-0.02
HbA1c	7	100292914/100695291	T	C	rs221783	11%	$3.58 \times 10^{-7}$	-0.02	-0.03	-0.01
Cholesterol	7	100292914/100695291	T	C	rs221783	11%	$7.00 \times 10^{-12}$	0.03	0.02	0.03
LDL	7	100292914/100695291	T	C	rs221783	11%	$6.25 \times 10^{-10}$	0.02	0.02	0.03
T2D	7	100292914/100695291	T	C	rs221783	11%	0.005	0.96	0.93	0.99
Hypothyroidism	7	100292914/100695291	T	C	rs221783	11%	$6.95 \times 10^{-7}$	0.92	0.88	0.95