

1
2
3
4
5 **Supplementary Information**
6
7
8

9 **Analysis of rare genetic variation underlying cardiometabolic diseases**
10 **and traits among 200,000 individuals in the UK Biobank**
11
12
13
14
15

16 Sean J. Jurgens^{1,2,†}, Seung Hoan Choi^{1,†}, Valerie N. Morrill^{1,†}, Mark Chaffin¹, James P. Pirruccello^{1,3},
17 Jennifer L. Halford¹, Lu-Chen Weng^{1,3}, Victor Nauffal^{1,3}, Carolina Roselli¹, Amelia W. Hall^{1,3},
18 Matthew T. Oetjens⁴, Braxton Lagerman⁵, David P. vanMaanen⁵, Regeneron Genetics Center^{6,*}, Krishna G.
19 Aragam^{1,3}, Kathryn L. Lunetta^{7,8}, Christopher M. Haggerty^{5,9}, Steven A. Lubitz^{1,3,10,#} & Patrick T. Ellinor^{1,3,10,#}
20
21

- 22 1. Cardiovascular Disease Initiative, The Broad Institute of MIT and Harvard, Cambridge, MA, USA
23 2. Department of Experimental Cardiology, Amsterdam UMC, University of Amsterdam, Amsterdam, NL
24 3. Cardiovascular Research Center, Massachusetts General Hospital, Boston, MA, USA
25 4. Autism & Developmental Medicine Institute, Geisinger, Danville, PA, USA
26 5. Department of Translational Data Science and Informatics, Geisinger, Danville, PA, USA
27 6. Regeneron Genetics Center, Tarrytown, NY, USA
28 7. NHLBI and Boston University's Framingham Heart Study, Framingham, MA, USA
29 8. Department of Biostatistics, Boston University School of Public Health, Boston, MA, USA
30 9. Heart Institute, Geisinger, Danville, PA, USA.
31 10. Demoulas Center for Cardiac Arrhythmias, Massachusetts General Hospital, Boston, MA, USA
32
33

34 † Contributed equally to this work.

35 * List of consortium authors and their affiliations appear at the end of the main manuscript document

36 # Jointly supervised this work.

Items	Pages
Supplementary Note: 1. Consortium Author List and Contribution Statements	
Regeneron Genetics Center	5
Supplementary Note: 2. Supplementary Methods	
Whole exome sequencing dataset	6
Quality control	6-7
Relationship inference, kinship matrix and principal component analysis	7-8
Variant annotation for missense variants	8
Evaluation of test statistic inflation in exome-wide gene-based testing	8
Analysis of rare synonymous variants	8-9
Analysis of common variation near rare variant signals in the UK Biobank	9
Common variant results from the Type 2 Diabetes Portal	9
GTEx expression-QTL and splice-QTL data for common variants	9
Clinical variants from the ClinVar database	9
<i>TTN</i> exons highly expressed in left ventricle tissue	9-10
Supplementary Note: 3. Supplementary Results and Discussion	
Evaluation of inflation in gene-based analyses	11
Associations between cardiac phenotypes and variants in <i>TTN</i> exons highly expressed in the heart	11
Common variants near genes with novel rare variant associations	11-12
Penetrance of predicted-deleterious variants in the UK Biobank	12-14
Yield of putatively pathogenic variants among disease cases	14
Supplementary Tables	
Supplementary Table 1: Definitions for curated disease phenotypes	Excel file
Supplementary Table 2: Baseline characteristics for each binary disease trait	Excel file
Supplementary Table 3: Baseline characteristics for each quantitative trait	Excel file
Supplementary Table 4: Gene-disease associations with Q-value<0.05	Excel file
Supplementary Table 5: Gene-quantitative trait associations with Q-value<0.05	Excel file
Supplementary Table 6: Single variant information and LOVO results for variants from significant binary disease associations	Excel file
Supplementary Table 7: Single variant information and LOVO results for variants from significant quantitative associations	Excel file
Supplementary Table 8: Gene-disease associations conditioned on nearby associated common variants	Excel file
Supplementary Table 9: Gene-quantitative trait associations conditioned on nearby associated common variants	Excel file

	Supplementary Table 10: Penetrance of variants in genes associated with increased risk of disease in the discovery phase and in the panel analyses.	Excel file
	Supplementary Table 11: Look-up of gene-based common-variant results in the Type 2 Diabetes Portal for genes with novel metabolic or anthropometric rare-variant-associations	Excel file
	Supplementary Table 12: Coding consequence, eQTL and sQTL data for most significant common variant near genes with novel rare-variant-associations	Excel file
	Supplementary Table 13: Association results for known and novel metabolic blood marker genes with cardiometabolic traits and conditions	Excel file
	Supplementary Table 14: Genes from cardiomyopathy and arrhythmia, hypercholesterolemia and monogenic diabetes panel and reported mode of inheritance for relevant disease in OMIM	Excel file
	Supplementary Table 15: Frequency of loss-of-function, pathogenic and likely pathogenic variants in cardiac disease and diabetes panel genes in the UK Biobank	Excel file
	Supplementary Table 16: Association results for loss-of-function, clinically pathogenic and likely pathogenic variants from cardiac disease and diabetes panels	Excel file
	Supplementary Table 17: Prevalence of putatively pathogenic variants in cardiovascular disease and diabetes genes among disease cases	Excel file
Supplemental Figures and Figure Legend		
	Supplementary Figure 1. Principal component analysis and self-reported ancestries for UK Biobank WES samples.	16
	Supplementary Figure 2. Quantile-quantile plots for exome-wide gene-based tests across all binary and all quantitative phenotypes.	17-18
	Supplementary Figure 3. Quantile-quantile plots for exome-wide gene-based tests for each individual binary trait.	19-26
	Supplementary Figure 4. Quantile-quantile plots for exome-wide gene-based tests for each individual quantitative trait.	27-30
	Supplementary Figure 5. Quantile-quantile plots for rare deleterious variants compared to rare synonymous variants for height, weight, BMI, and QTc.	31-32
	Supplementary Figure 6: Sensitivity analysis restricting to individuals of European ancestry only in the analysis of binary traits.	33-34
	Supplementary Figure 7. Sensitivity analysis restricting to individuals of European ancestry only in the analysis of quantitative traits.	35-36
	Supplementary Figure 8. Sensitivity analysis restricting to LOFs only in the primary analysis of binary traits.	37-38
	Supplementary Figure 9. Sensitivity analysis restricting to LOFs only in the primary analysis of quantitative traits.	39-40
	Supplementary Figure 10. Leave-one-variant-out (LOVO) analysis for novel rare variant associations.	41-43
	Supplementary Figure 11. Penetrance of predicted-damaging variants in genes associated with disease in the primary analyses.	44
	Supplementary Figure 12. Prevalence of predicted-damaging variants in genes identified in primary analysis among relevant disease cases.	45
	Supplementary Figure 13. Penetrance of putatively pathogenic variants in cardiovascular disease and diabetes panel genes for relevant phenotypes.	46

	Supplementary Figure 14. Prevalence of putatively pathogenic variants in cardiovascular disease and diabetes panel genes among disease cases.	47
Supplementary References		48-49

38 **Supplementary Note**

39
40 **1. Consortium Author List and Contribution Statements**

41
42 **Regeneron Genetics Center**

43 All authors/contributors are listed in alphabetical order.

44
45 **RGC Management and Leadership Team**

46 Goncalo Abecasis, Aris Baras, Michael Cantor, Giovanni Coppola, Andrew Deubler, Aris Economides, Luca A.
47 Lotta, John D. Overton, Jeffrey G. Reid, Alan Shuldiner, Katia Karalis and Katherine Siminovitch
48 Contribution: All authors contributed to securing funding, study design and oversight. All authors reviewed the
49 final version of the manuscript.

50
51 **Sequencing and Lab Operations**

52 Christina Beechert, Caitlin Forsythe, Erin D. Fuller, Zhenhua Gu, Michael Lattari, Alexander Lopez, John D.
53 Overton, Thomas D. Schleicher, Maria Sotiropoulos Padilla, Louis Widom, Sarah E. Wolf, Manasi Pradhan, Kia
54 Manoochehri, Ricardo H. Ulloa.
55 Contribution: C.B., C.F., A.L., and J.D.O. performed and are responsible for sample genotyping. C.B, C.F.,
56 E.D.F., M.L., M.S.P., L.W., S.E.W., A.L., and J.D.O. performed and are responsible for exome sequencing.
57 T.D.S., Z.G., A.L., and J.D.O. conceived and are responsible for laboratory automation. M.S.P., K.M., R.U.,
58 and J.D.O are responsible for sample tracking and the library information management system.

59
60 **Genome Informatics**

61 Xiaodong Bai, Suganthi Balasubramanian, Boris Boutkov, Gisu Eom, Lukas Habegger, Alicia Hawes, Shareef
62 Khalid, Olga Krasheninina, Rouel Lanche, Adam J. Mansfield, Evan K. Maxwell, Mona Nafde, Sean O’Keeffe,
63 Max Orelus, Razvan Panea, Tommy Polanco, Ayesha Rasool, Jeffrey G. Reid, William Salerno, Jeffrey C.
64 Staples,
65 Contribution: X.B., A.H., O.K., A.M., S.O., R.P., T.P., A.R., W.S. and J.G.R. performed and are responsible for
66 the compute logistics, analysis and infrastructure needed to produce exome and genotype data. G.E., M.O.,
67 M.N. and J.G.R. provided compute infrastructure development and operational support. S.B., S.K., and J.G.R.
68 provide variant and gene annotations and their functional interpretation of variants. E.M., J.S., R.L., B.B., A.B.,
69 L.H., J.G.R. conceived and are responsible for creating, developing, and deploying analysis platforms and
70 computational methods for analyzing genomic data.

71
72 **Clinical Informatics:**

73 Michael Cantor and Dadong Li

74 Contribution: All authors contributed to the clinical informatics of the project

75
76 **Translational Genetics:**

77 Niek Verweij, Jonas Nielsen, Tanima De and Manuel A. R. Ferreira

78 Contribution: All authors contributed to the review process for the final version of the manuscript.

79
80 **Research Program Management**

81 Marcus B. Jones, Jason Mighty and Lyndon J. Mitnaul

82 Contribution: All authors contributed to the management and coordination of all research activities, planning
83 and execution. All authors contributed to the review process for the final version of the manuscript.

2. Supplementary Methods

Whole exome sequencing dataset

Whole exome sequencing (WES)

Exomes were captured with the IDT xGen Exome Research Panel v1.0 including supplemental probes. The basic design targets 39Mbp of the human genome (19,396 genes). Multiplexed samples were sequenced with dual-indexed 75x75bp paired-end reads on the Illumina NovaSeq 6000 platform using S2 (initial 50k samples) and S4 flow cells (all subsequent samples). In each sample and among targeted bases, coverage exceeds 20X at 95% of sites on average. More information is available on the UK Biobank website (<https://biobank.ctsu.ox.ac.uk/showcase/label.cgi?id=170>).

Variant calling for the OQFE dataset

In the present analysis, we used the pVCF files from the OQFE dataset¹. Briefly, all reads were duplicate marked and aligned to genome build GRCh38 in an alt-aware manner as described in the Functional Equivalence protocol². Variants were called per-sample using DeepVariant, after which individual level VCF files were combined and joint-genotyped using GLnexus³. More information is available on the UK Biobank website (<https://biobank.ctsu.ox.ac.uk/showcase/label.cgi?id=170>).

Quality control

In addition to any quality-control that was performed centrally, we applied extensive additional genotype, variant and sample quality-control procedures to ensure a high-quality dataset for analyses. To this end, we utilized the OQFE WES pVCF files provided by the UK Biobank, which contained calls for 200,643 sequenced samples.

Genotype quality control

We applied genotype refinement to the raw genotype calls in the pVCF files using Hail. We first split multi-allelic sites to represent separate bi-allelic sites. All calls that did not pass the following hard filters were then set to no-call in our analysis:

- For homozygous reference calls: Genotype Quality < 20; Genotype Depth < 10; Genotype Depth > 200
- For heterozygous calls: (A1 Depth + A2 Depth)/Total Depth < 0.9; A2 Depth/Total Depth < 0.2; Genotype likelihood[ref/ref] < 20; Genotype Depth < 10; Genotype Depth > 200
- For homozygous alternative calls: (A1 Depth + A2 Depth)/Total Depth < 0.9; A2 Depth/Total Depth < 0.9; Genotype likelihood[ref/ref] < 20; Genotype Depth < 10; Genotype Depth > 200

These filters removed 9% of the 3,573,574,459,423 raw genotype calls leaving 3,214,727,581,104 genotype calls across 17,981,897 variant sites and 200,643 samples.

Variant quality control

We then performed variant level filters. We removed variants that failed the following filters:

- Call rate of < 90% (restricting to males for Y chromosomal markers) (N= 4,023,284)
- Failed a liberal Hardy-Weinberg Equilibrium test (HWE) at $P < 10^{-15}$ among unrelated samples (not applied to Y chromosomal markers) (N=136,869)
- Present in Ensembl low-complexity regions (N=748,116)
- Monomorphic in the final dataset (N=55,614)

After performing these variant filters, 13,003,057 variants remained of which 12,756,075 were autosomal.

High-quality variants for sample quality control and relationship inference

To perform sample level quality control and kinship inference, we defined three subsets of genetic variants that were independent and of very high-quality:

- 'High-quality independent autosomal variants subset' with MAF > 0.1%, missingness < 1%, HWE $P > 10^{-6}$ and two rounds of pruning using `--indep-pairwise 200 100 0.1` and `--indep-pairwise 200 100 0.05` in PLINK⁴ (81,121 variants).
- 'WES-vs-array independent autosomal variants subset' with MAF > 0.1%, missingness < 1% and HWE $P > 10^{-6}$ in both the WES dataset and in the genotyping array data provided by the UK Biobank⁵ (among participants who had both available). We further removed indels and ambiguous SNPs and performed two rounds of pruning (24,207 variants).
- 'High quality independent X-chromosomal subset' with missingness < 1%, HWE $P > 10^{-6}$, not within pseudo-autosomal regions, and two rounds of pruning.

Sample quality control

We computed a number of quality metrics to identify bad-quality or duplicated samples. We first used KING⁶ (version 2.2.5) to calculate pairwise heterozygote concordance rates for each pair of samples, using the high-quality independent autosomal markers. Then we used the high-quality autosomal variants present in both WES and array datasets to compute per-sample heterozygote concordance rates between WES calls and genotyping array calls. We inferred the genetic sex of each participants with the `--check-sex` option in PLINK, using the high-quality independent X-chromosomal markers. We set any sample with $F > 0.8$ to male, while samples with $F < 0.5$ were set to female. Finally, using all ~12.7M autosomal WES variants, we computed a number of additional metrics including sample call rate, transition/transversion ratio (Ti/Tv), heterozygote/homozygote ratio (Het/Hom), SNV/indel ratio (SNV/indel) and the number of singletons. After computing these metrics, we excluded participants based on the following criteria:

- Decided to revoke their consent (N=13)
- Sample duplicates based on heterozygote concordance rates > 0.8 (N=0) (27 putative genetic duplicates could be resolved as monozygotic twins and were not removed)
- Samples with blatant discordance between self-reported and genetically inferred sex (N=80)
- Discordance between WES and array calls with heterozygote concordance rates < 0.8 (N=0)
- Call rate < 90% (N=1)
- Samples further than 8 standard deviations from the mean for Ti/Tv (n=0), Het/Hom (N=100), SNV/indel (N=1) and number of singletons (N=111)

After applying these filters 200,337 samples remained for analysis.

Relationship inference, kinship matrix and principal component analysis

Kinship inference and kinship matrix

We used the KING-robust algorithm to compute pairwise kinship estimates for all samples in the dataset (using the high-quality independent autosomal variants). We then retained all information on pairs estimated to be genetically related to one another at 3rd degree or closer (kinship coefficient ≥ 0.0442). We used this data to construct a sparse kinship matrix in which all relationships with kinship coefficient < 0.0442 were set to 0. Finally, we scaled the values in this matrix so it had a diagonal of 1 (as opposed to 0.5 on the KING kinship scale).

Unrelated subset

We defined an unrelated subset of the WES cohort, where no relationships with kinship coefficient ≥ 0.0442 remained, a threshold that excludes any individuals related at 3rd degree or closer. To maximize the sample size of this unrelated subset, we first iteratively removed individuals related to multiple other individuals until none remained. We then removed one sample from each remaining pair at random, leaving 185,990 unrelated samples.

Principal component analysis

PCAir⁷ (from GENESIS version 2.18.0) was used to calculate the top 20 ancestral principal components (using the high-quality independent autosomal variants), while implementing a randomized algorithm⁸ for computational efficiency. We performed this analysis among the unrelated subset of the cohort, after which the remaining samples were projected onto the PCs.

Variant annotation for missense variants

Missense variants annotated from VEP incorporated 30 in-silico prediction tools from the dbNSFP database (version 4.1a.) These tools included qualitative prediction algorithms (SIFT, SIFT4G, Polyphen2 HDIV, Polyphen2 HVAR, LRT, MutationTaster, FATHMM, PROVEAN, MetaSVM, MetaLR, MCAP, PrimateAI, DEOGEN2, BayesDel addAF, BayesDel noAF, ClinPred, LIST-S2, fathmm-MKL coding, fathmm-XF coding, MutationAssessor, and Aloft) and quantitative algorithms (VEST4, REVEL, MutPred, MVP, MPC, DANN, CADD, Eigen, and Eigen-PC). When the qualitative prediction tools (except for MutationAssessor and Aloft) indicated "D" for a variant, the variant gained one score from each algorithm. An indicator for a deleterious variant of MutationAssessor was "H" and of Aloft was "R" or "D" with high confidence. For the quantitative algorithms, when the variant indicators were higher than 90% of predicted variants in the entire dataset, a variant gained one score from each quantitative algorithm. Then, if a variant was annotated with more than seven prediction tools (over 20% out of the 30 tools), and the proportion of the deleterious score (total gained score / # none missing prediction tools) was greater than or equal to 0.9, we included the variant in the gene-based analyses.

Evaluation of test statistic inflation in exome-wide gene-based testing

To inspect the calibration of test statistics in our analyses, we visually inspected quantile-quantile (QQ) plots across all performed tests for binary traits, and across all performed tests for quantitative traits. To evaluate the effect of the minimum carrier count we employed (≥ 20 rare variant carriers), we further made QQ plots for all test with ≥ 50 , and all tests with ≥ 200 rare variant carriers. Then we inspected per-trait QQ plots and computed per-trait λ values, representing a statistic of inflation where a value of 1 indicates perfect calibration of P -values. For quantitative traits, we computed λ at the median ($\lambda[\text{median}]$), as is conventional. For binary traits, the Saddle Point Approximation was only applied to tests reaching $P < 0.05$ in an initial regular score test. For this reason, $\lambda(\text{median})$ may not accurately represent test statistic inflation at the tails of the distribution. We therefore computed λ values for binary traits at the tail of the test statistic distribution, by comparing empirical values at the 95th quantile to the expected value at the 95th quantile ($\lambda[q0.95]$).

Analysis of rare synonymous variants

For any trait showing unexplained test statistic inflation, we analyzed rare synonymous variants. Synonymous variants are generally expected to have no protein consequence, and therefore represent a class of genetic variation that should produce a null distribution. As with the predicted-deleterious variants, we pooled rare synonymous variants variant by gene and performed a collapsing test. Variants were considered rare if they had $\text{MAF} < 0.1\%$ in the UK Biobank WES dataset and $\text{MAF} < 0.1\%$ in five major gnomAD⁹ populations.

Analysis of common variation near rare variant signals in the UK Biobank

To identify common variant associations near the identified rare variant signals, we ran common variant association analyses in the genomic region 500KB downstream and upstream of the identified gene. To this end, we utilized the UK Biobank version 3 imputed data. Details on genotyping and quality-control have been described previously⁵. Briefly, samples were genotyped using Affymetrix UK biobank Axiom (450,000 samples) and Affymetrix UK BiLEVE axion (50,000 samples) arrays. Genetic data were then imputed to the Haplotype Reference Consortium panel and UK10K + 1000 Genomes panels. For the common variant association analyses, we removed samples that were outliers for heterozygosity or missingness, samples with putative sex

227 chromosome aneuploidy, samples with a mismatch between self-reported and genetically inferred sex,
228 samples not included in the central kinship inference, and samples who had revoked their consent. Imputed
229 variants with $MAF < 0.5\%$ and $INFO < 0.3$ were removed. We ran two-sided common variant association tests
230 using PLINK⁴. Logistic regression was used for binary phenotypes and linear regression for continuous traits.
231 We analyzed all unrelated UK Biobank individuals with imputed data and relevant phenotypic data available.
232 We adjusted for age, sex, genotyping array and associated ancestral principal components ($P < 0.05$). Common
233 variants with $P < 1 \times 10^{-5}$ were considered significant.

234 **Common variant results from the Type 2 Diabetes Portal (T2DKP)**

235 For genes in which we identified novel rare variant associations for metabolic and anthropometric traits, we aimed to
236 find additional evidence for the role of these genes using publicly-available common variant results from the
237 T2DKP¹⁰. We used gene-based common variant results downloaded on the 7th of December 2020, displayed in
238 **Supplementary Table 11**. Gene-based results were based on single variant summary statistics from many large-
239 scale common variant GWAS. In short, the portal first filters summary statistics to include only biallelic markers with
240 no missing data, and then separates variants by frequency (common vs rare) and by ancestry. Then, it meta-
241 analyses GWAS results for common variants using METAL¹¹ (with OVERLAP ON) in an ancestry-specific manner,
242 after which it performs a trans-ancestry meta-analysis using METAL (with OVERLAP OFF). Gene-based common
243 variant analyses were subsequently performed using the Multi-marker Analysis of GenoMic Annotation method
244 (MAGMA)¹². We further identified index single variants using the 'explore region' option for a given gene; the most
245 significant single variant in a gene region for a given phenotype were extracted from the T2DKP data on the 7th of
246 June 2021, displayed in **Supplementary Table 12**.

247 **GTEx expression-QTL and splice-QTL data for common variants**

248 For each index single variant from the T2DKP mentioned above, we leveraged data from GTEx to identify
249 significant expression-QTL and splice-QTL associations. For a few phenotypes, no data in the T2DKP was
250 available; in these cases, we used the index variants from our imputed common variant analyses in the UK
251 Biobank (e.g. supraventricular tachycardia and *TTN*). We extracted expression-QTL and splice-QTL data for
252 index variants and the relevant gene from the GTEx version 8 dataset (<https://gtexportal.org/home/>), on the 7th
253 of June 2021. The dataset consists of RNA-sequencing and whole-genome sequencing data from 838 donors
254 after previously described quality-control¹³; 49 tissues or cell lines had at least 70 individuals with both data
255 sources available (15,201 total samples) for expression-QTL and splice-QTL analysis. We determined that a
256 variant was a significant expression-QTL for the given gene if i) the expression-QTL reached tissue-specific
257 FDR5%, as described previously¹³ and ii) had Bonferroni-corrected $P < 7.3 \times 10^{-5} = 0.05 / (14 \text{ variants} \times 49$
258 tissues) in our analysis. Splice-QTLs were determined to be significant if i) the splice-QTL reached tissue-
259 specific FDR5%, as described previously¹³ and ii) had Bonferroni-corrected $P < 1.6 \times 10^{-6} = 0.05 / (629 \text{ introns} \times$
260 49 tissues) in our analysis. Results for this lookup are displayed in **Supplementary Table 12**.

261 **Clinical variants from the ClinVar database**

262 To identify pathogenic rare variants, we used the ClinVar database. We downloaded the ClinVar dataset on
263 11/2020. Variants that were not submitted by clinical testing labs or which were evaluated before 2015 were
264 excluded from our analyses. We used the clinical significance interpretation at the most recent submission.
265 The clinical significance interpretation included Pathogenic, Likely-Pathogenic, Likely-Benign, Benign, Variant
266 of Uncertain Significance, and Conflicting data from submitters; we only used variants with the Pathogenic or
267 Likely-Pathogenic classification in the present study.

268 ***TTN* exons highly expressed in left ventricle tissue**

269 Previous work described that distinguishing highly expressed *TTN* exons in heart tissues is important to
270 understand phenotypic presentation^{14,15}. As post-hoc analyses, we performed association tests between
271

275 deleterious variant in highly expressed (Percentage Spliced-In [PSI] \geq 90%) in left ventricular tissue¹⁴ and
276 cardiac traits using the same model implemented in our primary analyses.

3. Supplementary Results and Discussion

Evaluation of inflation in gene-based analyses

QQ plots for P -values from all performed tests in the discovery phase (all quantitative and all binary) did not show any inflation (**Supplementary Figure 2**). We also made QQ plots restricting to tests with at least 50 variant carriers and tests with at least 200 rare variant carriers (**Supplementary Figure 2**). QQ plots showed a similar distribution of P -values without clear inflation. We then inspected QQ plots for individual traits (**Supplementary Figures 3-4**). Most traits did not show evidence of inflation ($\lambda < 1.1$); however for three traits, height, weight and QTc, lambda values were consistent with moderate inflation ($1.1 \leq \lambda_{GC} < 1.25$). Indeed, height and weight had visually inflated distributions of P -values. Such inflation could be due to biases such as population stratification, or alternatively due to a high degree of polygenicity. To distinguish between these causes, we analyzed rare synonymous variants for these traits. Seeing as most synonymous variants are expected to have no protein consequence, such an analysis should yield a null distribution. We found that rare synonymous variants indeed yielded a distribution of P -values consistent with the null ($\lambda_{GC} < 1.05$) for each traits (**Supplementary Figure 5**), implying that a large proportion of the observed inflation was due to polygenicity rather than bias.

Associations between cardiac phenotypes and variants in *TTN* exons highly expressed in the heart

Concordant with our prior knowledge *TTN* associations with heart failure, atrial fibrillation, dilated cardiomyopathy, left ventricle ejection fraction and left ventricular end systolic volume strengthened after restricting to variants in cardiac expressed exons. Supraventricular tachycardia ($P = 3.0 \times 10^{-12}$), ventricular arrhythmia ($P = 2.6 \times 10^{-10}$), and mitral valve disease ($P = 5.4 \times 10^{-15}$) also showed markedly stronger associations when restricting to cardiac exons of *TTN*. Furthermore, implantable cardioverter defibrillator ($P = 6.6 \times 10^{-9}$), tricuspid valve disease ($P = 9.7 \times 10^{-7}$), RR interval ($P = 2.6 \times 10^{-6}$), Pulse rate ($P = 1.1 \times 10^{-25}$) and LVESVi ($P = 1.8 \times 10^{-7}$) were significantly associated with variants in cardiac exons of *TTN*.

Common variants near genes with novel rare variant associations

Among our novel associations were 3 associations for rare variants in *GIGYF1*, namely for increased risk of type 2 diabetes, elevated glucose levels and lower low-density lipoprotein levels. In accordance, common variants near *GIGYF1* were associated with all these traits (**Supplementary Table 11**). The top common variants for each of these traits in the *GIGYF1* locus are expression-QTLs for *GIGYF1* in many tissues (**Supplementary Table 12**), including many relevant tissues such as adipose tissue, pancreas, skeletal muscle, thyroid and pituitary, as well as many other brain and gastro-intestinal tissues. The alleles associated with lower *GIGYF1* expression were consistently associated with increased risk of diabetes, higher glucose and lower low-density lipoprotein levels across tissues, in strong concordance with the observed LOF associations. These results suggest that higher *GIGYF1* levels may be protective for diabetes. *CCAR2* rare variants were associated with increased risk of diabetes, and common variants near the locus were as well (**Supplementary Table 11**). The top *CCAR2* common variant was a significant expression-QTL for *CCAR2* across many tissues, including adipose tissue, skeletal muscle, pancreas, thyroid and multiple brain and gastro-intestinal tissues (**Supplementary Table 12**). Generally, the alleles associated with higher *CCAR2* expression were associated with higher risk of diabetes, which is not directly consistent with the observed LOF associations, although the sign was flipped in certain tissues such as fibroblasts.

Rare variants in *TTN* were novelly associated with mitral valve disease and supraventricular tachycardia. A common variant near *TTN* was also found to be associated with supraventricular tachycardia (rs10167882, $P = 2.1 \times 10^{-6}$, OR 1.11; **Supplementary Table 12**). This variant was not found to be a significant expression-QTL or splice-QTL for *TTN*, although it is in LD with a missense variant that also shows evidence of

324 association with supraventricular tachycardia as well (p.Gln8542His, $P=0.0015$, OR 1.23; **Supplementary**
325 **Table 12**).

326
327 Rare variants in *NR1H3* were associated with high-density lipoprotein in our primary analysis, and a common
328 variant association was also found at this locus (**Supplementary Table 11**). The top common variant in this
329 locus was an expression-QTL for *NR1H3*, with the alleles associated with higher high-density lipoprotein being
330 associated with increased *NR1H3* expression in some tissues (for example subcutaneous adipose tissue) and
331 decreased expression in others (whole blood, brain cortex) (**Supplementary Table 12**). The top common
332 variant was also a significant splice-QTL for *NR1H3* across many tissues including adipose tissue, with
333 consistent tissue effects, and was also in LD with an *NR1H3* missense variant (p.Ala101Val, $P=1.9 \times 10^{-21}$)
334 (**Supplementary Table 12**).

335
336 Among our novel rare variant associations were 7 associations for height, namely *DTL*, *PIEZO1*, *SCUBE3*,
337 *ANGPTL2*, *PAPPA*, *IRS1* and *ZFAT*. All of these genes are supported by significant nearby common variant
338 associations (**Supplementary 11**). For *IRS1*, we found that the top common variant had two splice-QTL
339 associations with *IRS1*: one in fibroblasts and thyroid, and another in subcutaneous adipose tissue
340 (**Supplementary Table 12**). For *PAPPA*, we found that the top common variant was in LD with a *PAPPA*
341 missense variant (p.Ser1224Tyr). Finally, for *SCUBE3*, we found that the top common variant was a
342 suggestive ($P=0.00019$) expression-QTL for *SCUBE3* in fibroblasts, with the allele associated with lower
343 *SCUBE3* expression being associated with shorter stature, consistent with the observed LOF association. The
344 relative absence of additional expression-QTL and splice-QTL data for the remaining common variant height
345 loci might be a reflection of the adult population in GTEx; relevant expression-QTLs for height may be
346 predominantly developmental and possibly not present in adult tissue.

347 348 **Penetrance of predicted-deleterious and pathogenic variants in the UK Biobank**

349 In our primary analyses, 10 genes were significantly associated (Q-value < 0.01) with increased risk of a
350 disease or medical condition. For those 10 genes, 3371 participants (1.6% of the sample) carried predicted-
351 deleterious variants (LOF and predicted-deleterious missense variants). Among 3371 carriers, 621 (18.4%
352 penetrance) developed at least one medical condition. When we liberalize our significant threshold to FDR Q-
353 value 0.05, there were 15 genes associated with at least one medical condition. We found 3762 participants
354 (1.9% of the sample) who carried deleterious variants; meanwhile 693 (18.4% penetrance) developed an
355 associated disease. The penetrance of respective genes and traits are illustrated in **Supplementary Figure**
356 **11**. The highest penetrance was 71% [95%CI 61-79%] from *LDLR* for hypercholesterolemia. *PKD1* mutations
357 were associated with 47% penetrance for chronic kidney disease [95%CI 33-62%]. *PKD1* pathogenic
358 mutations are known for causing highly-penetrant autosomal dominant polycystic kidney disease, with end
359 stage kidney disease reached at around 58 years¹⁶. However, a higher-than-expected frequency of *PKD1*
360 mutations in healthy sequenced populations has recently been described, suggesting incomplete penetrance¹⁷.
361 We note, however, that *PKD1* has many pseudo-genes which may complicate read-mapping, and Sanger
362 sequencing validation is often performed in clinical settings to confirm *PKD1* variants. Despite this fact,
363 previous studies have mainly shown decreased sensitivity when utilizing next-generation sequencing;
364 specificity ranges from 90-100% when stringent QC filters are applied¹⁸⁻²⁰. Still, we cannot exclude the
365 possibility of some alignment issues, which may downward bias penetrance estimates for this gene.

366
367 The penetrance of putatively pathogenic variants in genes included in our panel analysis (*InVitae*
368 *Cardiomyopathy and Arrhythmia* panel, *InVitae hypercholesterolemia* panel and *InVitae Monogenic Diabetes*
369 panel) are shown in **Supplementary Figure 13**. The penetrance of cardiovascular disease variants was
370 generally modest (**Supplementary Table 10**). Of *TTNtv* carriers, 16% [95%CI 14-19%] had diagnoses of atrial
371 fibrillation, 9.6% [95%CI 7.9-12%] of heart failure, 4.0% [95%CI 2.8-5.6%] of dilated cardiomyopathy and 3.8%

[95%CI 2.6-5.3%] of ventricular arrhythmia (**Supplementary Table 10**), considerably lower than the incidence of these diseases in previous family-member based analyses^{21,22}, although for cardiomyopathy and atrial fibrillation estimates were not dissimilar to genome-first estimates from the Geisinger Health System²³. Penetrance of *MYBPC3* and *MYH7* putatively pathogenic variants for hypertrophic cardiomyopathy was 7.0% [95%CI 4.1-11%] and 4.8% [95%CI 2.9-7.0%], respectively (**Supplementary Table 10**). *MYBPC3* LOFs were associated with 9.7% [95%CI 4.5,18%] penetrance. Relative-based analyses have frequently yielded estimates over 30% for sarcomere mutations²⁴⁻²⁷, although we note that our OR estimate for *MYBPC3* LOFs is very consistent with a previous case-control study²⁸. Similarly, family-member analyses have reported 40% incidence of arrhythmogenic cardiomyopathy/dysplasia for pathogenic desmosome mutations²⁹, yet we find that fewer than 5% of *PKP2* and *DSP* variant carriers have diagnoses of dilated cardiomyopathy or ventricular arrhythmia; fewer than 12% of carriers had atrial fibrillation (**Supplementary Tables 10 and 15**).

The penetrance estimates of *GCK* and *HNF1A* putatively pathogenic variants for type 2 diabetes were large at 64% [95%CI 49-78%] and 45% [95%CI 26-64%], respectively, with an age dependent penetrance (**Supplementary Table 10**). Previous studies have suggested that *HNF1A* mutations have over 90% penetrance for progressive diabetes at 50 years of age³⁰, while *GCK* mutations are thought to cause a shift in glucose-sensing and mild hyperglycemia from birth³¹. Interestingly, *GCK* LOF mutations - such as those contributing strongly to our signal - are found in MODY patients³², and rare *GCK* mutations are also enriched in individuals diagnosed with type 2 diabetes³³. Further population-based assessment seems warranted to determine diabetes-related outcomes, given the conventional knowledge that *GCK* mutations cause hyperglycemia that often does not require medical intervention³¹.

Penetrance estimates for significant associations at different cut-offs for age-at-onset are shown in **Supplementary Table 10**, showing an age-dependent probability of diagnosis for most gene-phenotype pairs. We acknowledge that these penetrance estimates are based on age-at-diagnosis, which may be inaccurate for diseases defined at UK Biobank visits. However, for age-specific penetrance estimates, we did not include cases defined at baseline for this reason; we further found that electronic health records were the most important source of data for many phenotypes (**Supplementary Table 2**). This should be considered when interpreting the age-stratified penetrance estimates, as true age-at-onset may be earlier than age-at-diagnosis based on ICD codes. In addition, by excluding cases defined by self-report at baseline, some individuals with early-onset disease may have been excluded for age-stratified analyses. Despite these limitations, these analyses highlight how age is an important factor in disease presentation in carriers of pathogenic variation.

Overall, our penetrance results highlight - from a genome-first perspective - substantially lower penetrance for pathogenic variation than previously reported from family-based analyses. This finding is consistent with previous analyses in the UK Biobank that utilized well-genotyped likely-pathogenic rare variants from the genotyping array³⁴. There are various factors that should be taken into account when interpreting population-based penetrance estimates. First, some survivor and ascertainment bias are to be expected in our relatively healthy middle-aged population-based cohort, which may bias penetrance estimates downwards. Furthermore, it is possible that certain putatively pathogenic variants included in our analysis are not truly pathogenic variants; for example, the 'likely pathogenic' variants from ClinVar may include some non-pathogenic alleles, and certain LOF variants may not be truly LOF. To mitigate these issues, we only included ClinVar variants reported from 2015 onwards (which should therefore conform to stringent guidelines for pathogenicity assertions) and we used LOFTEE to filter out as many low-confidence or dubious LOF variants as possible. Third, for many of the diseases, cases were defined primarily by ICD codes, which may downward bias penetrance estimates for diseases that can go undiagnosed or that are diagnosed outside of the hospital (e.g. diabetes, dyslipidemias, chronic kidney disease). Therefore, our estimates may reflect more severe symptomatic cases, while not including subclinical and mildly symptomatic disease. However, given strikingly

420 high penetrance estimates for pathogenic variation in hypercholesterolemia and diabetes genes (>60% for
421 *LDLR* and *GCK*), this effect generally appears not to be large. On the other hand, the high penetrance
422 estimates from family-member analyses are likely biased upwards. First, since many family-member based
423 analyses are based on clinically ascertained index cases with severe disease, such analyses are strongly
424 biased towards families prone to more severe disease and higher disease penetrance. Second, in-depth
425 phenotyping in such studies may over-diagnose disease even though clinical symptoms may never have
426 arisen. In sum, true penetrance estimates likely lie somewhere in between population-based estimates and
427 family-based/clinical cohort estimates.

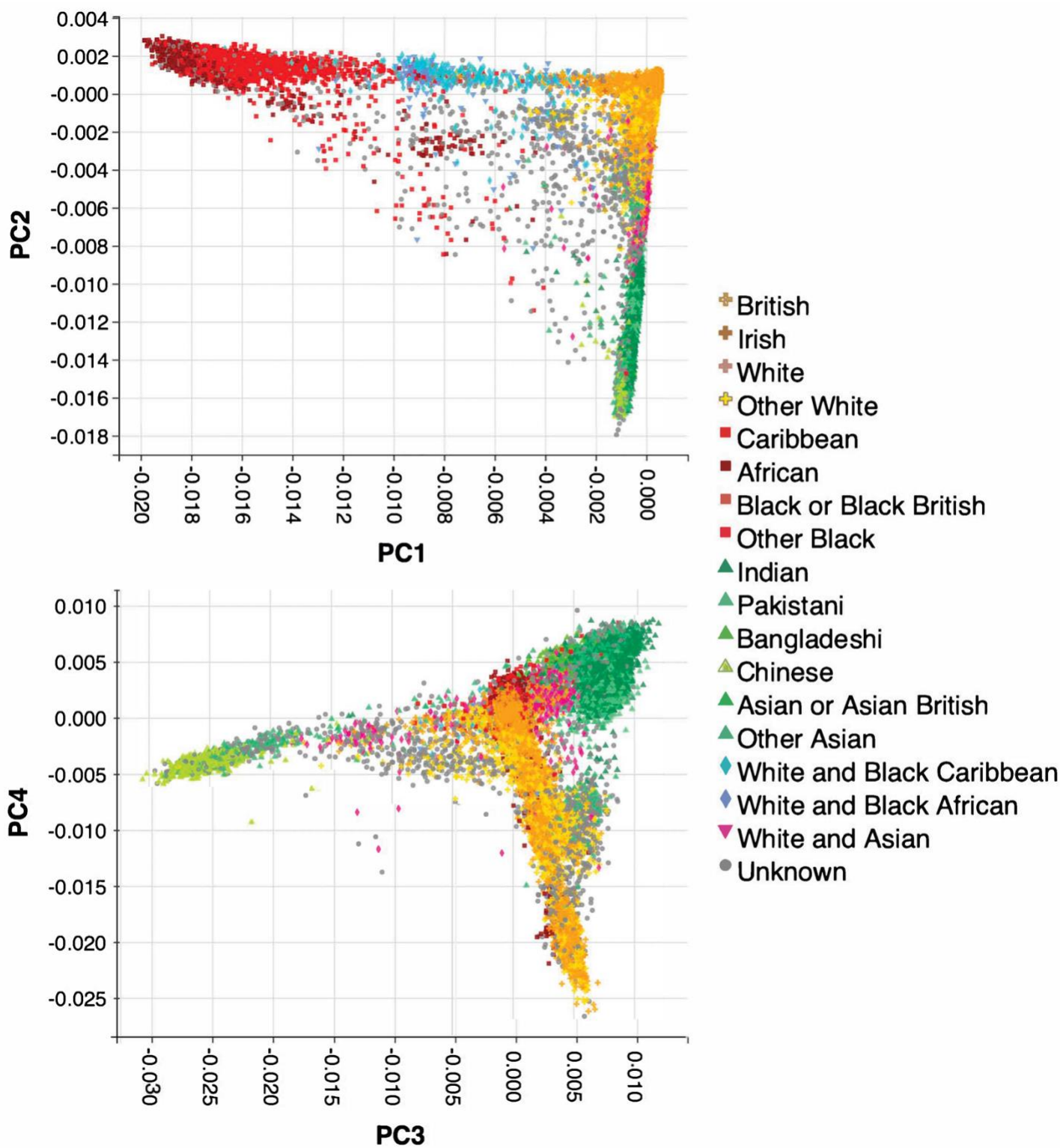
428 **Yield of putatively pathogenic variants among disease cases**

429 The yield of rare putatively pathogenic variants was generally low among disease cases. Among heart failure
430 and atrial fibrillation cases, the yield of associated pathogenic variants was ~1.5% and ~1.1%, respectively. For
431 hypertrophic and dilated cardiomyopathy, the yield of associated putatively pathogenic variants was ~11.0%
432 and ~10.1%, respectively (**Supplementary Table 17 and Supplementary Figure 14**). It should be noted that
433 rare variant yield in this case only represents the yield of LOFs, known likely pathogenic variants and known
434 pathogenic variants in genes showing evidence of association at $P < 0.005$. As such, this yield is a conservative
435 lower-bound estimate that should rise as more genes are included and more non-truncating pathogenic
436 variants are discovered.
437

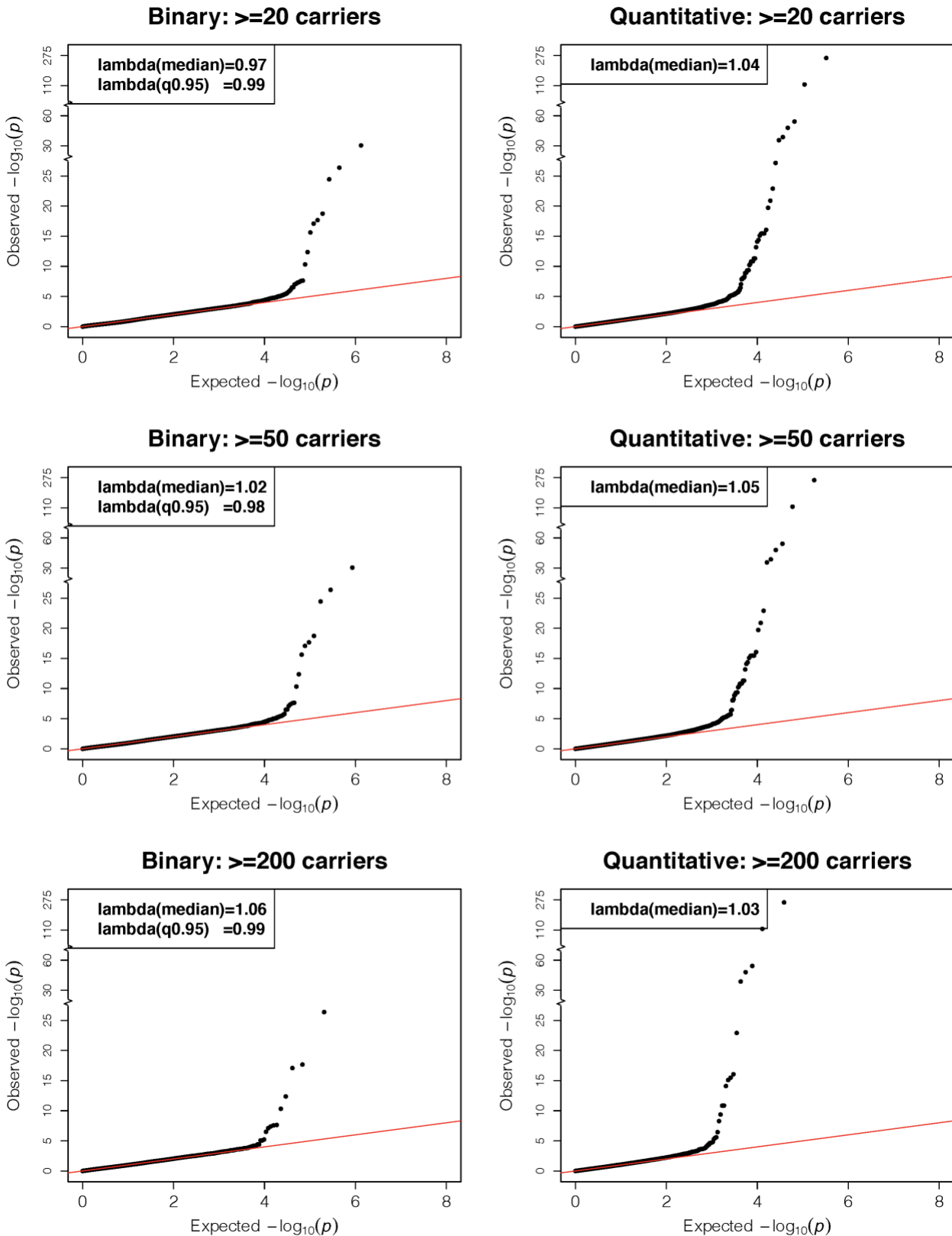
438 **Supplementary Tables**

439

440 **Supplementary Tables 1-17** can be found in the **Supplementary Excel File**.



442
 443
 444 **Supplementary Figure 1. Principal component analysis and self-reported ancestries for UK Biobank**
 445 **WES samples.** Samples are plotted for PC1-4, and self-reported ancestries are highlighted. The principal
 446 components stratify samples from different major ancestral groups.

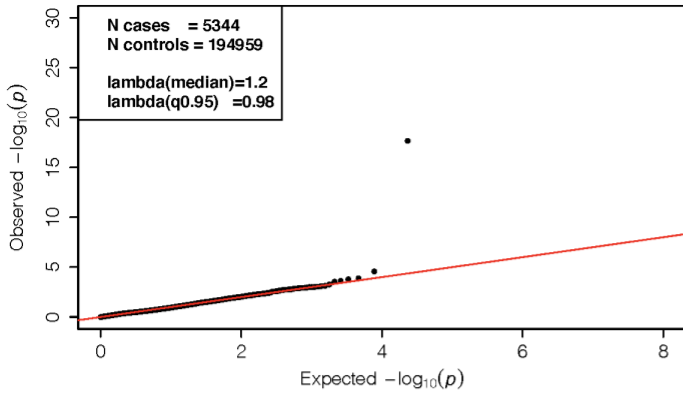


447
 448
 449
 450
 451
 452

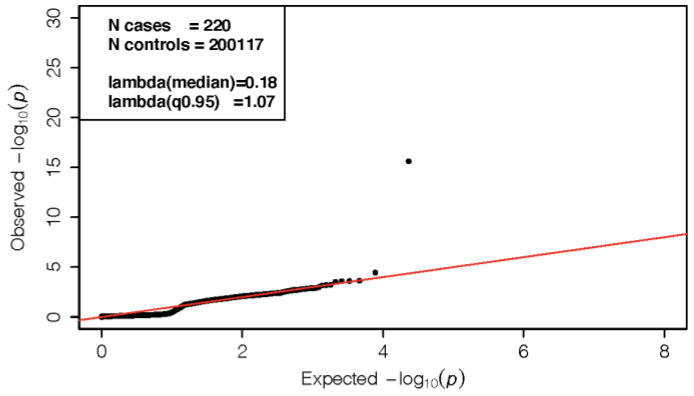
Supplementary Figure 2: Quantile-quantile plots for exome-wide gene-based tests across all binary and all quantitative phenotypes. The y-axis represents the observed $-\log_{10} P$ -values across all tests, while the x-axis represents the expected under the null-hypothesis. P -values were obtained from score tests in linear mixed effects models (quantitative traits) or saddle point approximation in logistic mixed effects models (binary traits), adjusting for sex, age, sequencing batch, MRI serial number (for MRI traits), associated principal

453 components (PCs) and a sparse kinship matrix. *P*-values shown are two-sided and unadjusted for multiple
454 testing, The left panels show the results for binary traits (for tests with ≥ 20 carriers, 50 carriers and 200
455 carriers), while the right panels show the results for quantitative traits. Across all performed tests, no systemic
456 inflation is observed.

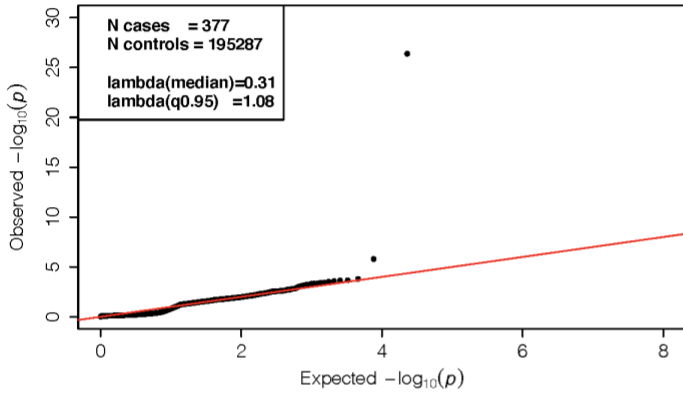
Heart failure



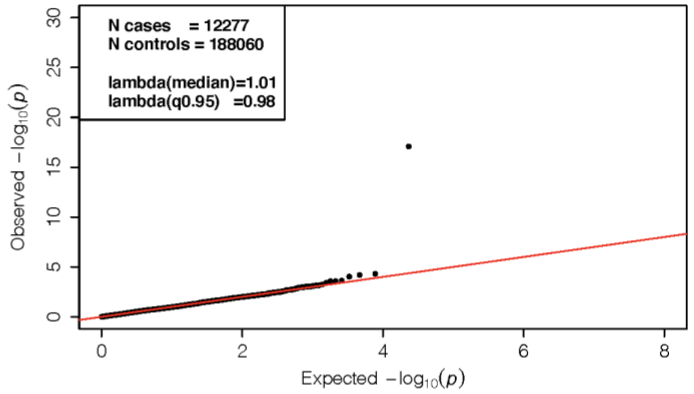
Hypertrophic cardiomyopathy



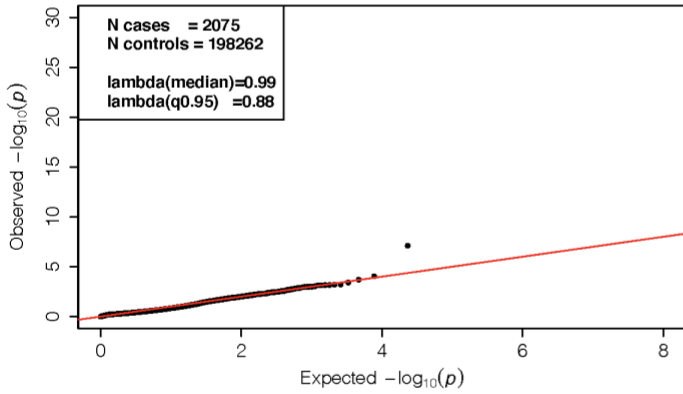
Dilated cardiomyopathy



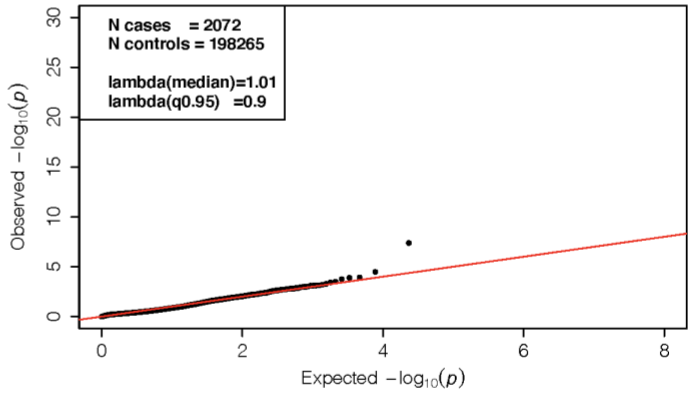
Atrial fibrillation



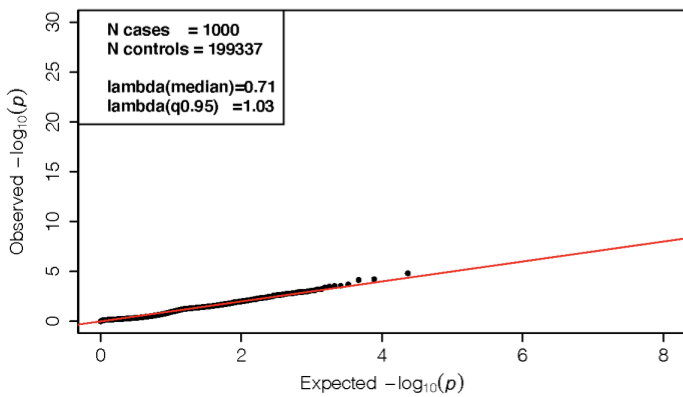
Supraventricular tachycardia



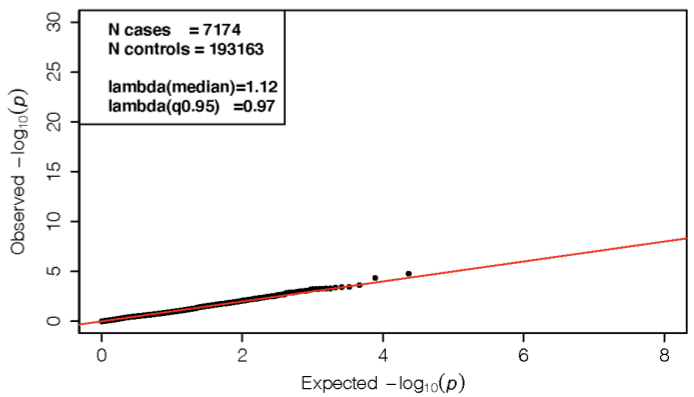
Ventricular arrhythmia



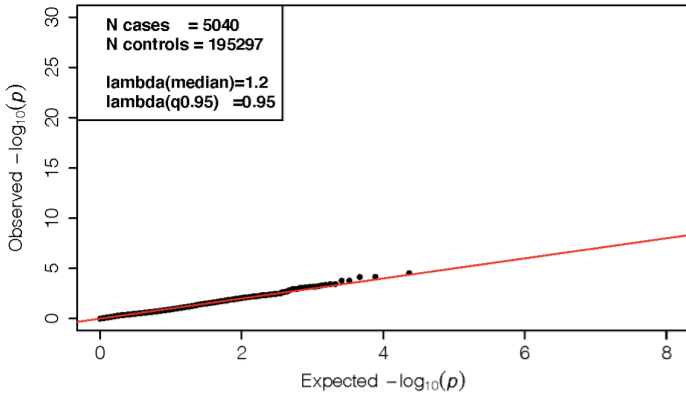
Cardiac arrest



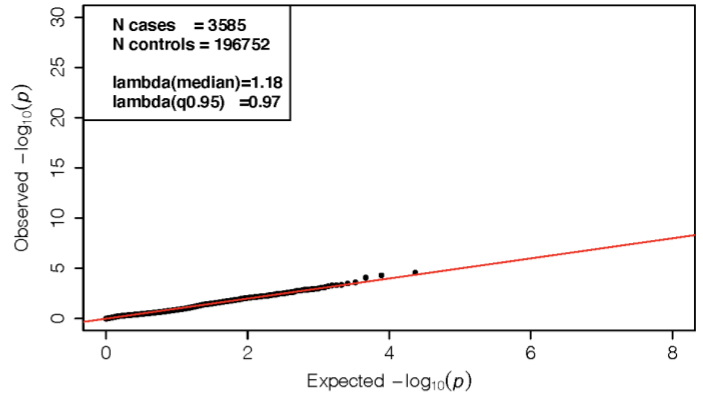
Bradycardia



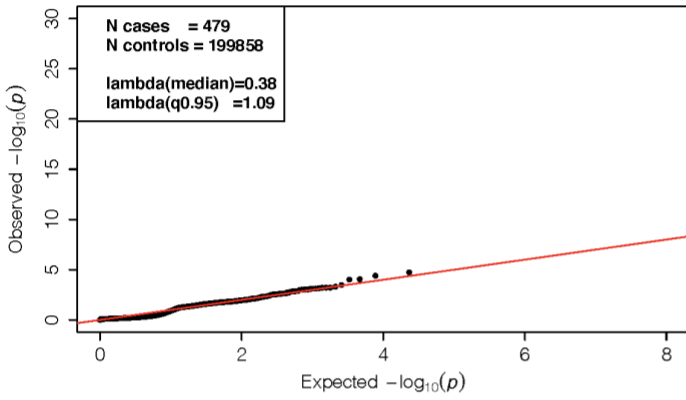
AV or bundle branch block



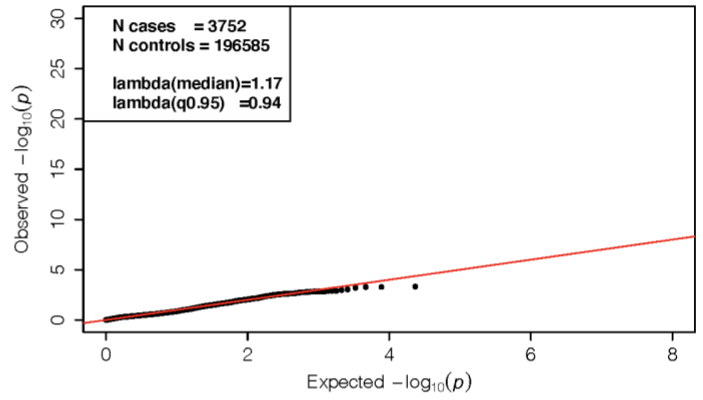
Pacemaker



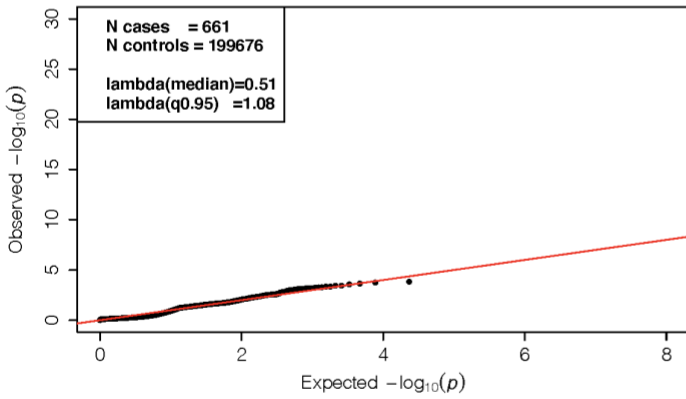
Implantable cardioverter defibrillator



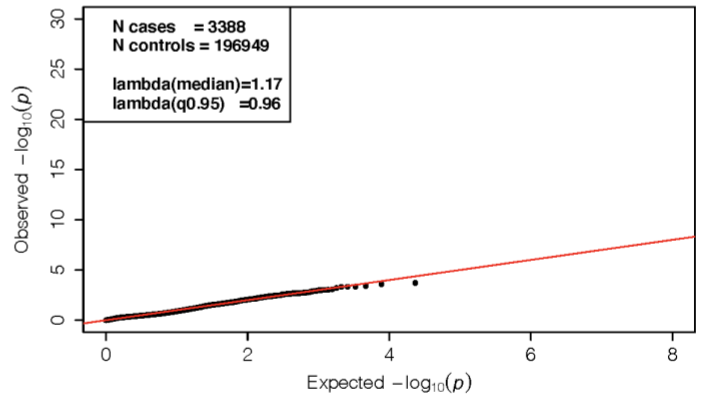
Heart surgery



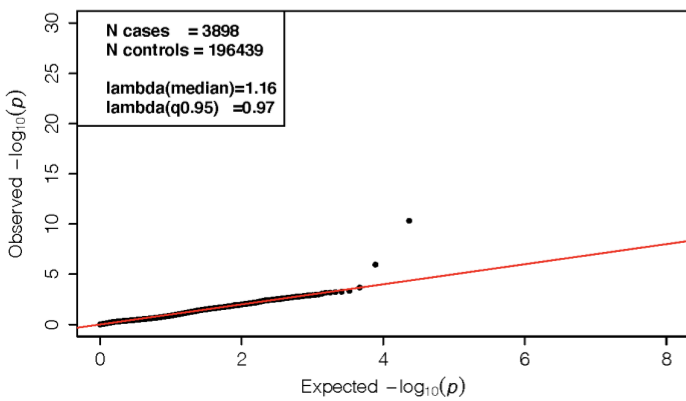
Congenital heart disease



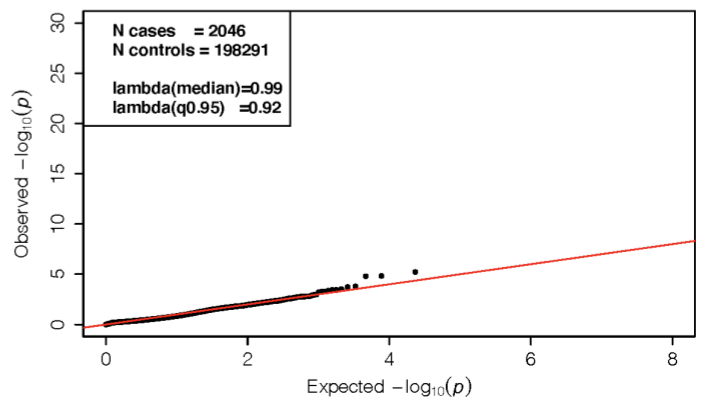
Aortic valve disease



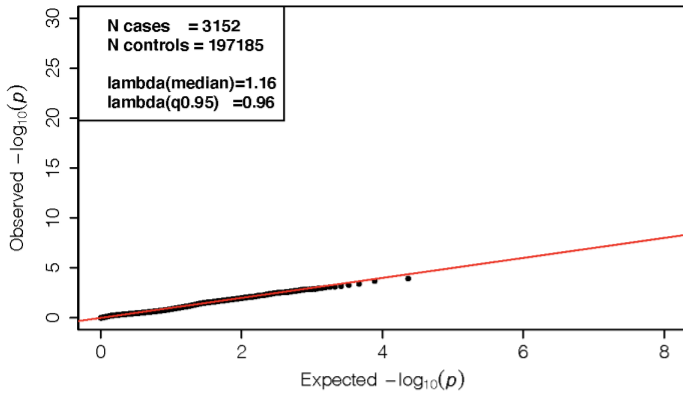
Mitral valve disease



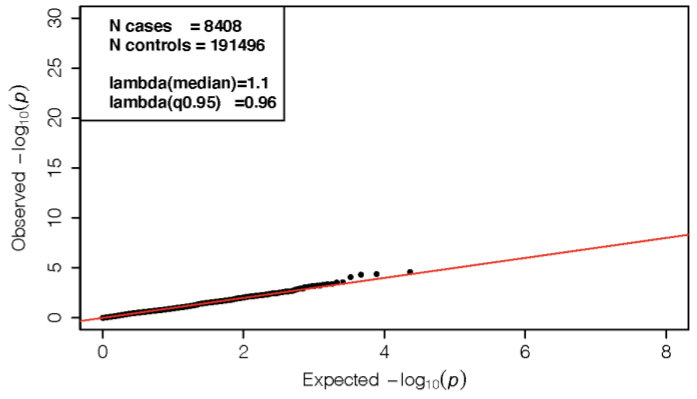
Tricuspid valve disease



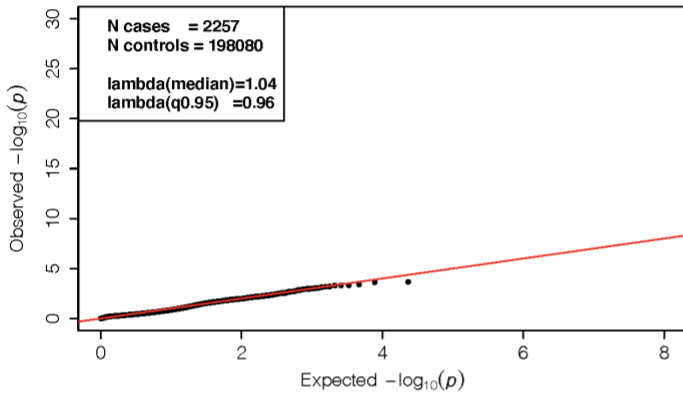
Peripheral vascular disease



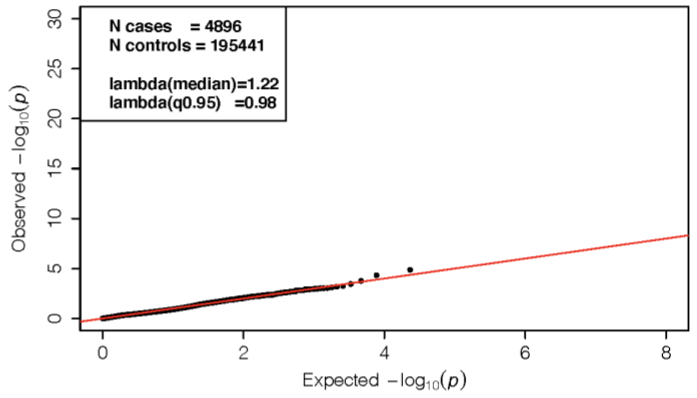
Venous thromboembolism



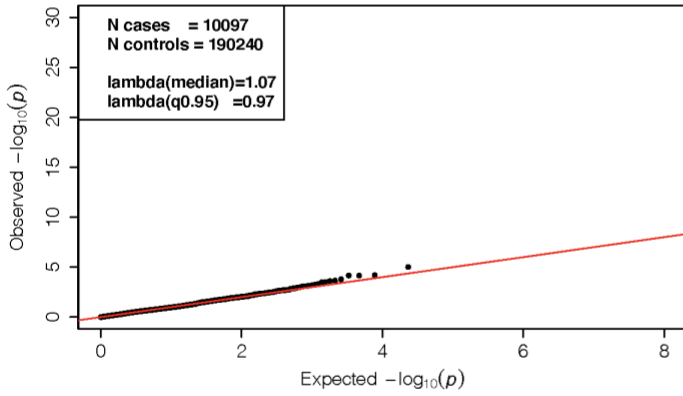
Ischemic stroke



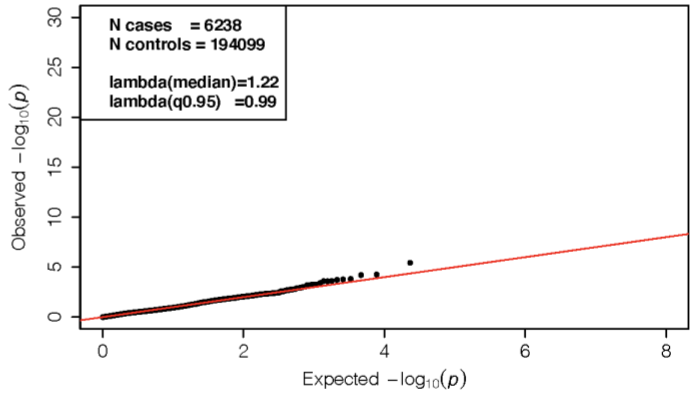
Stroke



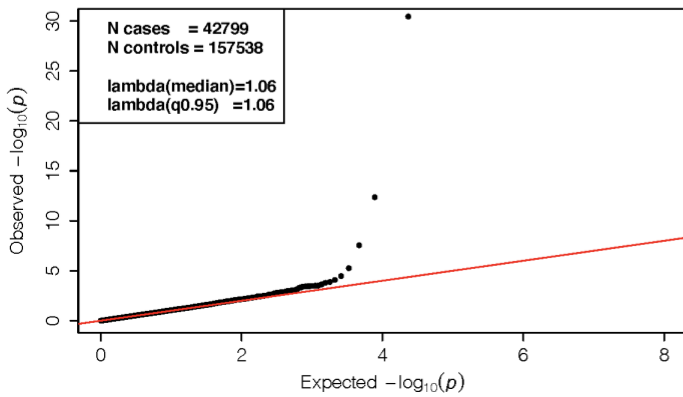
Coronary artery disease



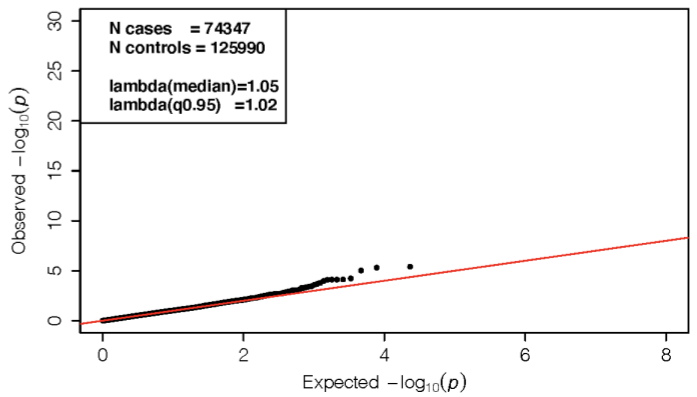
Myocardial Infarction



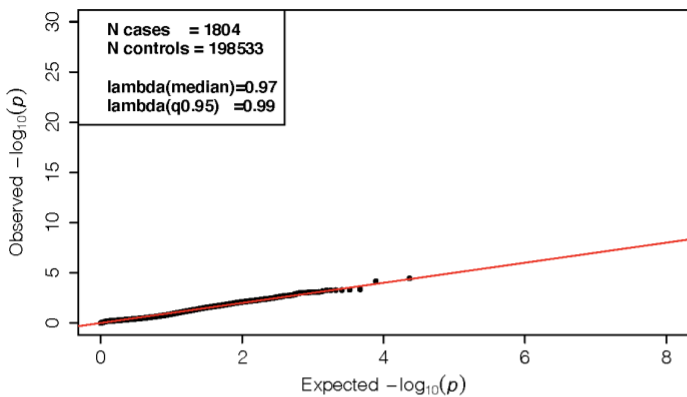
Hypercholesterolemia



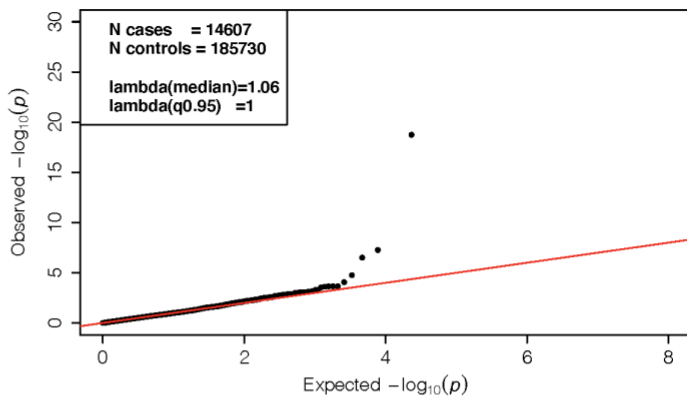
Hypertension



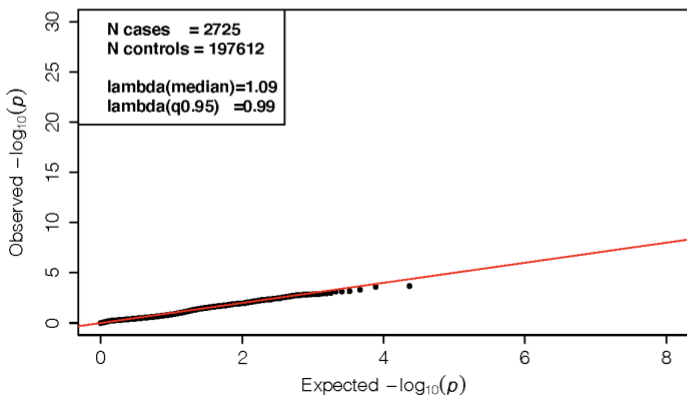
Diabetes type 1



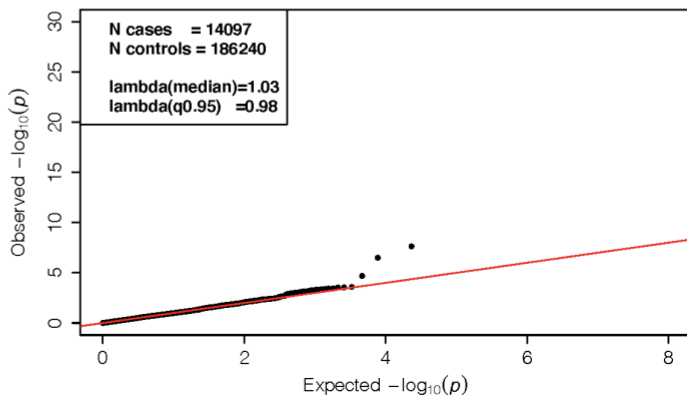
Diabetes type 2



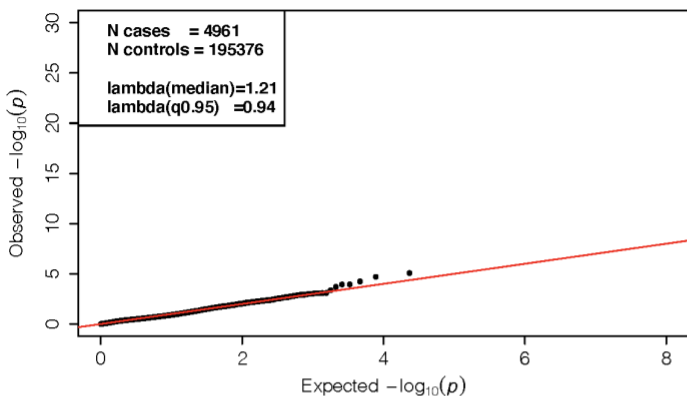
Hyperthyroidism



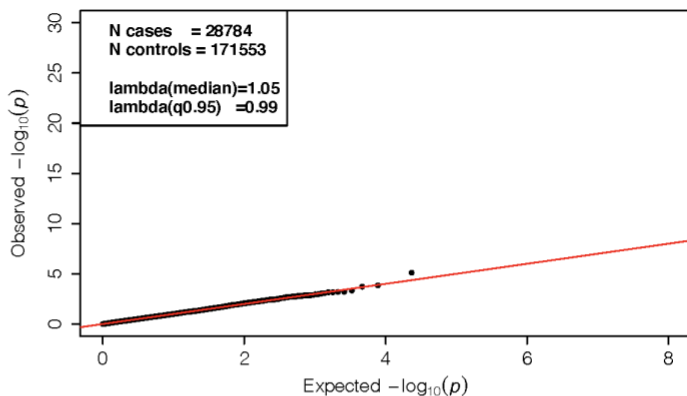
Hypothyroidism



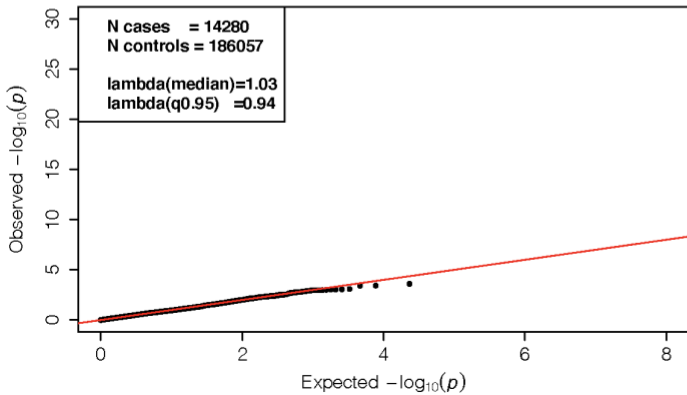
Gout



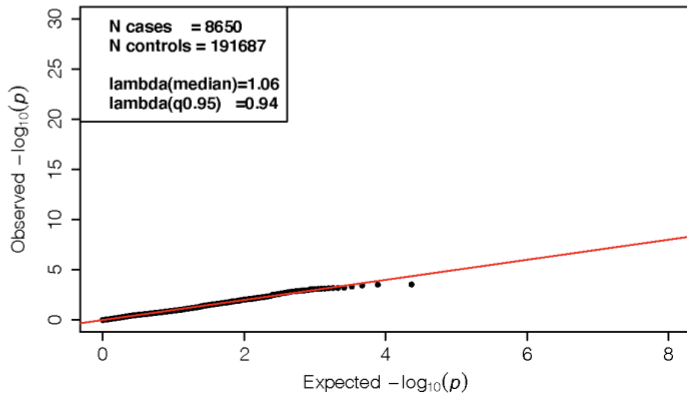
Asthma



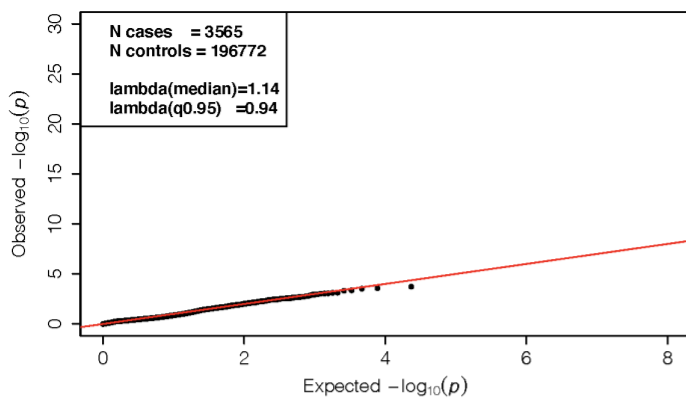
Allergic rhinitis



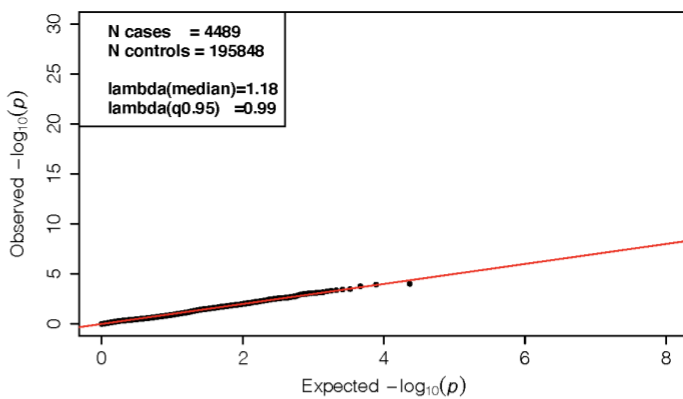
Dermatitis



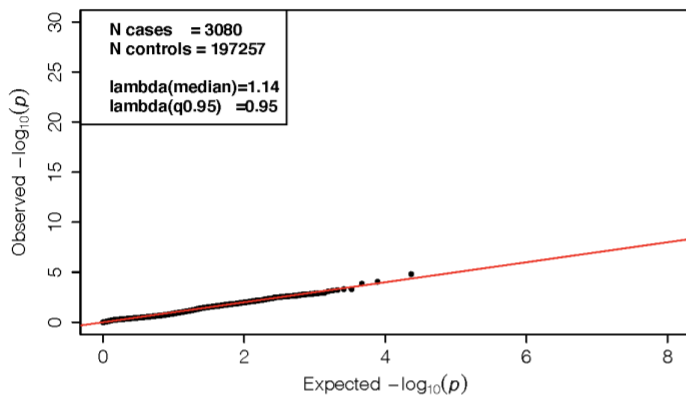
Psoriasis



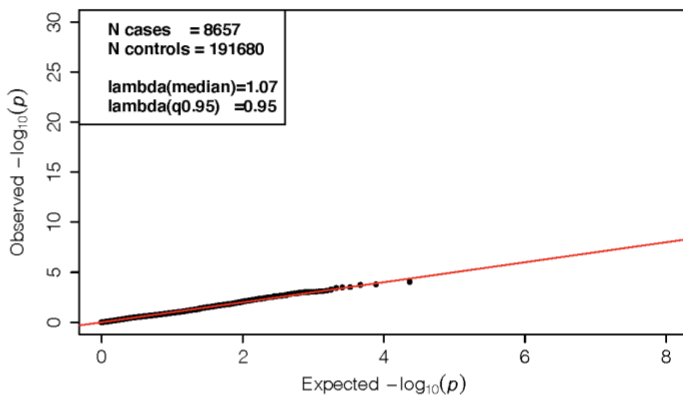
Rheumatoid arthritis



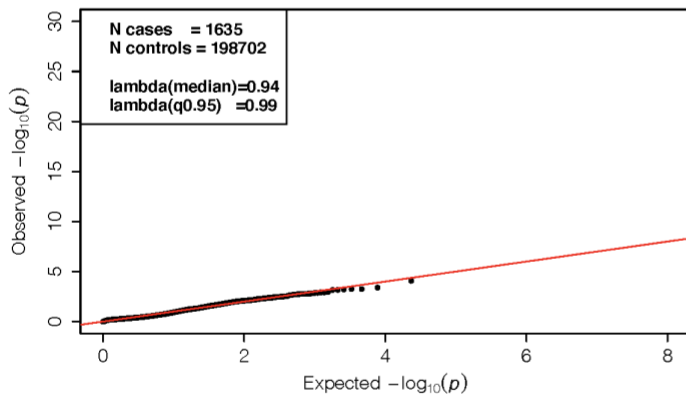
Inflammatory bowel disease



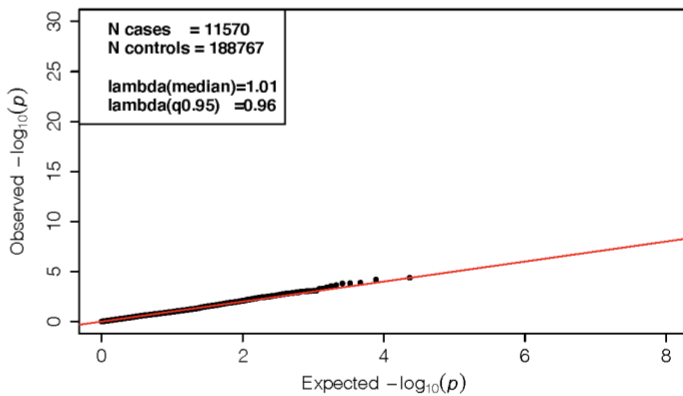
Irritable bowel syndrome



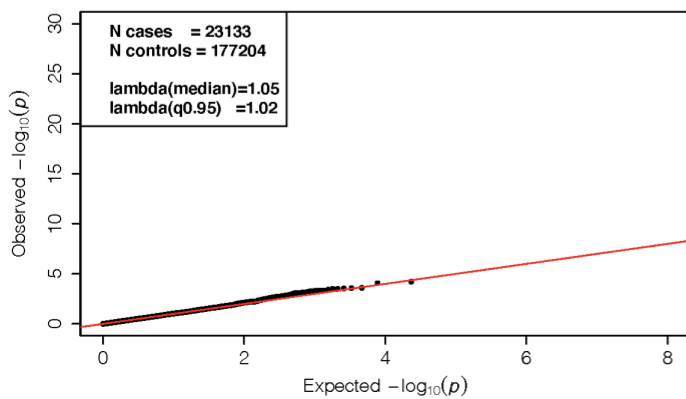
Pancreatitis



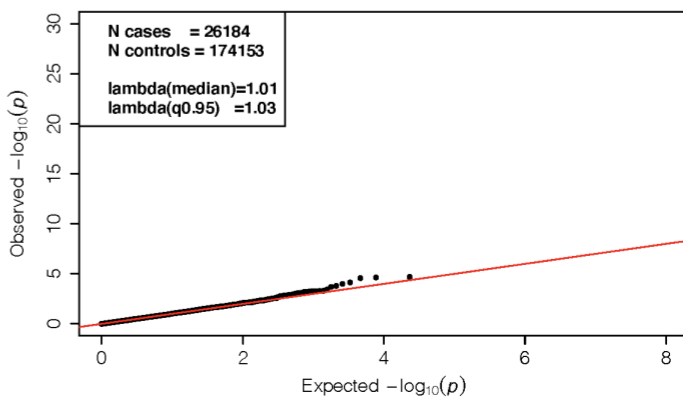
Gallstones



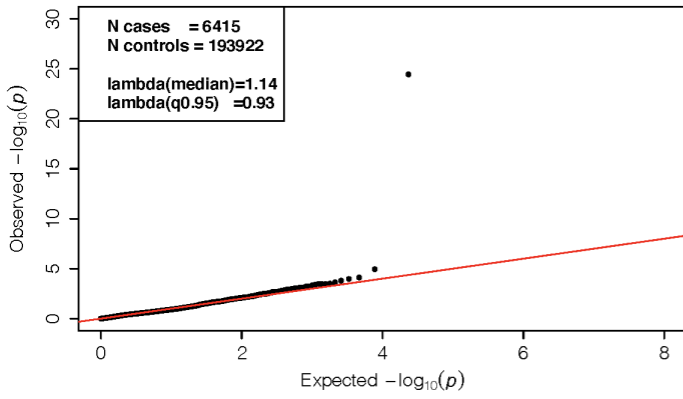
Diverticular disease



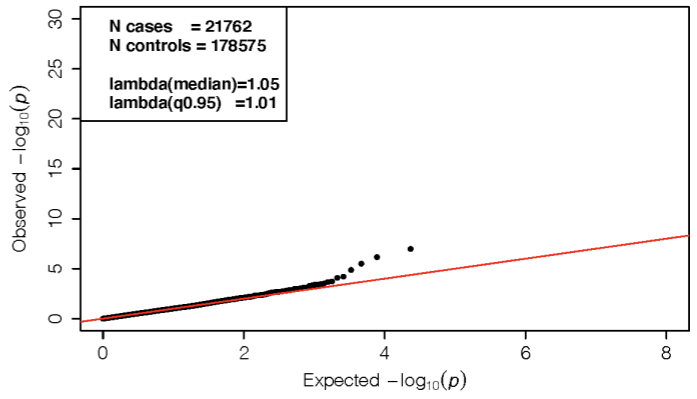
Gastroesophageal reflux disease



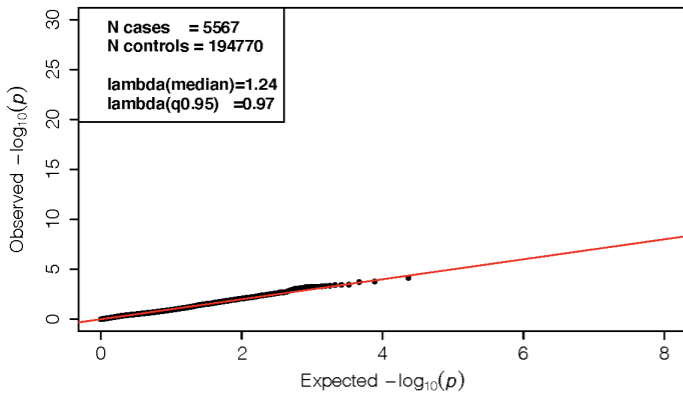
Chronic kidney disease



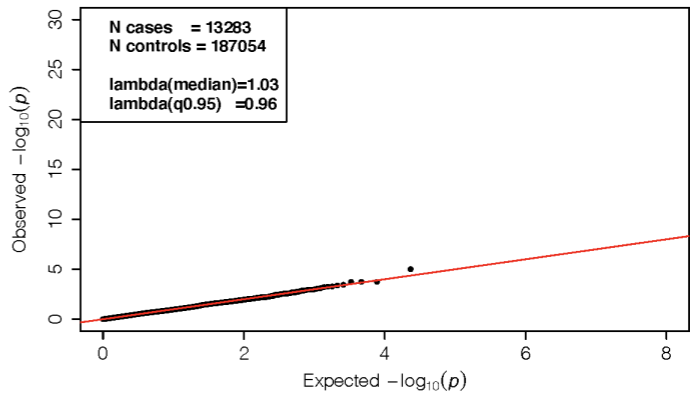
Cataract



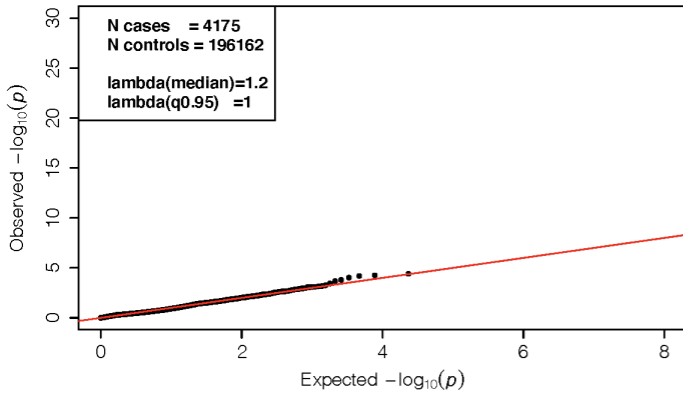
Glaucoma



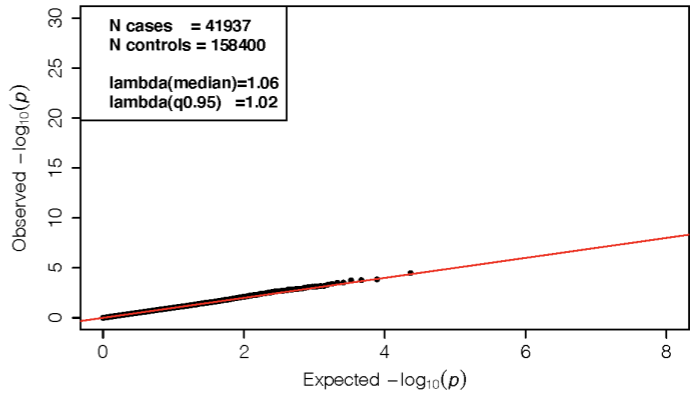
Back pain



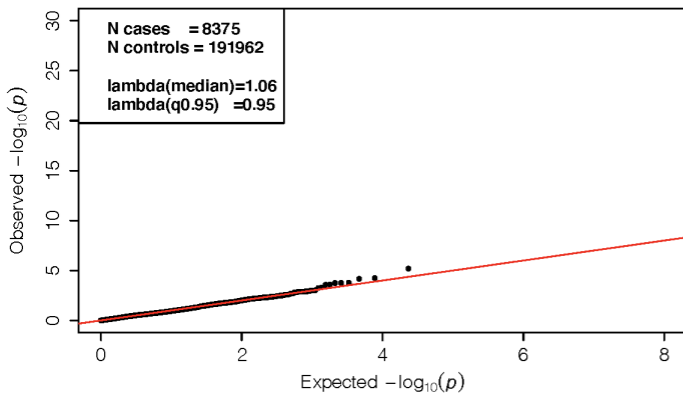
Sciatica



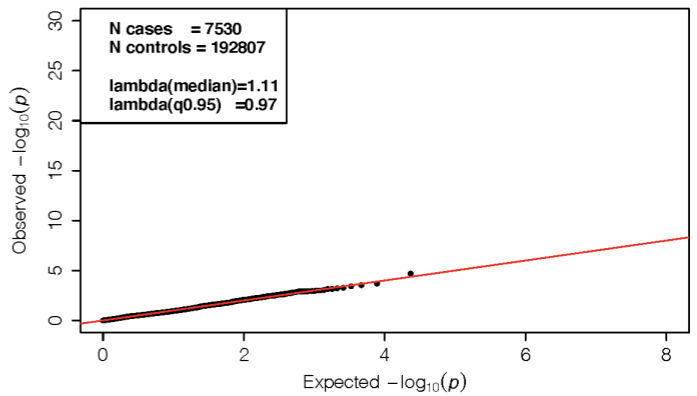
Osteoarthritis



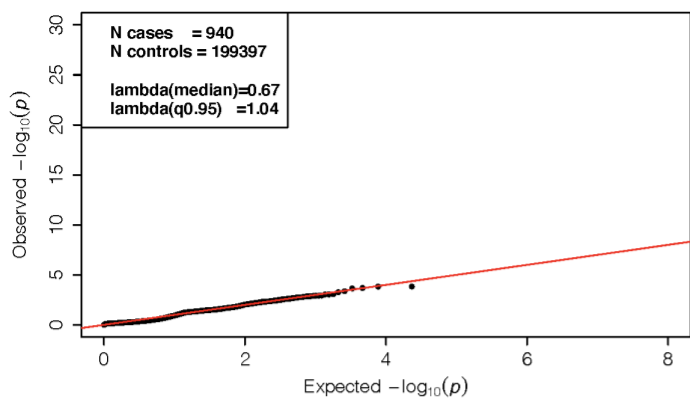
Osteoporosis



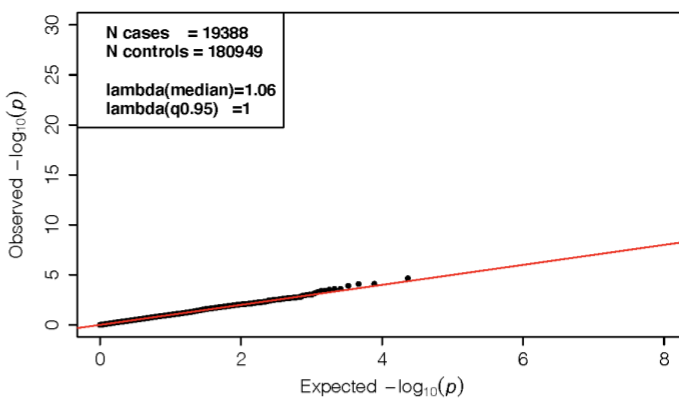
Anxiety



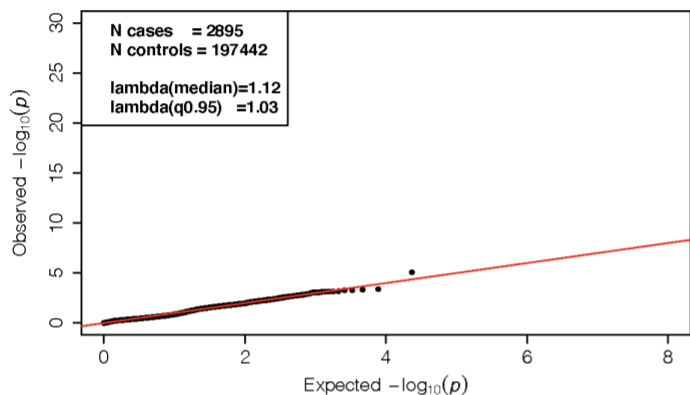
Bipolar disorder



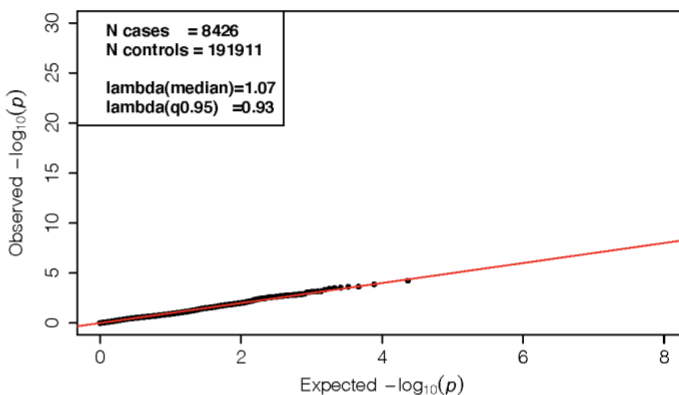
Depression



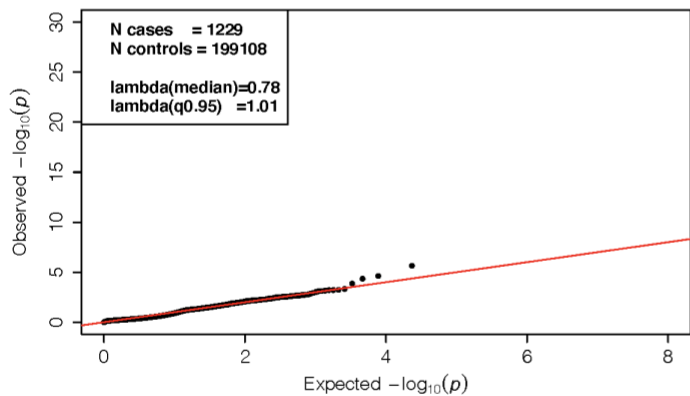
Epilepsy



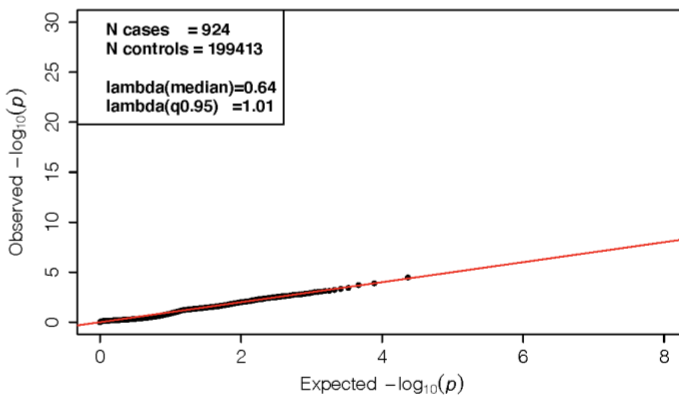
Migraine



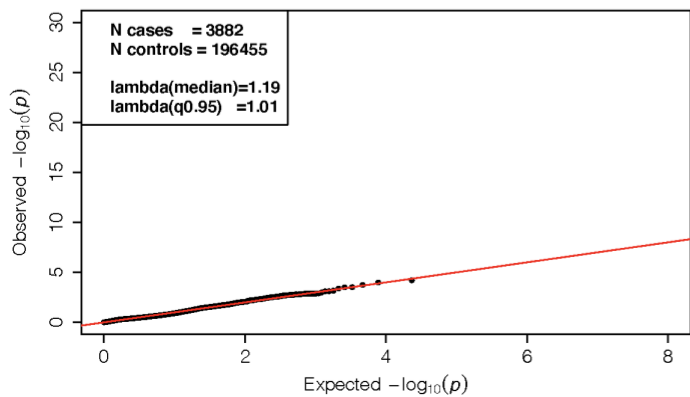
Parkinson's disease



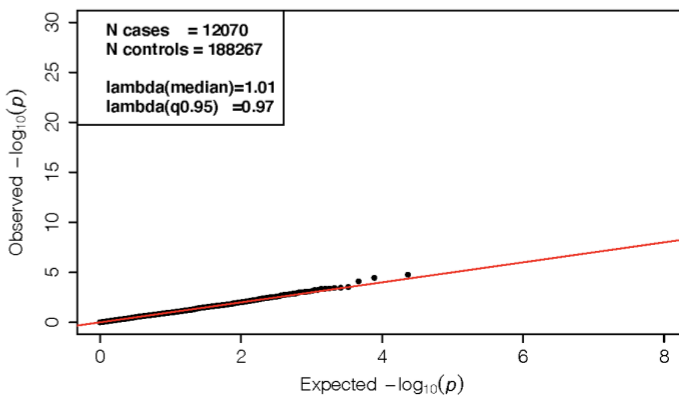
Multiple sclerosis



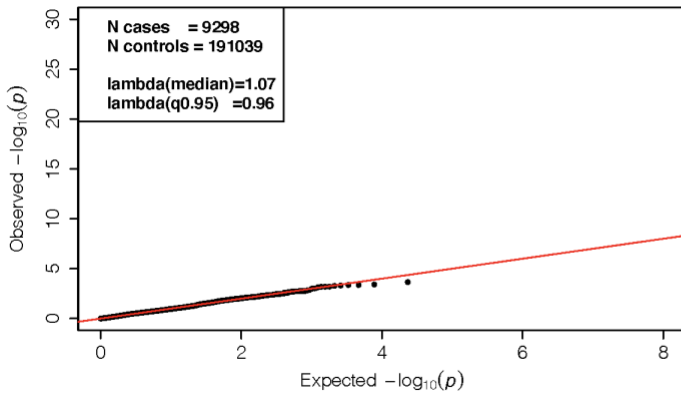
Sleep apnea



Pneumonia



COPD



464

465

466

467

468

469

470

471

472

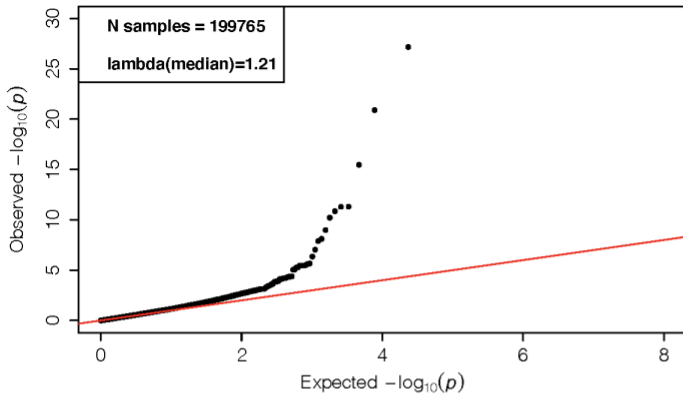
473

474

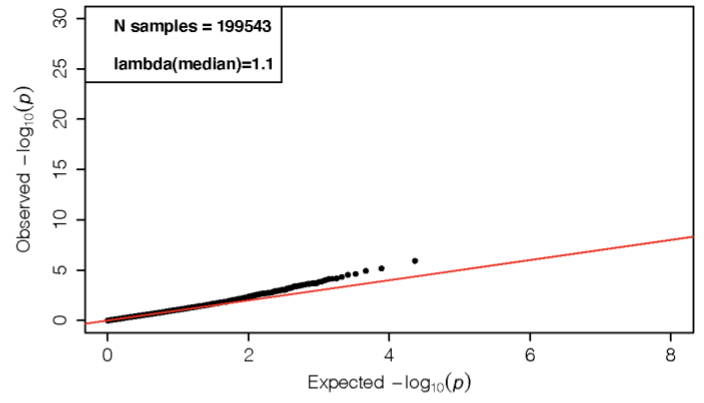
475

Supplementary Figure 3: Quantile-quantile plots for exome-wide gene-based tests for each individual binary trait. The y-axis represents the observed $-\log_{10} P$ -values across all tests, while the x-axis represents the expected under the null-hypothesis. P -values were obtained from saddle point approximation and were obtained from logistic mixed effects models, adjusting for sex, age, sequencing batch, associated principal components (PCs), a sparse kinship matrix. P -values shown are two-sided and unadjusted for multiple testing. The algorithm implemented in GENESIS applies the Saddle Point Approximation to raw P -values reaching $P < 0.05$ to account for case-control imbalance. As such, P -values larger than 0.05 might not be well calibrated and λ estimated at the median of the P -value distribution might not capture the calibration of the tests. We therefore estimated the λ values for binary traits at the tail of the distributions, $\lambda(q0.95)$, by comparing the test statistic at the 95th quantile to the expected statistic at this quantile under the null. Visually, and by judging $\lambda(q0.95)$, none of the binary traits showed major systemic inflation.

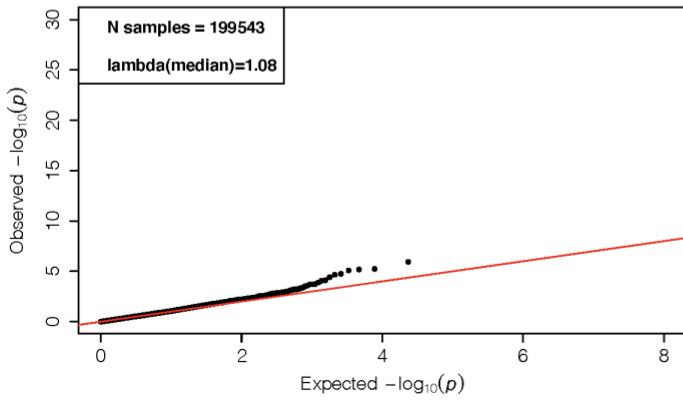
Height



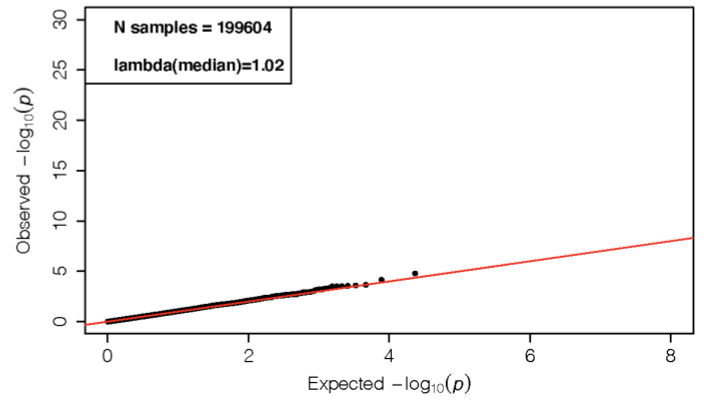
Weight



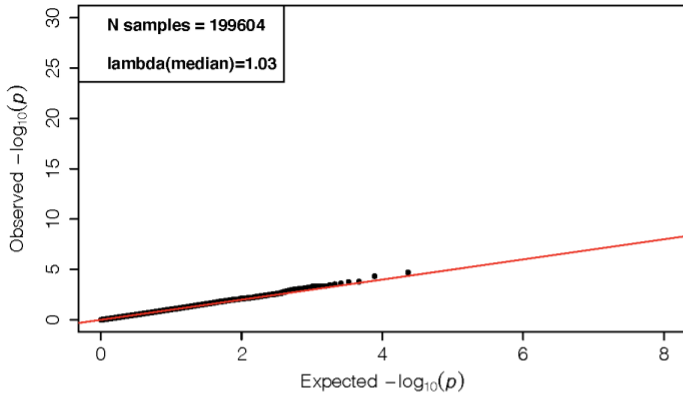
BMI



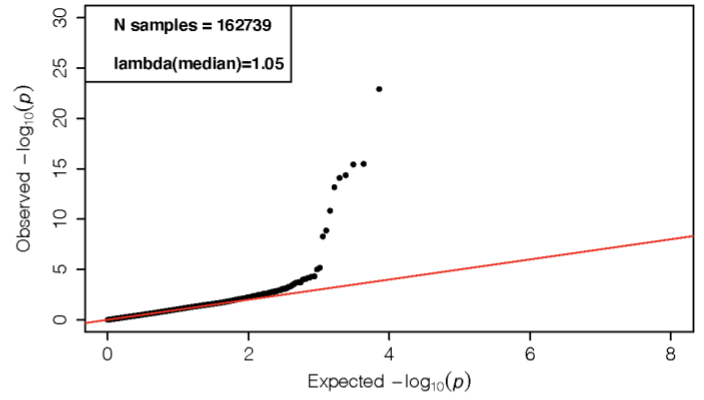
Systolic blood pressure



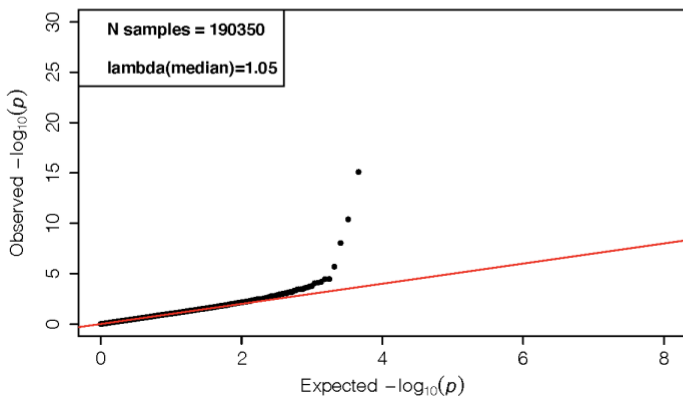
Diastolic blood pressure



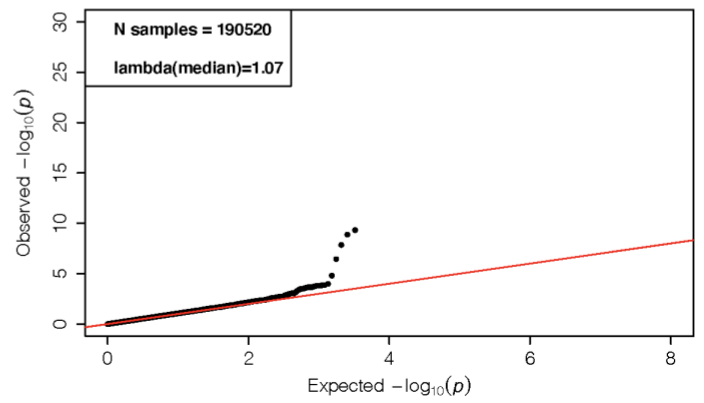
HDL



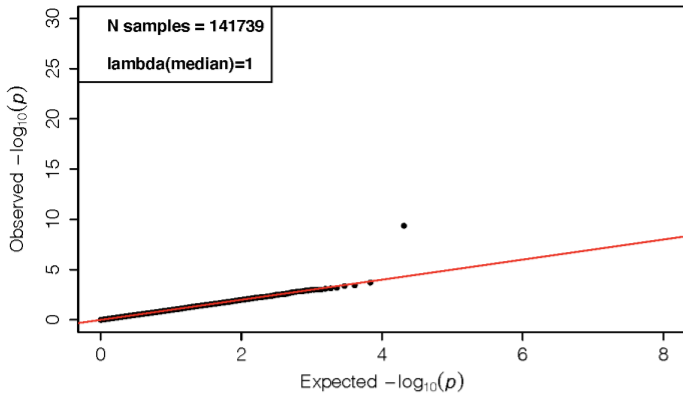
LDL



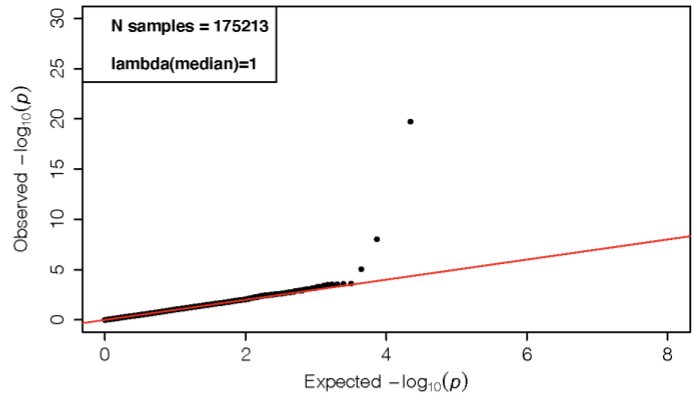
Triglycerides



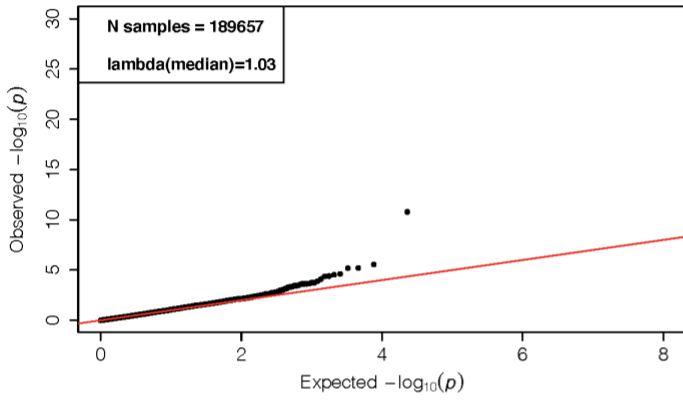
Lipoprotein(a)



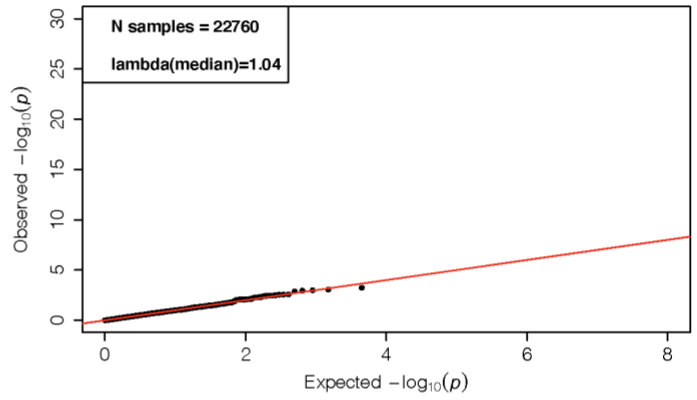
Glucose



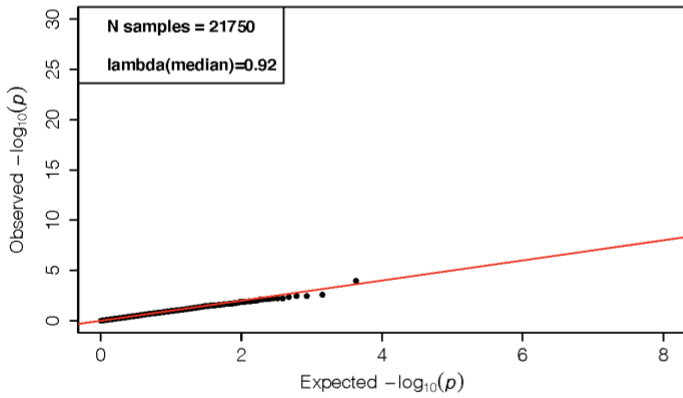
Igf-1



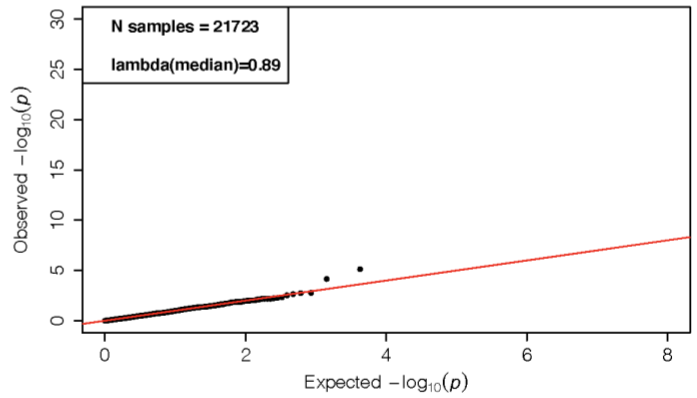
RR



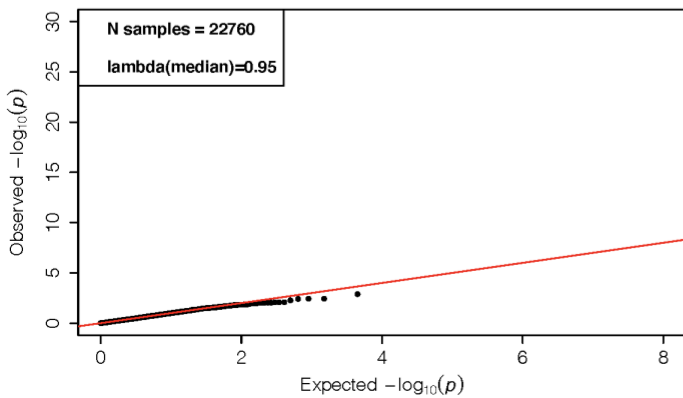
Pdur



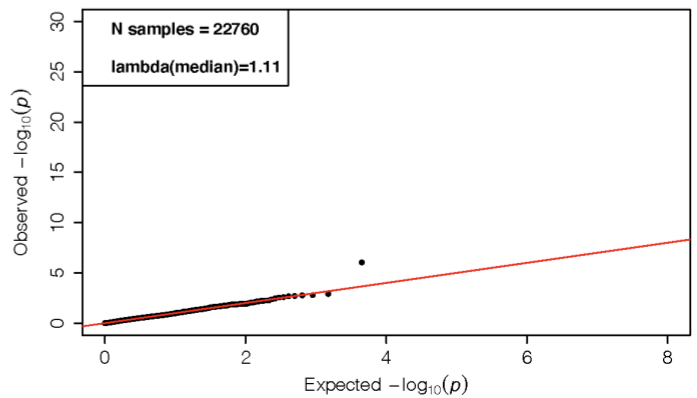
PQ



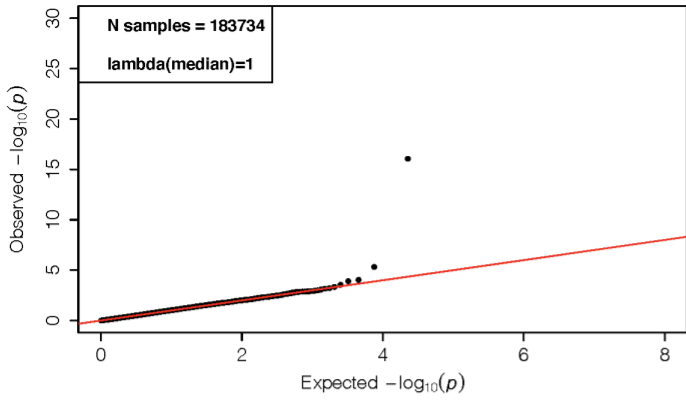
QRS



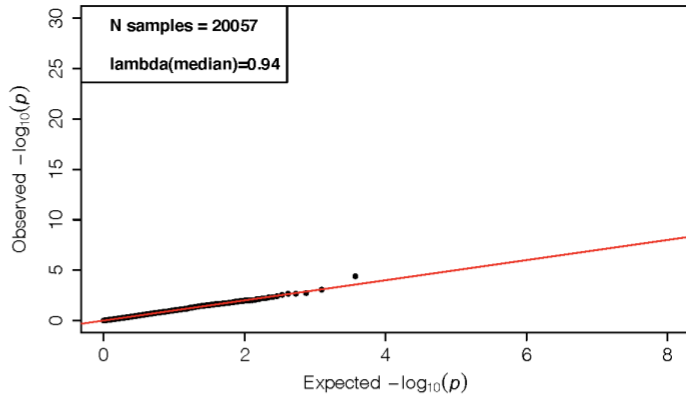
QTc



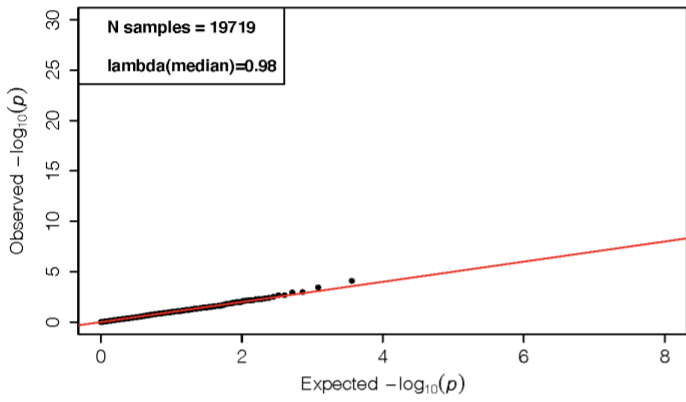
Pulse Rate



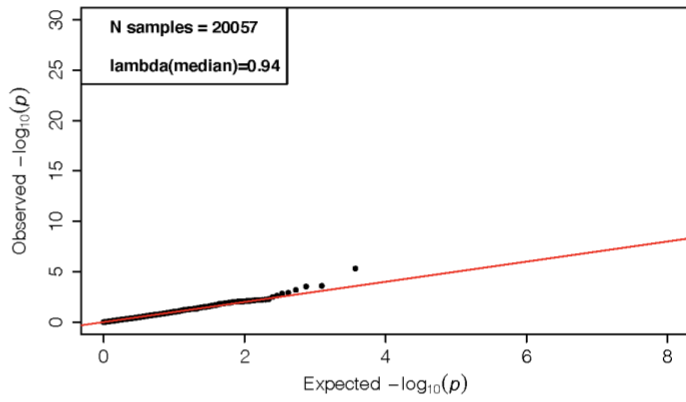
LVEDV



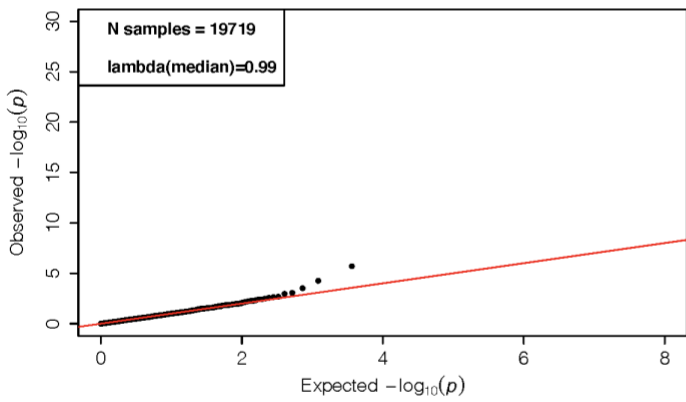
LVEDVi



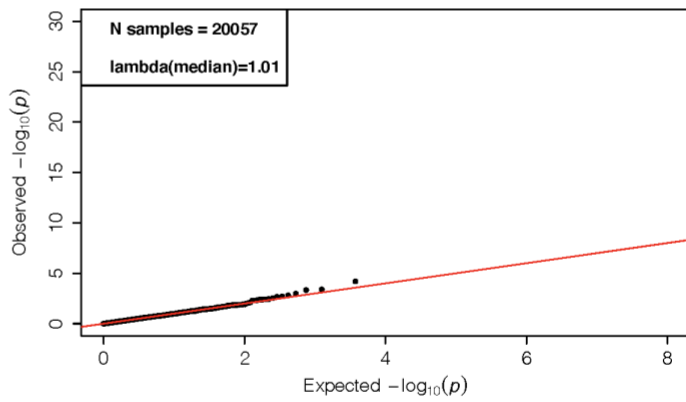
LVESV



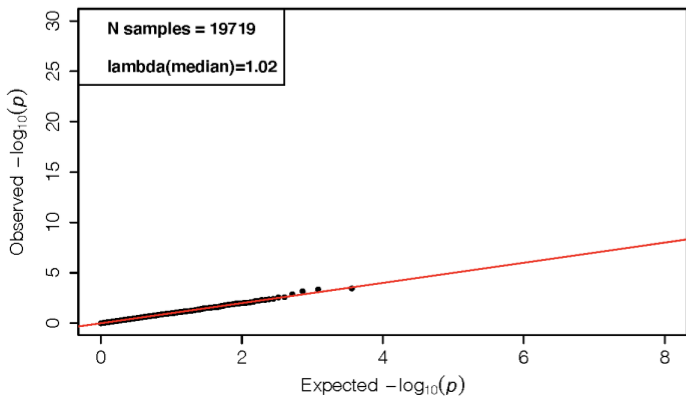
LVESVi



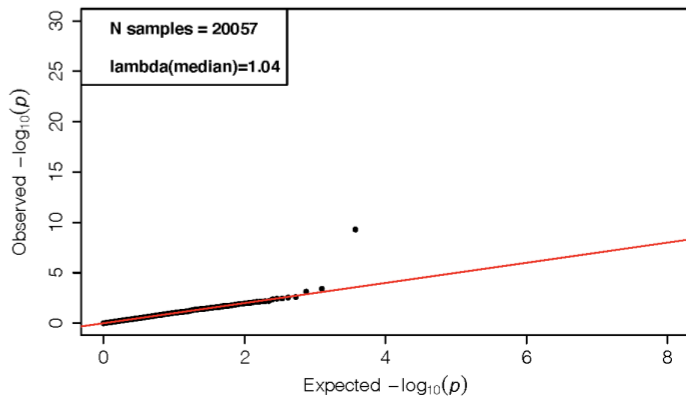
SV

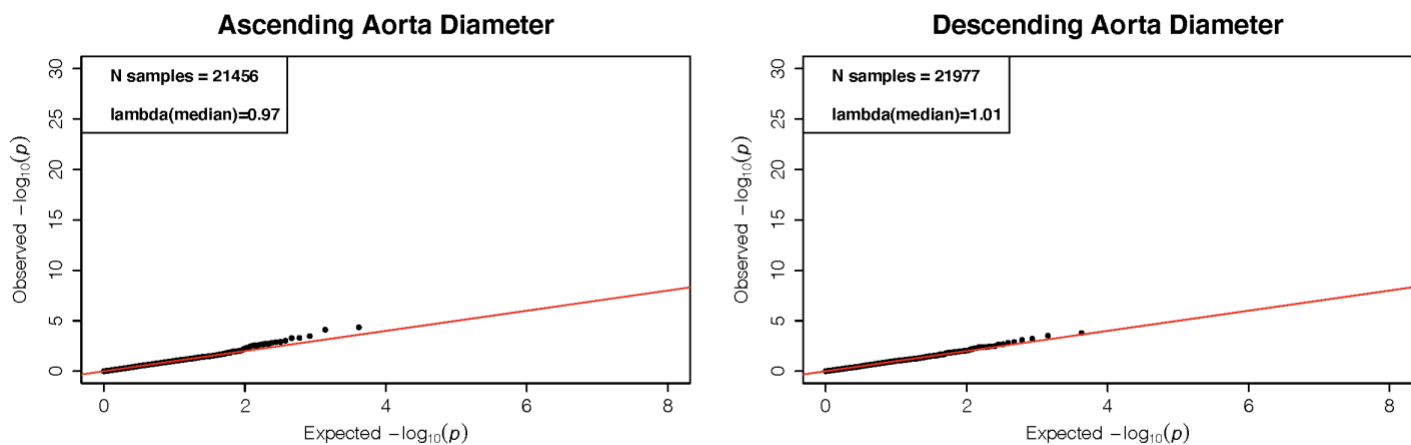


SVi



LVEF

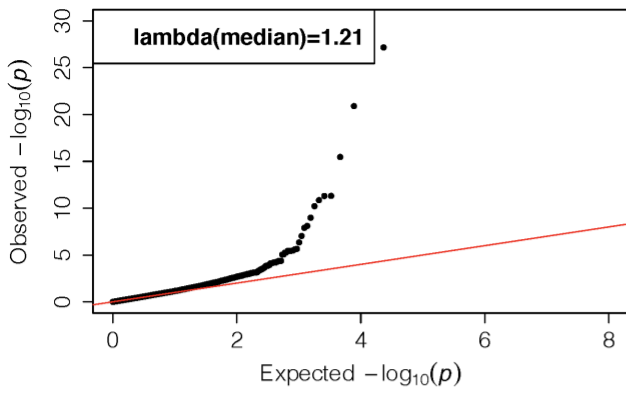




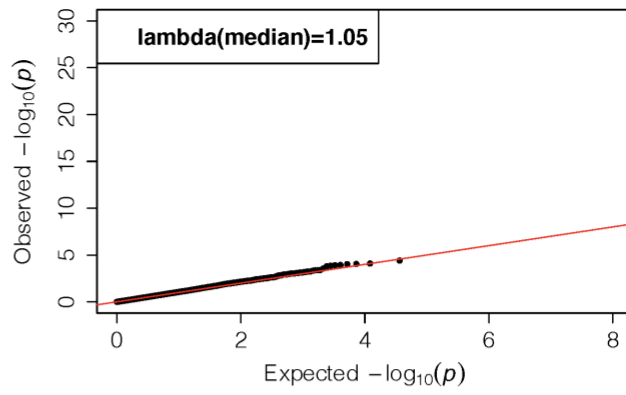
479
480
481
482
483
484
485
486

Supplementary Figure 4: Quantile-quantile plots for exome-wide gene-based tests for each individual quantitative trait. The y-axis represents the observed $-\log_{10} P$ -values across all tests, while the x-axis represents the expected under the null-hypothesis. P -values were obtained from score tests in linear mixed effects models, adjusting for sex, age, sequencing batch, associated principal components (PCs), MRI serial number (for MRI traits) and a sparse kinship matrix. P -values shown are two-sided and unadjusted for multiple testing. Values of λ were estimated at the median of the test statistic distribution. For three traits, height, weight and QTc, $\lambda(\text{median})$ was larger than 1.1. Height and weight indeed show a visual inflation.

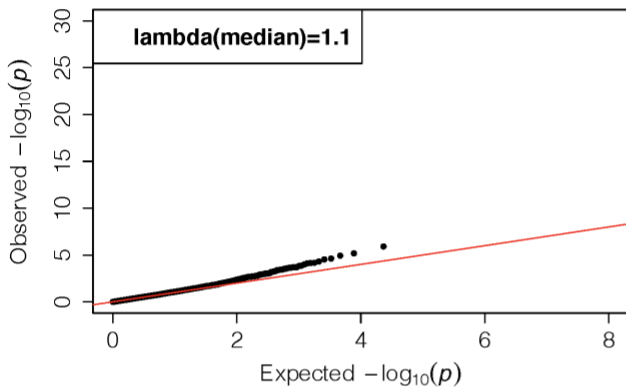
Height: LOF+Missense0.9



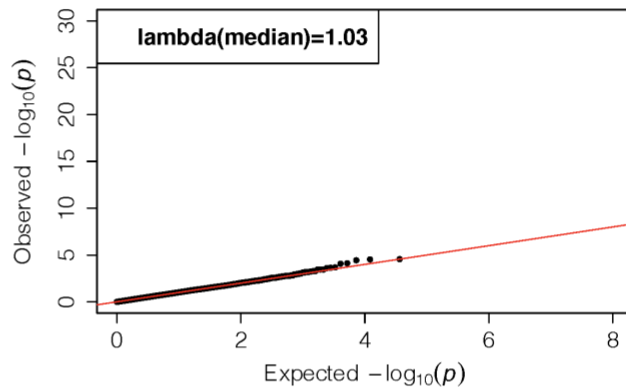
Height: synonymous



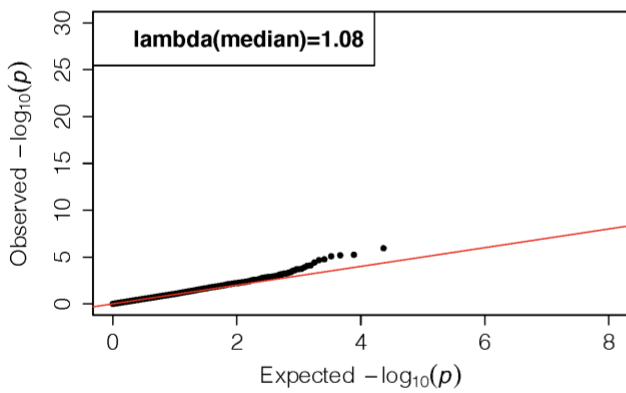
Weight: LOF+Missense0.9



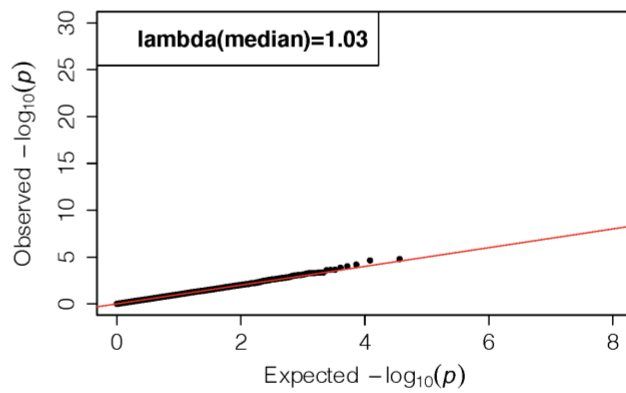
Weight: synonymous



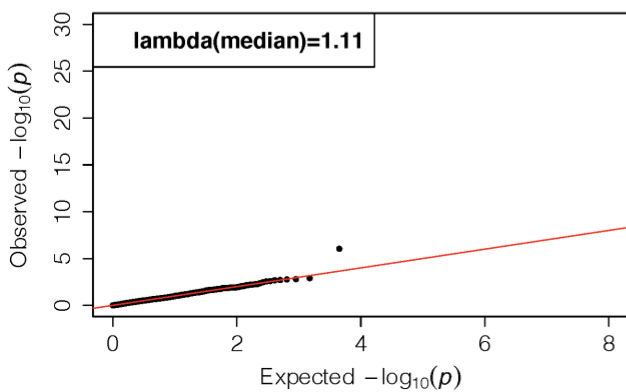
BMI: LOF+Missense0.9



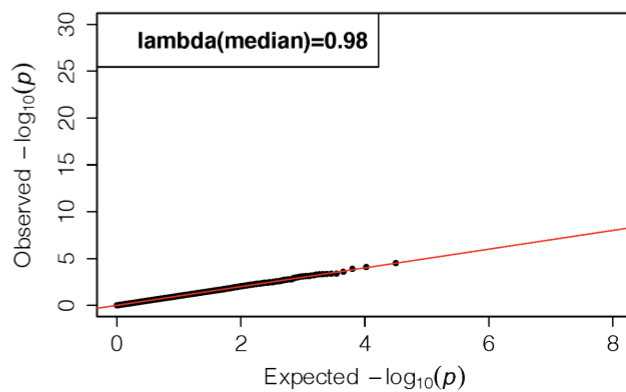
BMI: synonymous



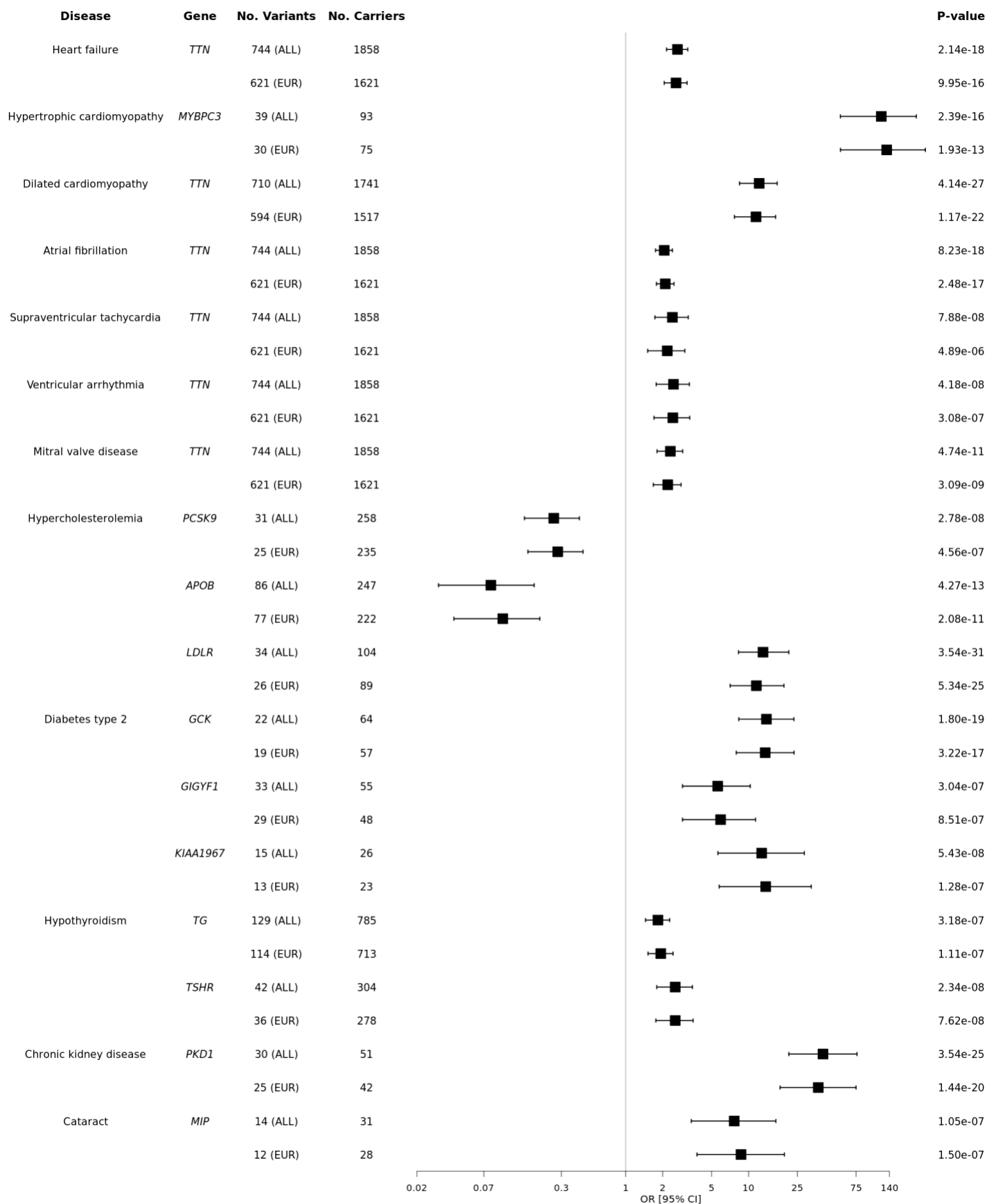
QTc: LOF+Missense0.9



QTc: synonymous

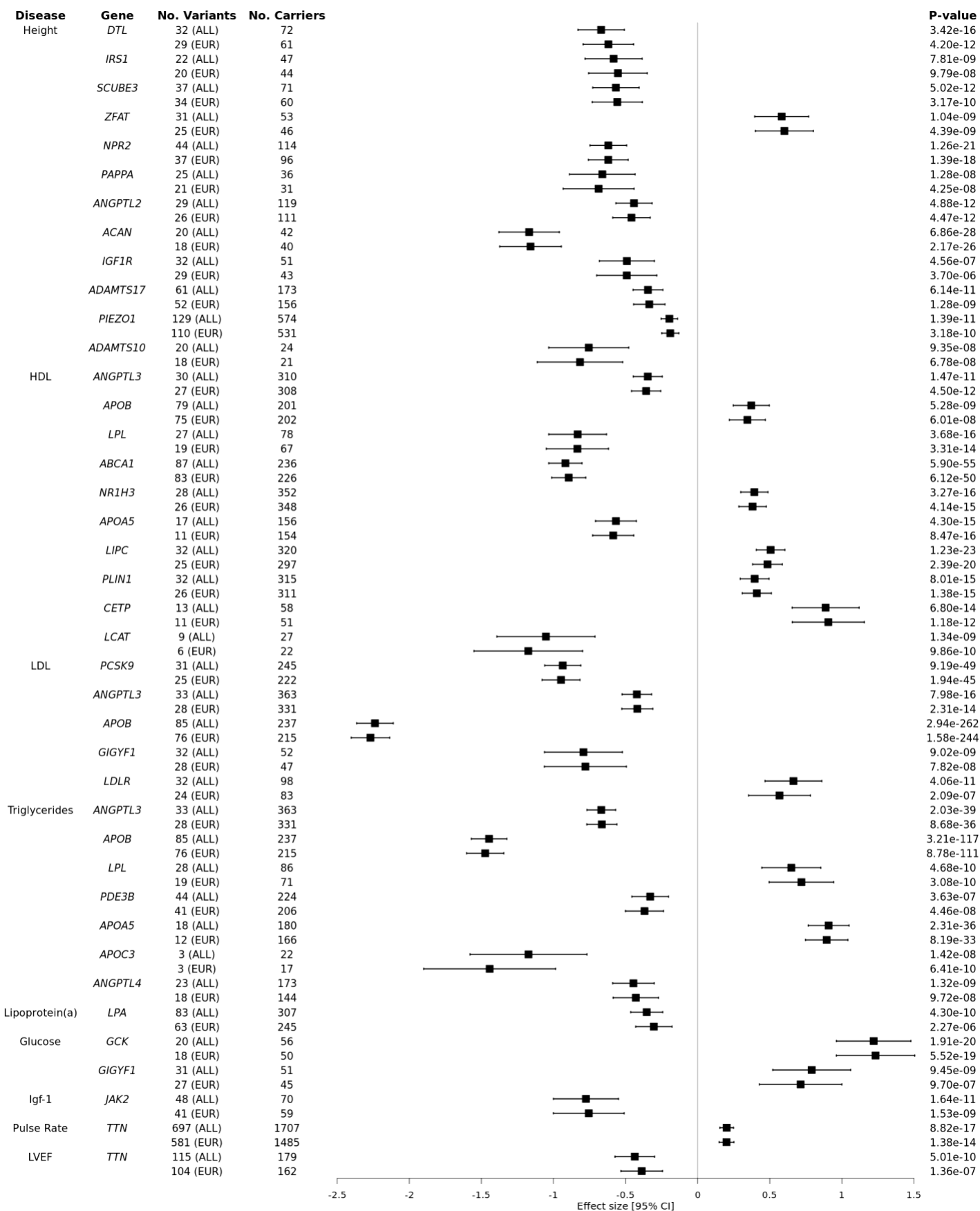


488 **Supplementary Figure 5: Quantile-quantile plots for rare deleterious variants compared to rare**
489 **synonymous variants for height, weight, BMI, and QTc.** The y-axis represents the observed $-\log_{10} P$ -
490 values across all tests, while the x-axis represents the expected under the null-hypothesis. P -values were
491 obtained from score tests in linear mixed effects models, adjusting for sex, age, sequencing batch, associated
492 principal components (PCs), MRI serial number (for MRI traits) and a sparse kinship matrix. P -values shown
493 are two-sided and unadjusted for multiple testing. Left panels represent the exome-wide discovery analysis
494 where we analyzed rare LOF and predicted deleterious missense variants, while the right panels show the
495 results for rare (MAF<0.1%) synonymous variants. As expected under the null, the distributions for the
496 synonymous variants do not show inflation.



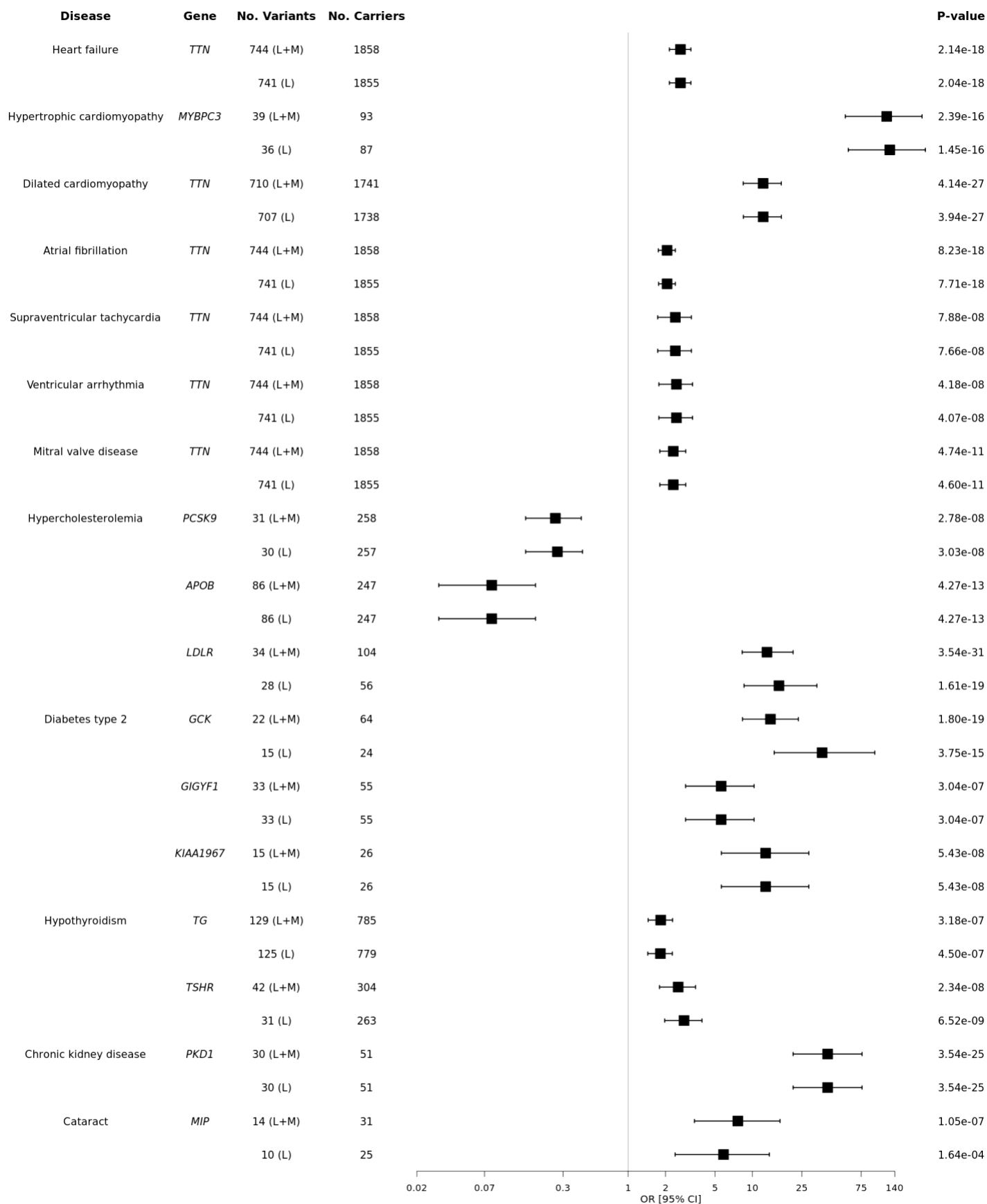
497
498

499 **Supplementary Figure 6. Sensitivity analysis restricting to individuals of European ancestry only in the**
500 **analysis of binary traits.** Data are presented as Odds Ratios (OR) with error bars representing 95%
501 confidence intervals (CI). *P*-values were computed using saddle point approximation and were obtained from
502 logistic mixed effects models, adjusting for sex, age, sequencing batch, associated principal components
503 (PCs), a sparse kinship matrix. *P*-values shown are two-sided and unadjusted for multiple testing. ORs and CIs
504 were obtained from Firth's regression models adjusting for sex, age, sequencing batch and associated PCs
505 among unrelated samples. *P*-values are two-sided and unadjusted for multiple testing. Exome-wide significant
506 associations for binary traits were largely consistent when restricting to a homogenous subset of the cohort
507 consisting of European individuals only. Abbreviations: ALL, all ancestry individuals included; EUR, European
508 ancestry individuals only.



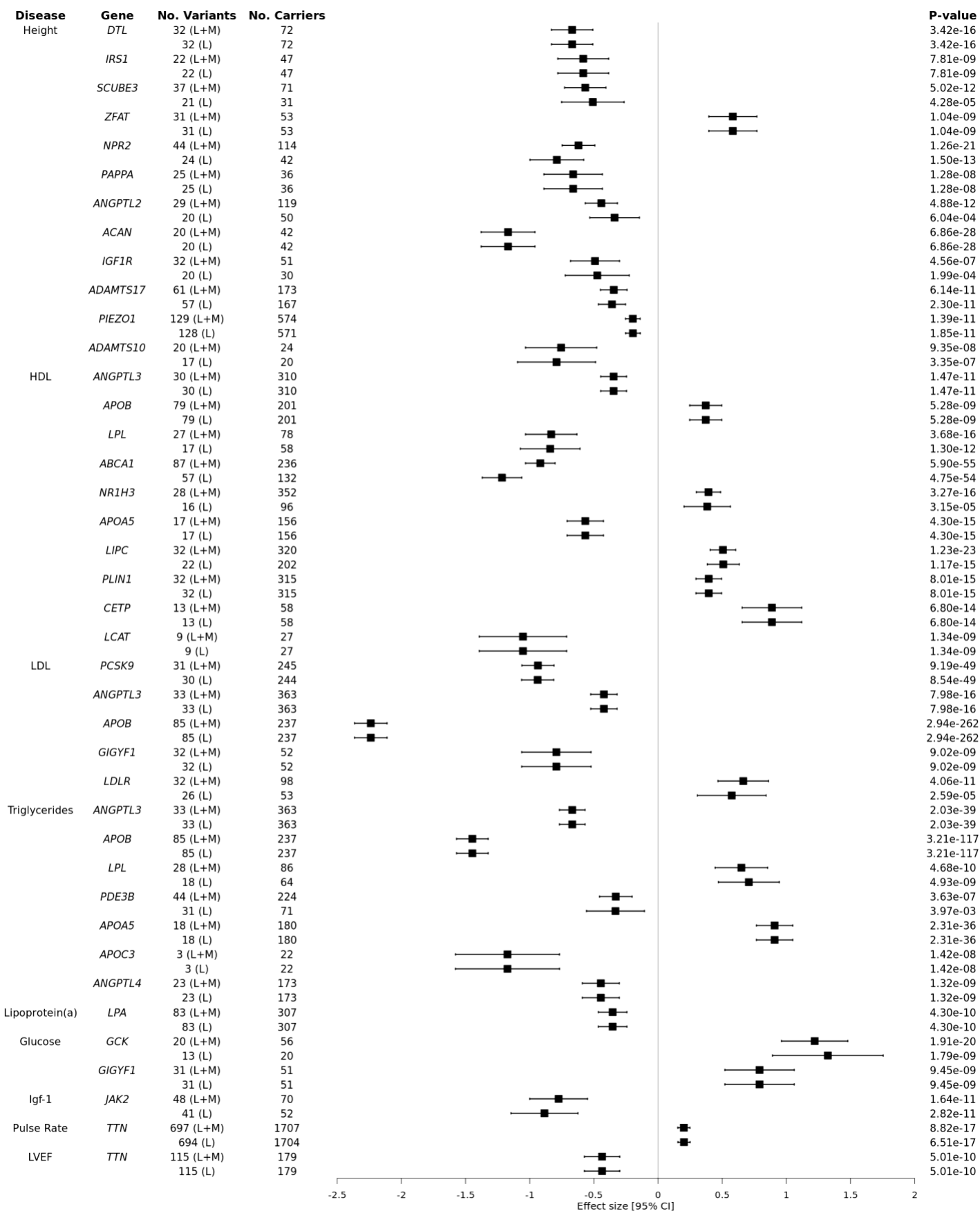
509
510
511

512 **Supplementary Figure 7. Sensitivity analysis restricting to individuals of European ancestry only in the**
513 **analysis of quantitative traits.** Data are presented as effect size (β) estimates per standard deviation with
514 error bars representing 95% confidence intervals (CI). *P*-values, effect sizes and 95% CIs were obtained from
515 score tests in linear mixed effects models, adjusting for sex, age, sequencing batch, associated principal
516 components (PCs), MRI serial number (for MRI traits) and a sparse kinship matrix. *P*-values shown are
517 unadjusted for multiple testing. Associations were largely consistent when restricting to samples from a
518 homogenous subset of European individuals only. Abbreviations: ALL, all ancestry individuals included; EUR,
519 European ancestry individuals only; CI, confidence interval; HDL, high-density lipoprotein; LDL, low-density
520 lipoprotein; Igf-1, insulin-like growth factor-1; QTc, Bazett-corrected QT interval; LVEF, left ventricular ejection
521 fraction; CI, confidence interval.



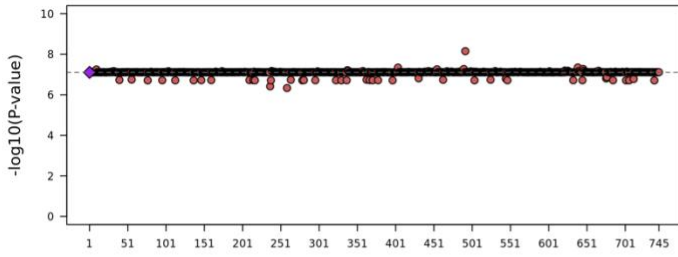
522
523

524 **Supplementary Figure 8. Sensitivity analysis restricting to LOFs only in the primary analysis of binary**
525 **traits.** Data are presented as Odds Ratios (OR) with error bars representing 95% confidence intervals (CI). *P*-
526 values were computed using the saddle point approximation and were obtained from logistic mixed effects
527 models, adjusting for sex, age, sequencing batch, associated principal components (PCs), a sparse kinship
528 matrix. *P*-values shown are two-sided and unadjusted for multiple testing. ORs and CIs were obtained from
529 Firth's regression models adjusting for sex, age, sequencing batch and associated PCs among unrelated
530 samples. *P*-values are two-sided and unadjusted for multiple testing. Effect estimates for analysis of LOFs
531 were largely consistent with effect estimates from LOFs and predicted-damaging missense combined,
532 indicating that in general effect sizes from our discovery analysis our not diluted by the included missense
533 variants. However, effect sizes were attenuated by including missense variants for *GCK*/type 2 diabetes.
534 Interestingly, the *GCK*/diabetes association also dropped in significance after removing missense variants; this
535 indicates that a number of these missense variants were functional, possibly with smaller effect sizes than
536 LOFs. Abbreviations: L, high-confidence loss-of-function variants only; L+M, high-confidence loss-of-function
537 and predicted-damaging missense variants combined; CI, confidence interval.

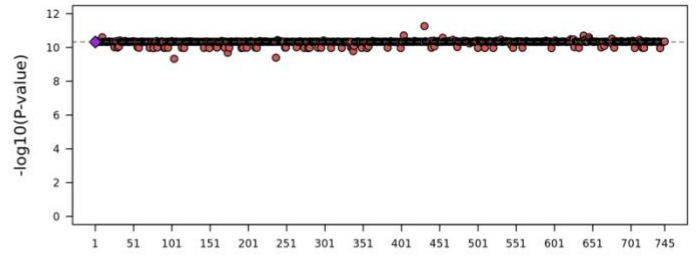


540 **Supplementary Figure 9. Sensitivity analysis restricting to LOFs only in the primary analysis of**
541 **quantitative traits.** Data are presented as effect size (β) estimates per standard deviation with error bars
542 representing 95% confidence intervals (CI). *P*-values, effect sizes and 95% CIs were obtained from score tests
543 in linear mixed effects models, adjusting for sex, age, sequencing batch, associated principal components
544 (PCs), MRI serial number (for MRI traits) and a sparse kinship matrix. *P*-values shown are two-sided and
545 unadjusted for multiple testing. Effect estimates for analysis of LOFs were largely consistent with effect
546 estimates from LOFs and predicted-damaging missense combined. Abbreviations: L, high-confidence loss-of-
547 function variants only; L+M, high-confidence loss-of-function and predicted-damaging missense variants
548 combined; CI, confidence interval; HDL, high-density lipoprotein; LDL, low-density lipoprotein; Igf-1, insulin-like
549 growth factor-1; QTc, Bazett-corrected QT interval; LVEF, left ventricular ejection fraction; CI, confidence
550 interval.

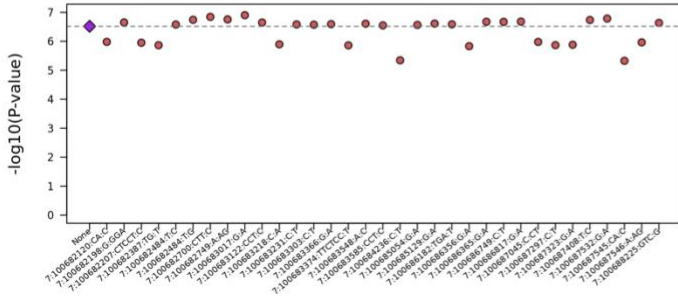
Supraventricular tachycardia and *TTN*



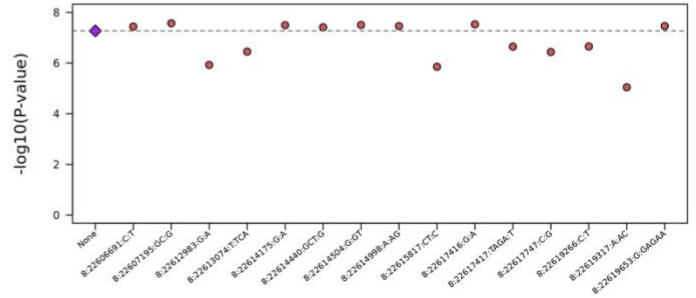
Mitral valve disease and *TTN*



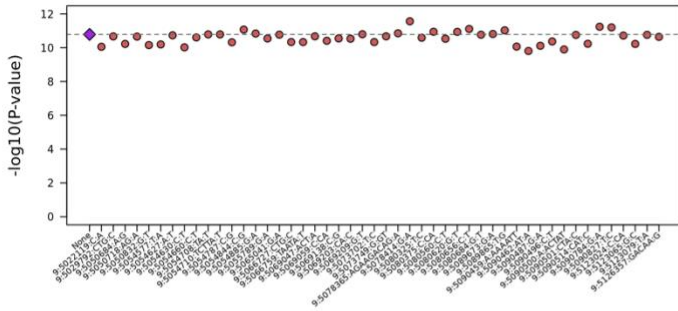
Diabetes type 2 and *GIGYF1*



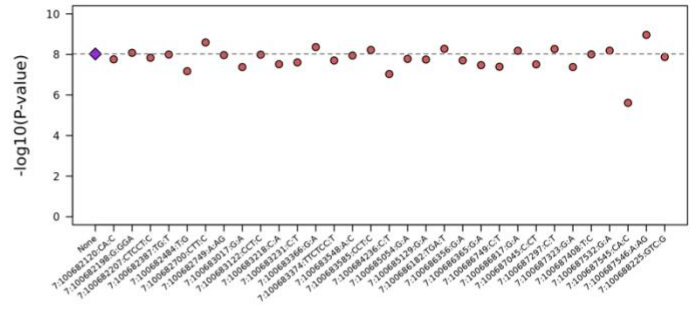
Diabetes type 2 and *CCAR2*



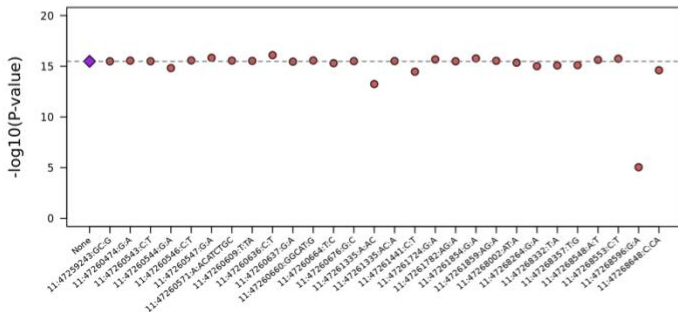
Igf-1 and *JAK2*



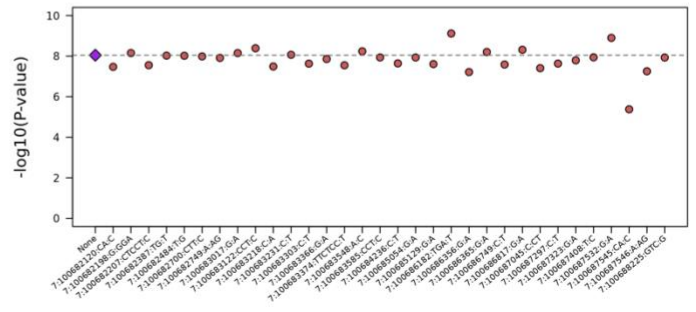
Glucose and *GIGYF1*



HDL and *NR1H3*

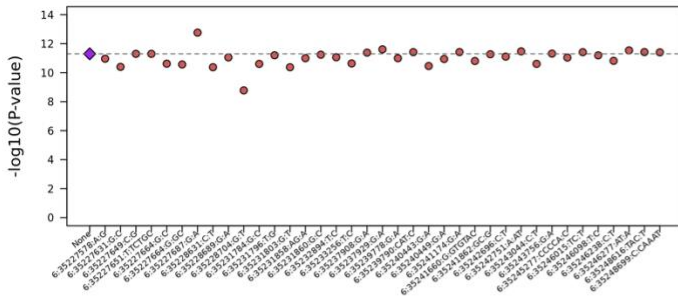


LDL and *GIGYF1*

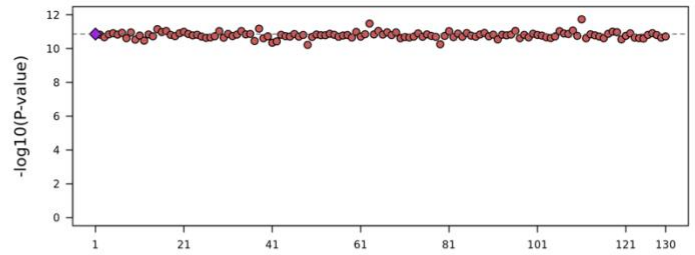


551
552
553
554
555
556

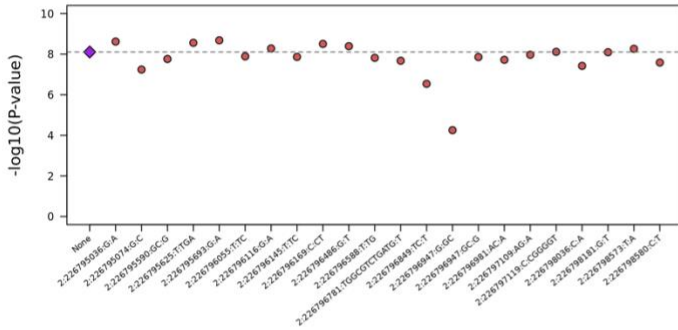
Height and SCUBE3



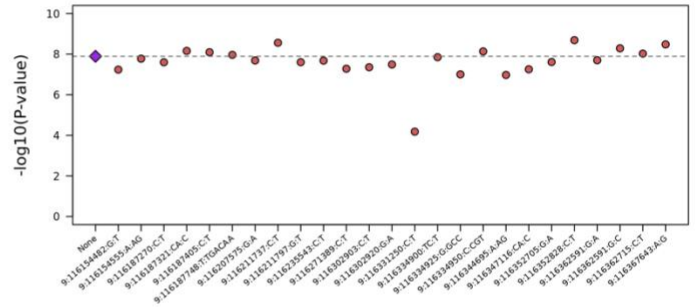
Height and PIEZO1



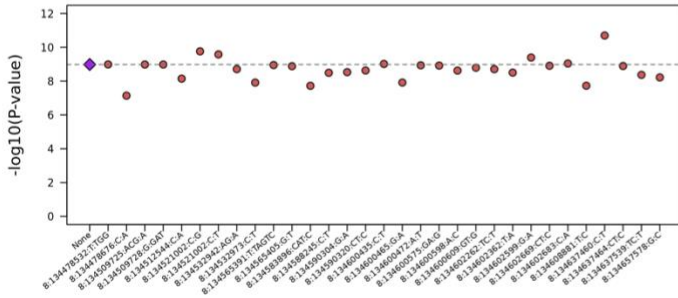
Height and IRS1



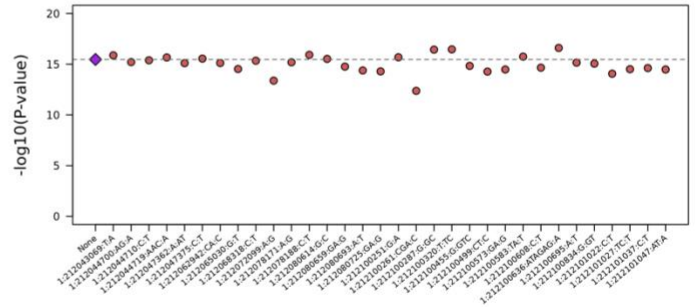
Height and PAPA



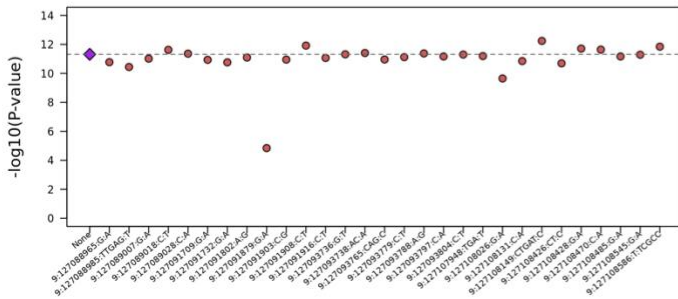
Height and ZFAT



Height and DTL

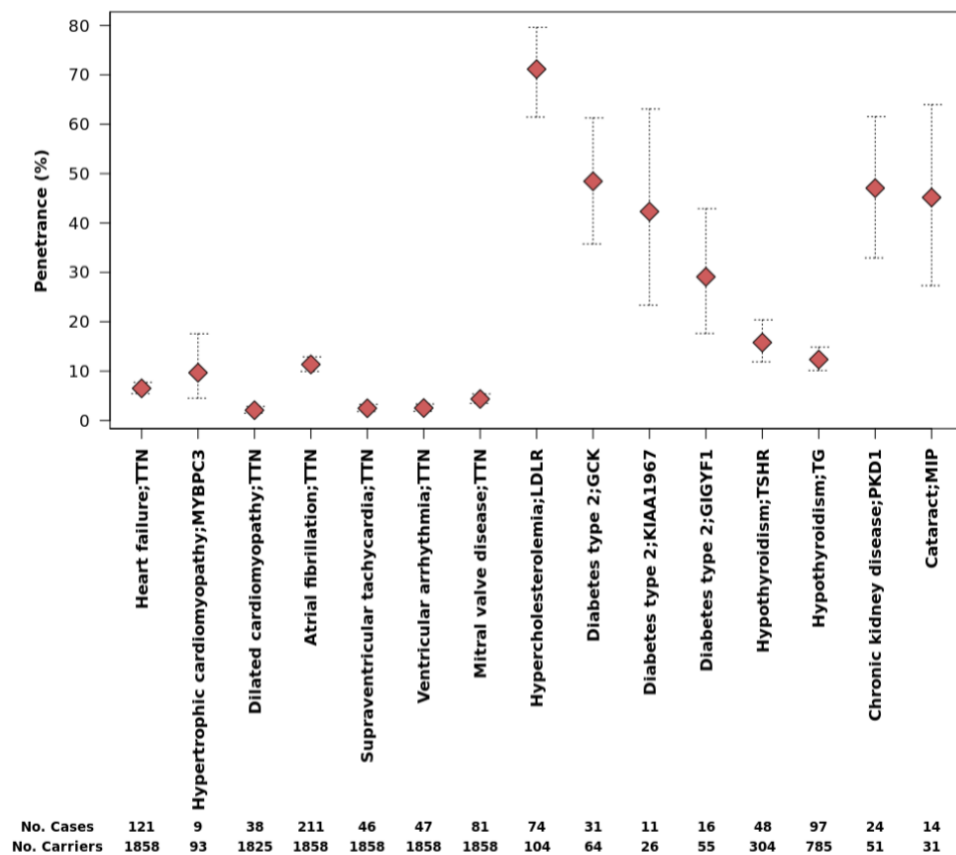


Height and ANGPTL2

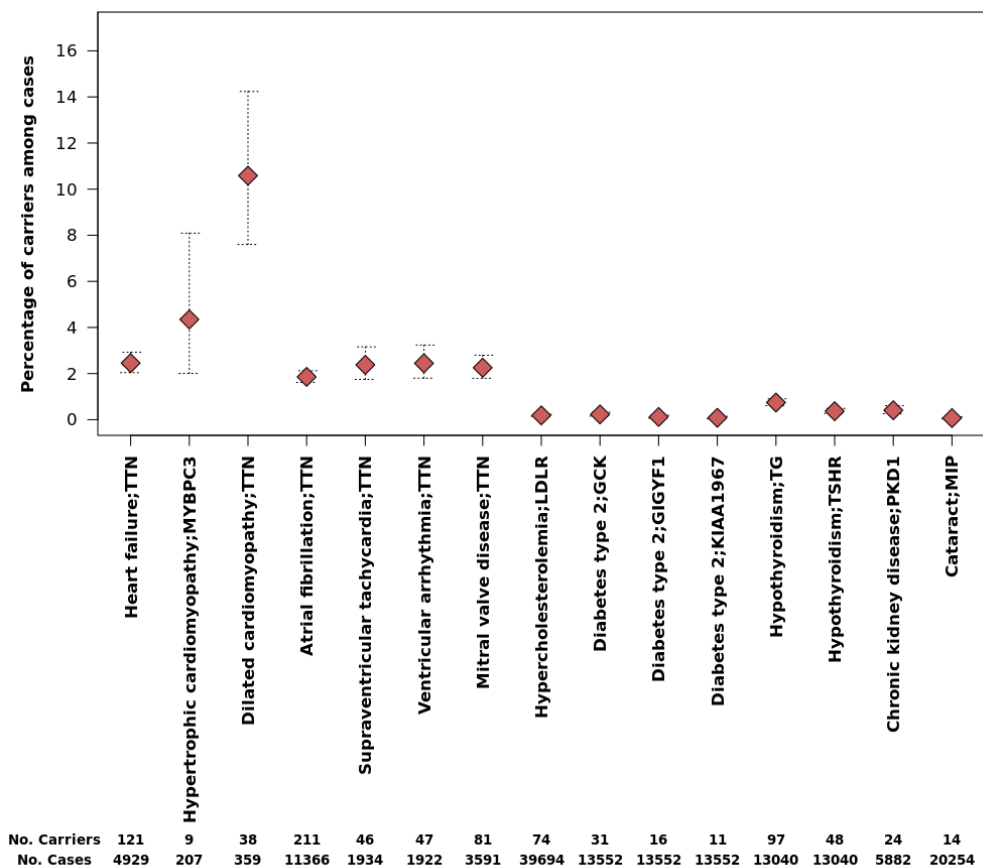


Supplementary Figure 10. Leave-one-variant-out (LOVO) analysis for novel rare variant associations. The x-axis represents a single variant removed from the gene-based analysis, while the y-axis shows the $-\log_{10}$ P -value of the association without that given variant. P -values were obtained from score tests in linear mixed effects models (quantitative traits) or saddle point approximation in logistic mixed effects models (binary traits), adjusting for sex, age, sequencing batch, associated principal components (PCs) and a sparse kinship matrix. P -values shown are two-sided and unadjusted for multiple testing. The first result (diamond) shows the original result without any variant removed. Variants are annotated with the variant name in format

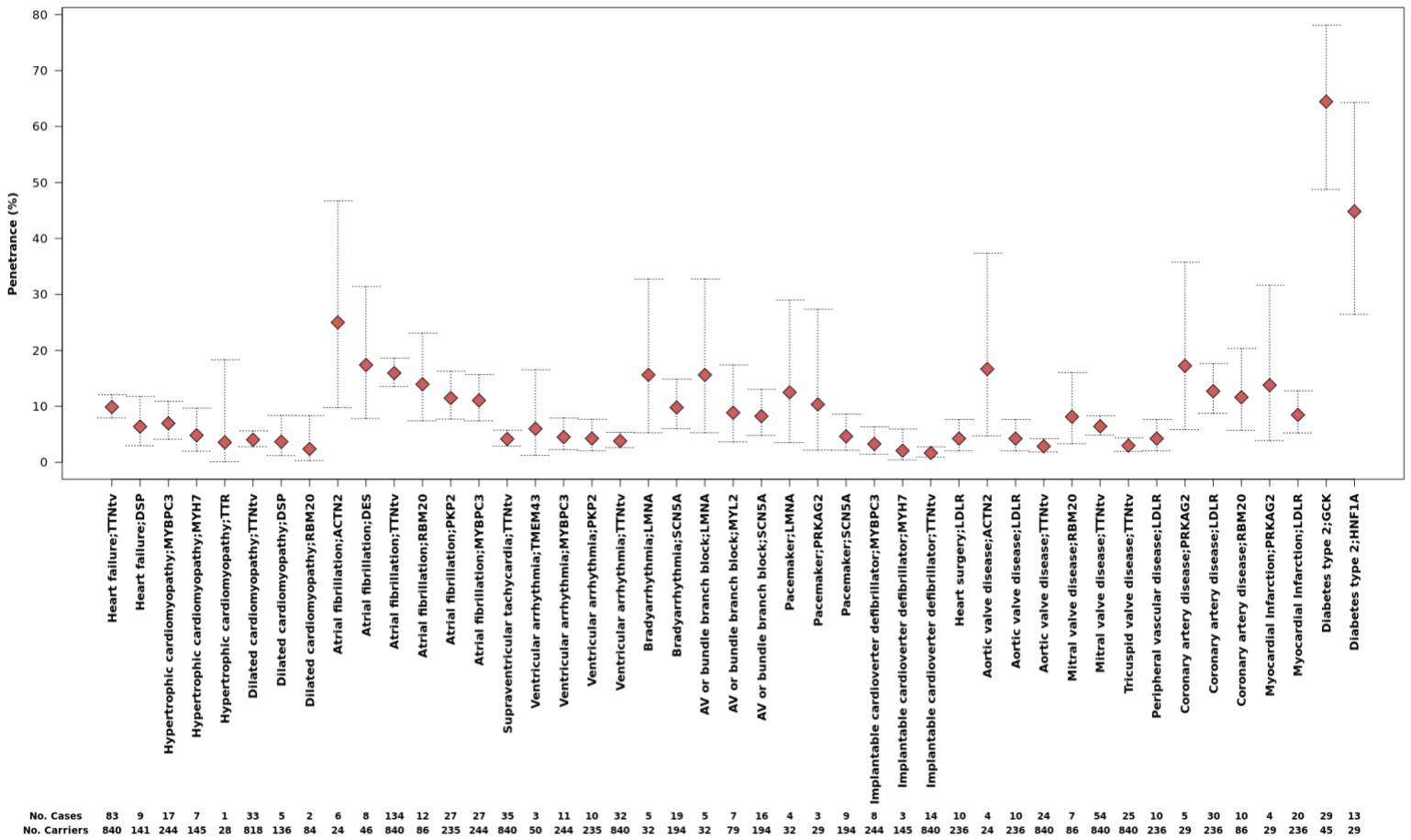
565 chromosome:position:reference:alternative; for *TTN* variant names are not shown given the many variants in
566 the masks. Associations are never abolished upon removing the most important variant from the mask.



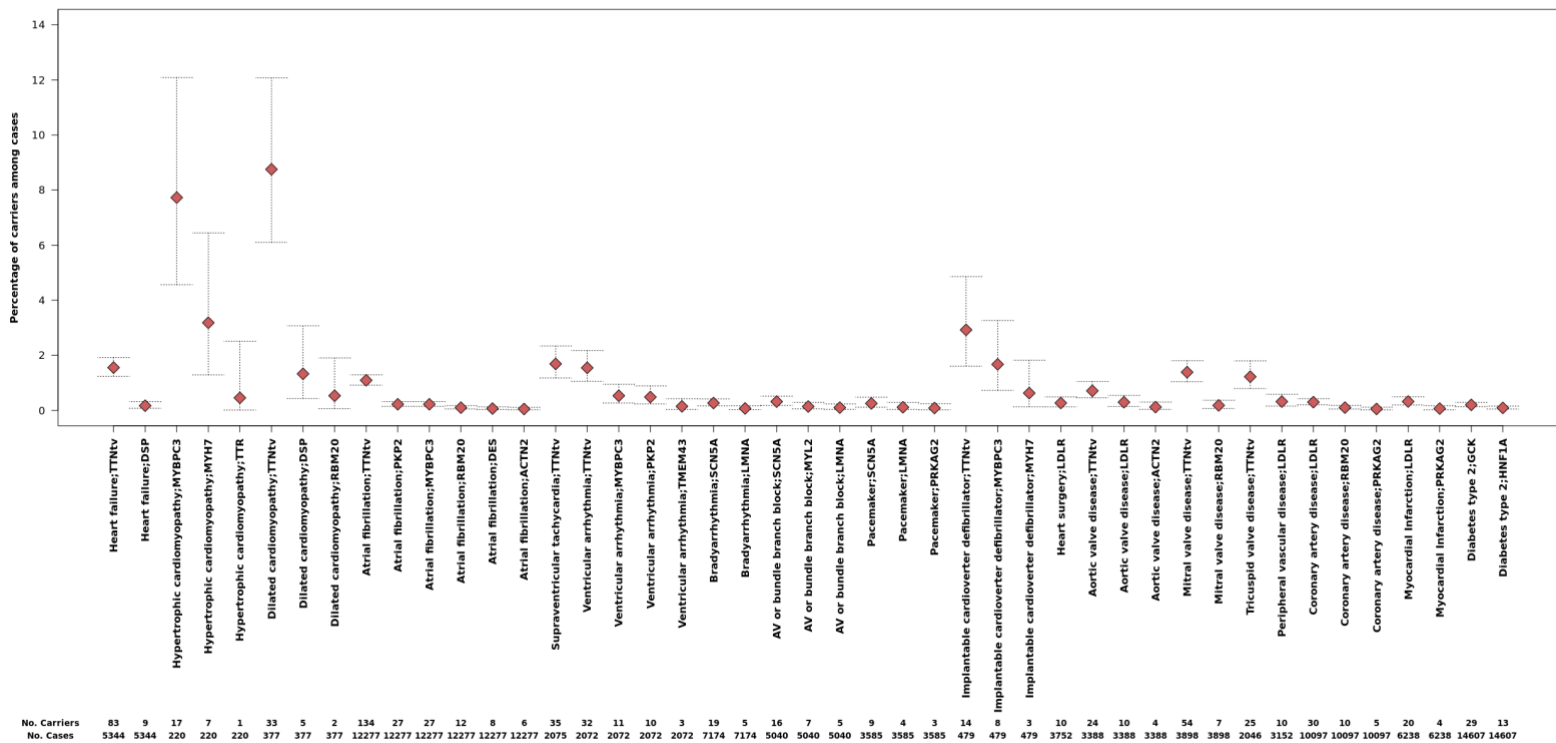
567
568 **Supplementary Figure 11. Penetrance of predicted-damaging variants in genes associated with disease**
569 **in the primary analyses.** The x-axis presents gene-phenotype pairs significantly associated with increased
570 disease risk in the primary analysis of loss-of-function and predicted-deleterious missense variants. 'Significant
571 increased disease risk' was defined as Benjamini-Hochberg two-sided Q-value < 0.01 (computed from *P*-
572 values from all tests in the discovery phase using saddle point approximation in logistic mixed effects models)
573 and Odds Ratio (OR) > 1 (computed from Firth's regression models among unrelated samples). Data on the y-
574 axis are presented as penetrance - calculated as the number of rare variant carriers who were disease cases
575 divided by the total number of carriers times 100% - with dotted lines representing 95% exact binomial
576 confidence intervals. Based on effect sizes, the penetrance estimates for diabetes type 2/*GCK* likely
577 underestimate true loss-of-function, while the other associations should be comparable to loss-of-function
578 variants estimates (**Supplementary Figure 8**).



Supplementary Figure 12. Prevalence of predicted-damaging variants in genes identified in primary analysis among relevant disease cases. The x-axis presents gene-phenotype pairs significantly associated with increased disease risk in the primary analysis of loss-of-function and predicted-deleterious missense variants. ‘Significant increased disease risk’ was defined as Benjamini-Hochberg two-sided Q-value < 0.01 (computed from *P*-values from all tests in the discovery phase using saddle point approximation in logistic mixed effects models) and Odds Ratio (OR) > 1 (computed from Firth’s regression models among unrelated samples). Data on the y-axis are presented as the percentage of rare variant carriers among disease cases - calculated as the number of rare variant carriers who were disease cases divided by the total number of disease cases times 100% - with dotted lines representing 95% exact binomial confidence intervals. Among individuals with dilated cardiomyopathy, up to 12% may carry rare variants in *TTN*. In general, however, rare high-impact variants are rare among common adult-onset disease cases.



592 **Supplementary Figure 13. Penetrance of putatively pathogenic variants in cardiovascular disease and**
593 **diabetes panel genes for relevant phenotypes.** The x-axis presents gene-phenotype pairs showing at least
594 suggestive evidence of association with increased disease risk in the analysis of putatively pathogenic variants
595 in cardiovascular and diabetes panel genes. ‘Suggestive increased disease risk’ was defined as two-sided *P*-
596 value < 0.005 (unadjusted for multiple testing; computed using saddle point approximation in logistic mixed
597 effects models) and Odds Ratio (OR) > 1 (computed from Firth’s regression models among unrelated
598 samples). Data on the y-axis are presented as penetrance - calculated as the number of rare variant carriers
599 who were disease cases divided by the total number of carriers times 100% - with dotted lines representing
600 95% exact binomial confidence intervals.
601



603

604

605

606

607

608

609

610

611

612

Supplementary Figure 14. Prevalence of putatively pathogenic variants in cardiovascular disease and diabetes panel genes among disease cases. The x-axis presents gene-phenotype pairs showing at least suggestive evidence of association in the analysis of putatively pathogenic variants in cardiovascular and diabetes panel genes. ‘Suggestive increased disease risk’ was defined as two-sided P -value < 0.005 (unadjusted for multiple testing; computed using saddle point approximation in logistic mixed effects models) and Odds Ratio (OR) > 1 (computed from Firth’s regression models among unrelated samples). Data on the y-axis are presented as the percentage of rare variant carriers among disease cases - calculated as the number of rare variant carriers who were disease cases divided by the total number of disease cases times 100% - with dotted lines representing 95% exact binomial confidence intervals.

Supplementary References

1. Szustakowski, J.D. *et al.* Advancing Human Genetics Research and Drug Discovery through Exome Sequencing of the UK Biobank. *medRxiv* (2020).
2. Regier, A.A. *et al.* Functional equivalence of genome sequencing analysis pipelines enables harmonized variant calling across human genetics projects. *Nat Commun* **9**, 4038 (2018).
3. Yun, T. *et al.* Accurate, scalable cohort variant calls using DeepVariant and GLnexus. *Bioinformatics* **36**, 5582–5589 (2021).
4. Chang, C.C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience* **4**, 7 (2015).
5. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).
6. Manichaikul, A. *et al.* Robust relationship inference in genome-wide association studies. *Bioinformatics* **26**, 2867–73 (2010).
7. Conomos, M.P., Miller, M.B. & Thornton, T.A. Robust inference of population structure for ancestry prediction and correction of stratification in the presence of relatedness. *Genet Epidemiol* **39**, 276–93 (2015).
8. Galinsky, K.J. *et al.* Fast Principal-Component Analysis Reveals Convergent Evolution of ADH1B in Europe and East Asia. *Am J Hum Genet* **98**, 456–472 (2016).
9. Karczewski, K.J. *et al.* The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **581**, 434–443 (2020).
10. Type 2 Diabetes Knowledge Portal. Accessed in December 2020 and June 2021; <http://www.type2diabetesgenetics.org>
11. Willer, C.J., Li, Y. & Abecasis, G.R. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* **26**, 2190–1 (2010).
12. de Leeuw, C.A., Mooij, J.M., Heskes, T. & Posthuma, D. MAGMA: generalized gene-set analysis of GWAS data. *PLoS Comput Biol* **11**, e1004219 (2015).
13. Consortium, G. The GTEx Consortium atlas of genetic regulatory effects across human tissues. *Science* **369**, 1318–1330 (2020).
14. Roberts, A.M. *et al.* Integrated allelic, transcriptional, and phenomic dissection of the cardiac effects of titin truncations in health and disease. *Sci Transl Med* **7**, 270ra6 (2015).
15. Choi, S.H. *et al.* Monogenic and Polygenic Contributions to Atrial Fibrillation Risk: Results From a National Biobank. *Circ Res* **126**, 200–209 (2020).
16. Cornec-Le Gall, E. *et al.* Type of PKD1 mutation influences renal outcome in ADPKD. *J Am Soc Nephrol* **24**, 1006–13 (2013).
17. Mallawaarachchi, A.C., Furlong, T.J., Shine, J., Harris, P.C. & Cowley, M.J. Population data improves variant interpretation in autosomal dominant polycystic kidney disease. *Genet Med* **21**, 1425–1434 (2019).
18. Ali, H. *et al.* PKD1 Duplicated regions limit clinical Utility of Whole Exome Sequencing for Genetic Diagnosis of Autosomal Dominant Polycystic Kidney Disease. *Sci Rep* **9**, 4141 (2019).
19. Rossetti, S. *et al.* Identification of gene mutations in autosomal dominant polycystic kidney disease through targeted resequencing. *J Am Soc Nephrol* **23**, 915–33 (2012).
20. Trujillano, D. *et al.* Diagnosis of autosomal dominant polycystic kidney disease using efficient PKD1 and PKD2 targeted next-generation sequencing. *Mol Genet Genomic Med* **2**, 412–21 (2014).
21. Akhtar, M.M. *et al.* Clinical Phenotypes and Prognosis of Dilated Cardiomyopathy Caused by Truncating Variants in the. *Circ Heart Fail* **13**, e006832 (2020).
22. Franaszczyk, M. *et al.* Titin Truncating Variants in Dilated Cardiomyopathy - Prevalence and Genotype-Phenotype Correlations. *PLoS One* **12**, e0169007 (2017).
23. Haggerty, C.M. *et al.* Genomics-First Evaluation of Heart Disease Associated With Titin-Truncating Variants. *Circulation* **140**, 42–54 (2019).
24. Michels, M. *et al.* Disease penetrance and risk stratification for sudden cardiac death in asymptomatic hypertrophic cardiomyopathy mutation carriers. *Eur Heart J* **30**, 2593–8 (2009).
25. Viswanathan, S.K. *et al.* Hypertrophic cardiomyopathy clinical phenotype is independent of gene mutation and mutation dosage. *PLoS One* **12**, e0187948 (2017).

- 667 26. Adalsteinsdottir, B. *et al.* Hypertrophic cardiomyopathy in myosin-binding protein C (MYBPC3) Icelandic
668 founder mutation carriers. *Open Heart* **7**, e001220 (2020).
- 669 27. Lorenzini, M. *et al.* Penetrance of Hypertrophic Cardiomyopathy in Sarcomere Protein Mutation
670 Carriers. *J Am Coll Cardiol* **76**, 550-559 (2020).
- 671 28. Walsh, R. *et al.* Reassessment of Mendelian gene pathogenicity using 7,855 cardiomyopathy cases
672 and 60,706 reference samples. *Genet Med* **19**, 192-203 (2017).
- 673 29. Groeneweg, J.A. *et al.* Clinical Presentation, Long-Term Follow-Up, and Outcomes of 1001
674 Arrhythmogenic Right Ventricular Dysplasia/Cardiomyopathy Patients and Family Members. *Circ*
675 *Cardiovasc Genet* **8**, 437-46 (2015).
- 676 30. Patel, K.A. *et al.* Heterozygous RFX6 protein truncating variants are associated with MODY with
677 reduced penetrance. *Nat Commun* **8**, 888 (2017).
- 678 31. Chakera, A.J. *et al.* Recognition and Management of Individuals With Hyperglycemia Because of a
679 Heterozygous Glucokinase Mutation. *Diabetes Care* **38**, 1383-92 (2015).
- 680 32. Wang, Z. *et al.* Identification and functional analysis of GCK gene mutations in 12 Chinese families with
681 hyperglycemia. *J Diabetes Investig* **10**, 963-971 (2019).
- 682 33. Bansal, V. *et al.* Spectrum of mutations in monogenic diabetes genes identified from high-throughput
683 DNA sequencing of 6888 individuals. *BMC Med* **15**, 213 (2017).
- 684 34. Wright, C.F. *et al.* Assessing the Pathogenicity, Penetrance, and Expressivity of Putative Disease-
685 Causing Variants in a Population Setting. *Am J Hum Genet* **104**, 275-286 (2019).
- 686