

Determinants of penetrance and variable expressivity in monogenic metabolic conditions across 77,184 exomes

Julia Goodrich¹, Moriel Singer-Berk¹, Rachel Son¹, Abigail Sveden¹, Jordan Wood¹, Eleina England¹, Joanne B. Cole¹, Ben Weisburd¹, Nick Watts¹, Zachary Zappala¹, Haichen Zhang², Kristin A. Maloney², Andy Dahl³, Carlos A. Aguilar-Salinas⁴, Gil Atzmon⁵⁻⁷, Francisco Barajas-Olmos⁸, Nir Barzilai^{5,7}, John Blangero⁹, Eric Boerwinkle^{10,11}, Lori L. Bonnycastle¹², Erwin Bottinger¹³, Donald W Bowden¹⁴⁻¹⁶, Federico Centeno- Cruz⁸, John C. Chambers^{17,18}, Nathalie Chami^{13,19}, Edmund Chan²⁰, Juliana Chan²¹⁻²⁴, Ching-Yu Cheng^{25,26,27}, Yoon Shin Cho²⁸, Cecilia Contreras-Cubas⁸, Emilio Córdova⁸, Adolfo Correa²⁹, Ralph A. DeFronzo³⁰, Ravindranath Duggirala⁹, Josée Dupuis³¹, Ma. Eugenia Garay-Sevilla³², Humberto García-Ortiz⁸, Christian Gieger³³, Benjamin Glaser³⁴, Clicerio González-Villalpando³⁵, Ma Elena Gonzalez³⁶, Niels Garup³⁷, Leif Groop^{38,39}, Myron Gross⁴⁰, Christopher Haiman⁴¹, Sohee Han⁴², Craig L Hanis¹⁰, Torben Hansen³⁷, Nancy L. Heard-Costa^{43,44}, Brian E Henderson⁴¹, Juan Manuel Malacara Hernandez³², Mi Yeong Hwang⁴², Sergio Islas-Andrade⁸, Marit E Jørgensen⁴⁵⁻⁴⁷, Hyun Min Kang⁴⁸, Bong-Jo Kim⁴², Young Jin Kim⁴², Heikki A. Koistinen⁴⁹⁻⁵¹, Jaspal Singh Kooner⁵²⁻⁵⁵, Johanna Kuusisto⁵⁶, Soo-Heon Kwak⁵⁷, Markku Laakso⁵⁶, Leslie Lange⁵⁸, Jong-Young Lee⁵⁹, Juyoung Lee⁴², Donna M. Lehman³⁰, Allan Linneberg⁶⁰⁻⁶², Jianjun Liu^{63,20,64}, Ruth J.F. Loos^{13,19}, Valeriya Lyssenko^{38,65}, Ronald C. W. Ma²¹⁻²⁴, Angélica Martínez-Hernández⁸, James B. Meigs^{1,66,67}, Thomas Meitinger^{68,69}, Elvia Mendoza- Caamal⁸, Karen L. Mohlke⁷⁰, Andrew D. Morris^{71,72}, Alanna C. Morrison¹⁰, Maggie CY Ng¹⁴⁻¹⁶, Peter M. Nilsson⁷³, Christopher J. O'Donnell⁷⁴⁻⁷⁷, Lorena Orozco⁸, Colin N. A. Palmer⁷⁸, Kyong Soo Park^{57,79,80}, Wendy S. Post⁸¹, Oluf Pedersen³⁷, Michael Preuss¹³, Bruce M. Psaty^{82,83}, Alexander P. Reiner⁸⁴, Cristina Revilla-Monsalve⁸⁵, Stephen S Rich⁸⁶, Jerome I Rotter⁸⁷, Danish Saleheen⁸⁸⁻⁹⁰, Claudia Schurmann^{13,91,92}, Xueling Sim⁶³, Rob Sladek⁹³⁻⁹⁵, Kerrin S Small⁹⁶, Wing Yee So²¹⁻²³, Xavier Soberón⁸, Timothy D Spector⁹⁶, Konstantin Strauch^{97,98}, Tim M Strom^{99,68}, E Shyong Tai^{63,20,27}, Claudia H.T. Tam²¹⁻²³, Yik Ying Teo^{63,100,101}, Farook Thameem¹⁰², Brian Tomlinson¹⁰³, Russell P. Tracy^{104,105}, Tiinamaija Tuomi^{39,106-108}, Jaakko Tuomilehto¹⁰⁹⁻¹¹², Teresa Tusié-Luna^{113,114}, Rob M. van Dam^{63,20,115}, Ramachandran S. Vasan^{43,116}, James G Wilson¹¹⁷, Daniel R Witte^{118,119}, Tien-Yin Wong^{25,26,27}, Lizz Caulkins¹, Noël P. Burt¹, Noah Zaitlen³, Mark I. McCarthy¹²⁰⁻¹²², Michael Boehnke⁴⁸, Toni I. Pollin², Jason Flannick^{1,123,124}, Josep M. Mercader^{1,125,126}, Anne O'Donnell-Luria^{1,123,124}, Samantha Baxter¹, Jose C. Florez^{1,125,126}, Daniel MacArthur^{1,127,128}, Miriam S. Udler-Aubrey^{1,125,126}, for AMP-T2D-GENES Consortia

Affiliations

¹ Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA.

² School of Medicine, University of Maryland Baltimore, Baltimore, MD, USA.

³ Department of Neurology, UCLA, Los Angeles, California.

⁴ Instituto Nacional de Ciencias Medicas y Nutricion, Mexico City, Mexico.

⁵ Department of Medicine, Albert Einstein College of Medicine, New York, NY, USA.

⁶ Faculty of Natural Science, University of Haifa, Haifa, Israel.

⁷ Department of Genetics, Albert Einstein College of Medicine, New York, NY, USA.

⁸ Instituto Nacional de Medicina Genómica, Mexico City, Mexico.

⁹ Department of Human Genetics and South Texas Diabetes and Obesity Institute, University of Texas Rio Grande Valley, Brownsville and Edinburg, TX, USA

- ¹⁰ Human Genetics Center, School of Public Health, The University of Texas Health Science Center at Houston, Houston, TX, USA
- ¹¹ Human Genome Sequencing Center, Baylor College of Medicine, Houston, TX, USA
- ¹² Medical Genomics and Metabolic Genetics Branch, National Human Genome Research Institute, National Institutes of Health, Bethesda, MD, USA.
- ¹³ The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, NY, USA.
- ¹⁴ Center for Diabetes Research, Wake Forest School of Medicine, Winston-Salem, NC, USA.
- ¹⁵ Center for Genomics and Personalized Medicine Research, Wake Forest School of Medicine, Winston-Salem, NC, USA.
- ¹⁶ Department of Biochemistry, Wake Forest School of Medicine, Winston-Salem, NC, USA
- ¹⁷ Department of Epidemiology and Biostatistics, Imperial College London, UK
- ¹⁸ Lee Kong Chian School of Medicine, Nanyang Technological University, Singapore, Singapore
- ¹⁹ The Mindich Child Health and Development Institute, Ichan School of Medicine at Mount Sinai, New York, NY, USA.
- ²⁰ Department of Medicine, Yong Loo Lin School of Medicine, National University of Singapore and National University Health System, Singapore, Singapore
- ²¹ Department of Medicine and Therapeutics, The Chinese University of Hong Kong, Hong Kong, China.
- ²² Chinese University of Hong Kong-Shanghai Jiao Tong University Joint Research Centre in Diabetes Genomics and Precision Medicine, The Chinese University of Hong Kong, Hong Kong, China.
- ²³ Hong Kong Institute of Diabetes and Obesity, The Chinese University of Hong Kong, Hong Kong, China.
- ²⁴ Li Ka Shing Institute of Health Sciences, The Chinese University of Hong Kong, Hong Kong, China.
- ²⁵ Singapore Eye Research Institute, Singapore National Eye Centre, Singapore, Singapore
- ²⁶ Department of Ophthalmology, Yong Loo Lin School of Medicine, National University of Singapore and National University Health System, Singapore, Singapore.
- ²⁷ Duke-NUS Medical School, Singapore, Singapore.
- ²⁸ Department of Biomedical Science, Hallym University, Chuncheon, South Korea.
- ²⁹ Department of Medicine, University of Mississippi Medical Center, Jackson, MS, USA.
- ³⁰ Department of Medicine, University of Texas Health San Antonio (aka University of Texas Health Science Center at San Antonio), San Antonio, TX, USA
- ³¹ Department of Biostatistics, Boston University School of Public Health, Boston, MA, USA.
- ³² Department of Medical Science, División of Health Science. University of Guanajuato. Campus León. León, Gto. México.
- ³³ Institute of Human Genetics, Helmholtz Zentrum München, German Research Center for Environmental Health, Neuherberg, Germany.
- ³⁴ Endocrinology and Metabolism Service, Hadassah-Hebrew University Medical Center, Jerusalem, Israel.
- ³⁵ Unidad de Investigación en Diabetes y Riesgo Cardiovascular, Instituto Nacional de Salud Pública, Cuernavaca, Mexico.
- ³⁶ Centro de Estudios en Diabetes, Mexico City, Mexico.
- ³⁷ Novo Nordisk Foundation Center for Basic Metabolic Research, Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen, Denmark.
- ³⁸ Department of Clinical Sciences, Diabetes and Endocrinology, Lund University Diabetes Centre, Malmö, Sweden.
- ³⁹ Institute for Molecular Genetics Finland, University of Helsinki, Helsinki, Finland.
- ⁴⁰ Department of Laboratory Medicine and Pathology, University of Minnesota, Minneapolis, MN, USA.
- ⁴¹ Department of Preventive Medicine, Keck School of Medicine of USC, Los Angeles, CA, USA.
- ⁴² Division of Genome Research, Center for Genome Science, National Institute of Health, Chungcheongbuk-do, South Korea.

- ⁴³ Boston University and National Heart Lung and Blood Institute's Framingham Heart Study, Framingham, MA, USA.
- ⁴⁴ Department of Neurology, Boston University School of Medicine, Boston, MA, USA.
- ⁴⁵ Steno Diabetes Center Copenhagen, Gentofte, Denmark.
- ⁴⁶ National Institute of Public Health, University of Southern Denmark, Copenhagen, Denmark.
- ⁴⁷ Greenland Centre for Health Research, University of Greenland, Nuuk, Greenland.
- ⁴⁸ Department of Biostatistics and Center for Statistical Genetics, University of Michigan, Ann Arbor, MI, USA.
- ⁴⁹ Department of Public Health Solutions, Finnish Institute for Health and Welfare, Helsinki, Finland.
- ⁵⁰ University of Helsinki and Department of Medicine, Helsinki University Central Hospital, Helsinki, Finland.
- ⁵¹ Minerva Foundation Institute for Medical Research, Helsinki, Finland.
- ⁵² Department of Cardiology, Ealing Hospital, London North West Healthcare NHS Trust, London, UK.
- ⁵³ MRC-PHE Centre for Environment and Health, Imperial College London, London, UK.
- ⁵⁴ Imperial College Healthcare NHS Trust, Imperial College London, London, UK.
- ⁵⁵ National Heart and Lung Institute, Imperial College London, London, UK.
- ⁵⁶ Institute of Clinical Medicine, Internal Medicine, University of Eastern Finland and Kuopio University Hospital, Kuopio, Finland
- ⁵⁷ Department of Internal Medicine, Seoul National University Hospital, Seoul, South Korea.
- ⁵⁸ Department of Medicine, University of Colorado Denver, Anschutz Medical Campus, Aurora, CO, USA.
- ⁵⁹ Oneomics Soonchunhyang Mirae Medical Center, Bucheon-si Gyeonggi-do 14585, Republic of Korea
- ⁶⁰ Department of Clinical Medicine, Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen, Denmark
- ⁶¹ Center for Clinical Research and Prevention, Bispebjerg and Frederiksberg Hospital, Copenhagen, Denmark.
- ⁶² Department of Clinical Experimental Research, Rigshospitalet, Copenhagen, Denmark.
- ⁶³ Saw Swee Hock School of Public Health, National University of Singapore and National University Health System, Singapore, Singapore.
- ⁶⁴ Genome Institute of Singapore, Agency for Science Technology and Research, Singapore, Singapore.
- ⁶⁵ Department of Clinical Science, University of Bergen, Bergen, Norway.
- ⁶⁶ Department of Medicine, Harvard Medical School, Boston, MA, USA.
- ⁶⁷ Division of General Internal Medicine, Massachusetts General Hospital, Boston, MA, USA.
- ⁶⁸ Institute of Human Genetics, Technical University of Munich, Munich, Germany.
- ⁶⁹ German Centre for Cardiovascular Research (DZHK), Partner Site Munich Heart Alliance, Munich, Germany.
- ⁷⁰ Department of Genetics, University of North Carolina Chapel Hill, Chapel Hill, NC, USA.
- ⁷¹ Wellcome Centre for Human Genetics, Nuffield Department of Medicine, University of Oxford, Oxford, UK.
- ⁷² Department of Biostatistics, University of Liverpool, Liverpool, UK.
- ⁷³ Department of Clinical Sciences, Medicine, Lund University, Malmö, Sweden
- ⁷⁴ Department of Pediatrics, Harvard Medical School, Boston, MA, USA.
- ⁷⁵ Section of Cardiology, Department of Medicine, VA Boston Healthcare, Boston, MA, USA.
- ⁷⁶ Brigham and Women's Hospital, Boston, MA, USA.
- ⁷⁷ Intramural Administration Management Branch, National Heart Lung and Blood Institute, NIH, Framingham, MA, USA.
- ⁷⁸ Pat Macpherson Centre for Pharmacogenetics and Pharmacogenomics, University of Dundee, Dundee, UK.
- ⁷⁹ Department of Internal Medicine, Seoul National University College of Medicine, Seoul, South Korea.
- ⁸⁰ Department of Molecular Medicine and Biopharmaceutical Sciences, Graduate School of Convergence Science and Technology, Seoul National University, Seoul, South Korea.
- ⁸¹ Division of Cardiology, Department of Medicine, Johns Hopkins University, Baltimore, MD, USA.
- ⁸² Cardiovascular Health Research Unit, Departments of Medicine, Epidemiology, and Health Services, University of Washington, Seattle, WA, USA.
- ⁸³ Kaiser Permanente Washington Research Institute, Seattle, WA, USA.
- ⁸⁴ Fred Hutchinson Cancer Research Center, Seattle, WA, USA.

- ⁸⁵ Instituto Mexicano del Seguro Social SXXI, Mexico City, Mexico.
- ⁸⁶ Center for Public Health Genomics, University of Virginia School of Medicine, Charlottesville, VA, USA.
- ⁸⁷ The Institute for Translational Genomics and Population Sciences, Department of Pediatrics, The Lundquist Institute for Biomedical Innovation (formerly Los Angeles Biomedical Research Institute) at Harbor-UCLA Medical Center, Torrance, CA, USA.
- ⁸⁸ Division of Translational Medicine and Human Genetics, University of Pennsylvania, Philadelphia, PA, USA.
- ⁸⁹ Department of Biostatistics and Epidemiology, University of Pennsylvania, Philadelphia, PA, USA.
- ⁹⁰ Center for Non-Communicable Diseases, Karachi, Pakistan.
- ⁹¹ Digital Health Center, Hasso Plattner Institute, University of Potsdam, Prof.-Dr.-Helmert-Str. 2-3, 14482 Potsdam, Germany
- ⁹² Hasso Plattner Institute for Digital Health at Mount Sinai, Icahn School of Medicine at Mount Sinai, One Gustave L. Levy Place, New York, NY 10029, USA
- ⁹³ Department of Human Genetics, McGill University, Montreal, Quebec, Canada.
- ⁹⁴ Division of Endocrinology and Metabolism, Department of Medicine, McGill University, Montreal, Quebec, Canada.
- ⁹⁵ McGill University and Génome Québec Innovation Centre, Montreal, Quebec, Canada.
- ⁹⁶ Department of Twin Research and Genetic Epidemiology, King's College London, London, UK.
- ⁹⁷ Institute of Genetic Epidemiology, Helmholtz Zentrum Munchen, German Research Center for Environmental Health, Neuherberg, Germany.
- ⁹⁸ Institute for Medical Informatics Biometry and Epidemiology, Ludwig-Maximilians University, Munich, Germany.
- ⁹⁹ Institute of Human Genetics, Helmholtz Zentrum München, German Research Center for Environmental Health, Neuherberg, Germany.
- ¹⁰⁰ Life Sciences Institute, National University of Singapore, Singapore, Singapore.
- ¹⁰¹ Department of Statistics and Applied Probability, National University of Singapore, Singapore, Singapore.
- ¹⁰² Department of Biochemistry, Faculty of Medicine, Health Science Center, Kuwait University, Safat, Kuwait.
- ¹⁰³ Faculty of Medicine, Macau University of Science & Technology, Macau, China.
- ¹⁰⁴ Department of Pathology and Laboratory Medicine, The Robert Larner M.D. College of Medicine, University of Vermont, Burlington, VT, USA.
- ¹⁰⁵ Department of Biochemistry, The Robert Larner M.D. College of Medicine, University of Vermont, Burlington, VT, USA.
- ¹⁰⁶ Folkhälsan Research Centre, Helsinki, Finland.
- ¹⁰⁷ Department of Endocrinology, Abdominal Centre, Helsinki University Hospital, Helsinki, Finland.
- ¹⁰⁸ Research Programs Unit, Clinical and Molecular Medicine, University of Helsinki, Helsinki, Finland.
- ¹⁰⁹ Diabetes Prevention Unit, National Institute for Health and Welfare, Helsinki, Finland.
- ¹¹⁰ Center for Vascular Prevention, Danube University Krems, Krems, Austria.
- ¹¹¹ Diabetes Research Group, King Abdulaziz University, Jeddah, Saudi Arabia.
- ¹¹² Instituto de Investigacion Sanitaria del Hospital Universitario LaPaz (IdiPAZ), University Hospital LaPaz, Autonomous University of Madrid, Madrid, Spain.
- ¹¹³ Unidad de Biología Molecular y Medicina Genómica, Instituto Nacional de Ciencias Médicas y Nutrición Salvador Zubirán, Mexico City, Mexico.
- ¹¹⁴ Departamento de Medicina Genómica y Toxicología Ambiental, Instituto de Investigaciones Biomédicas, UNAM, Mexico City, Mexico.
- ¹¹⁵ Department of Nutrition, Harvard School of Public Health, Boston, MA, USA.
- ¹¹⁶ Preventive Medicine & Epidemiology, and Cardiovascular Medicine, Medicine, Boston University School of Medicine, and Epidemiology, Boston University School of Public Health, Boston, MA, USA.
- ¹¹⁷ Department of Physiology and Biophysics, University of Mississippi Medical Center, Jackson, MS, USA
- ¹¹⁸ Department of Public Health, Aarhus University, Aarhus, Denmark.
- ¹¹⁹ Danish Diabetes Academy, Odense, Denmark.

¹²⁰ Oxford Centre for Diabetes, Endocrinology and Metabolism, Radcliffe Department of Medicine, University of Oxford, Oxford, UK.

¹²¹ Wellcome Centre for Human Genetics, Nuffield Department of Medicine, University of Oxford, Oxford, UK.

¹²² Present address: Genentech, 1 DNA Way, South San Francisco, CA, USA.

¹²³ Division of Genetics and Genomics, Boston Children's Hospital, Boston, Massachusetts, USA

¹²⁴ Department of Pediatrics, Harvard Medical School, Boston, MA, USA.

¹²⁵ Diabetes Unit and Center for Genomic Medicine, Massachusetts General Hospital, Boston, MA, USA.

¹²⁶ Department of Medicine, Harvard Medical School, Boston, MA, USA.

¹²⁷ Centre for Population Genomics, Garvan Institute of Medical Research, and UNSW Sydney, Sydney, New South Wales, Australia

¹²⁸ Centre for Population Genomics, Murdoch Children's Research Institute, Melbourne, Victoria, Australia

Corresponding author:

Miriam S. Udler

mudler@mgh.harvard.edu

List of Tables

Table 1 - Impact of clinically significant variants on traits.

List of Supplementary Tables

Supplementary Table 1 - Summary characteristics of study populations.

Supplementary Table 2 - Detailed characteristics of study cohorts. AMP-T2D-GENES Cohort Information.

Supplementary Table 3 - Counts of clinically significant variants and carriers across conditions.

Supplementary Table 4 - Effect size and penetrance estimates of clinically significant monogenic variants across conditions and by gene.

Supplementary Table 5 - Effect size and penetrance in filtered out variants (designated benign, likely benign, or not pLOF)

Supplementary Table 6 – Comparison of Top 1% gePS with interquartile range and monogenic carriers.

Supplementary Table 7 - Mean serum LDL values based on ascertainment.

Supplementary Table 8 - Clinical characteristics of monogenic diabetes variant carriers.

Supplementary Table 9 - Impact of polygenic score on trait expressivity in monogenic carriers.

Supplementary Table 10 - Frequency cut-offs used for ClinVar variant curation.

Supplementary Table 11 - pLoF curation categories and guidelines for classification.

Supplementary Table 12 - Variant curation assessments and carrier counts.

List of Supplementary Figures

Supplementary Figure 1. Distribution of clinically significant variants across ancestries.

Supplementary Figure 2. Carriers of clinically significant variants in MODY genes show a younger age of diabetes diagnosis compared the rest of the AMP-T2D-GENES cohorts and the UK Biobank population.

Abstract

Hundreds of thousands of genetic variants have been reported to cause severe monogenic diseases, but the probability that a variant carrier will develop the disease (termed penetrance) is unknown for virtually all of them. Additionally, the clinical utility of common polygenic variation remains uncertain. Using exome sequencing from 77,184 adult individuals (38,618 multi-ancestral individuals from a type 2 diabetes case-control study and 38,566 participants from the UK Biobank, for whom genotype array data were also available), we applied clinical standard-of-care gene variant curation for eight monogenic metabolic conditions. Rare variants causing monogenic diabetes and dyslipidemias displayed effect sizes significantly larger than the top 1% of the corresponding polygenic scores. Nevertheless, penetrance estimates for monogenic variant carriers averaged below 60% in both studies for all conditions except monogenic diabetes. We assessed additional epidemiologic and genetic factors contributing to risk prediction, demonstrating that inclusion of common polygenic variation significantly improved biomarker estimation for two monogenic dyslipidemias.

Introduction

Healthcare providers and researchers are increasingly faced with interpreting genetic sequence data collected from individuals who are asymptomatic or for whom limited clinical information is available. Standard clinical practice for reporting whole exome and genome sequencing results may involve risk assessment for genetic variation causing conditions of known relevance to the individual and also potentially impactful variants unrelated to the primary indication for testing (termed “secondary genetic findings,” for example the American College of Medical Genetics and Genomics (ACMG) list of 59 medically actionable genes¹⁻⁴). Thus, predicting the risk conferred by genetic findings in individuals who are not known to have the relevant conditions is of critical importance, but remains a challenge⁵. Furthermore, the scope of genetic variation interpreted in current clinical genetics practice is predominantly limited to rare monogenic “Mendelian” disease variants with large predicted effect sizes, leaving the vast majority of the genome, including common variants, unassessed. Recent studies have suggested that a high burden of common genetic variation may confer increased disease risk equivalent in magnitude to carrying rare monogenic variants⁶; however, this equivalency has recently been called into question⁷, and it remains uncertain whether and how to integrate polygenic scores capturing common genetic variation into medical care⁸.

Clinical application of genomic sequence data requires identification of medically significant genetic variants and estimation of their impact. In recent years, detailed guidelines from the ACMG and the Association for Molecular Pathology (AMP)⁹ have provided standards for reporting clinically significant variants which have been implemented by approximately 95% of clinical laboratories internationally¹⁰. However, the probability that carriers of such variants will manifest the given condition (termed “penetrance”) is unknown or uncertain for the vast majority of reported pathogenic variants⁵. Furthermore, individuals with the same genotype may exhibit variable degrees of phenotype expression (termed “variable expressivity”) ^{11,12}. Estimates of penetrance and expressivity traditionally have been derived from studies focusing on individuals with a given condition and their family members; this approach suffers from ascertainment bias, since the proband, who came to clinical attention due to having the condition, may share other genetic and/or environmental factors influencing manifestation of the condition with their family members^{12,13}. Interpretation of rare variants identified by sequencing is further complicated by limited or no data available from any source, including families, to assess penetrance⁵.

Large-scale population-based and cohort studies with both sequence and phenotype data offer an opportunity to estimate penetrance and expressivity with less upward bias compared to family or case-control studies. In fact, population-based studies may have a healthy-participant bias which could provide downwardly biased estimates of penetrance¹⁴. Recent studies attempting to connect large-scale genetic and phenotypic data have noted reduced penetrance estimates compared to those previously reported; however, these recent studies were limited by sample size and/or application of less stringent curation of genetic variants than the current clinical standard of care ACMG/AMP guideline approach^{7,14-17}. Additionally, further characterization of additional epidemiologic and genetic factors, such as phenotypic ascertainment and polygenic risk, is needed for accurate prediction of penetrance and expressivity for rare monogenic variants.

Here we present analyses performed in two separate datasets: 38,618 exomes from individuals ascertained as part of multi-ancestral type 2 diabetes (T2D) case-control studies, and 38,566 exomes from individual volunteers in the UK Biobank (UKB). Our analyses focused on traits with complex genetic architectures, involving rare and common genetic contribution, and well-defined biomarkers. These included diabetes (maturity-onset diabetes of the young (MODY), neonatal diabetes, autosomal dominant lipodystrophy) and disorders of LDL cholesterol, HDL cholesterol, triglycerides, and obesity. In addition to performing stringent curation using the ACMG/AMP criteria⁹ to generate a set of clinically significant genetic variants, we have also calculated polygenic scores in the UKB dataset to assess the cumulative impact of common variation on the same phenotypes. These data allow us to make a direct comparison between monogenic and polygenic risk, and to assess the contribution of polygenic risk to expressivity for carriers of monogenic variants.

Results

Identification of high confidence clinically significant variants enhances risk stratification

We studied two distinct datasets for which both individual-level exome sequence and phenotypic data were available (N=77,184): a compilation of multi-ancestral case control-studies for T2D, involving 22,875 T2D (or prediabetes) cases (see **Methods**) and 15,743 controls from the T2D-GENES and AMP-T2D consortia¹⁸, (referred to subsequently as AMP-T2D-GENES); and 38,566 unrelated individuals of European origin from the UKB¹⁹ (see **Methods, Supplementary Table 1**). Our analyses focused on 26 genes offered by clinical laboratories the United States for evaluation of monogenic forms of diabetes or diabetes-related traits through autosomal dominant modes of inheritance: MODY most commonly offered in panel testing (*GCK*, *HNF1A*, *HNF1B*, *HNF4A*, *PDX1*), an extended set of purported MODY genes less frequently offered in panel testing (*AKT2*, *KLF11*, *APPL1*, *ABCC8*, *KCNJ11*, *NEUROD1*, *CEL*, *INS*), neonatal diabetes (*ABCC8*, *GATA4*, *GATA6*, *HNF1B*, *INS*, *KCNJ11*), lipodystrophy (*AKT2*, *LMNA*, *PLIN1*, *PPARG*), elevated LDL cholesterol (*LDLR*, *APOB*), low serum LDL cholesterol (*APOB*, *PCSK9*), elevated serum HDL cholesterol (*CETP*), hypertriglyceridemia (*APOA5*, *LPL*), and monogenic obesity (*MC4R*).

We performed stringent variant curation using the clinical gold standard ACMG/AMP criteria, blinded to carrier phenotypic data for two classes of variants: 276 variants previously reported to be clinically significant (designated “pathogenic” or “likely pathogenic”) in the ClinVar database²⁰ or designated as disease-causing in review articles²¹⁻²³; and 218 predicted loss of function (pLoF) variants in genes with supported loss-of-function mechanism of action (see **Methods**). Our approach was intended to capture high-confidence clinically significant variants, although notably excluded missense

variants beyond those in the ClinVar database because of the low prior probability of disease relevance and the challenges of inferring pathogenicity for this variant class. In total across the AMP-T2D-GENES and UKB study exomes, 238 variants, representing 51% of all 463 variants curated, were determined by ACMG/AMP criteria to be clinically significant and were found in 626 carriers (**Figure 1, Supplementary Table 3, Supplementary Table 12**). Across the conditions, the clinically significant variants were observed in all represented ancestral groups (**Supplementary Figure 1**).

We next assessed the impact of clinically significant monogenic variants on corresponding biomarkers, restricting analyses to conditions with at least ten carriers of variants in relevant genes (**Supplementary Table 3**). Monogenic variant carriers for dyslipidemias had significantly more extreme lipid trait values compared to non-carriers, with shifts of ~55 mg/dL for both high and low LDL cholesterol conditions, ~130 mg/dL for high triglycerides, and ~16 mg/dL for high HDL cholesterol ($P < 10^{-5}$ for all; adjusted for age, sex, and 10 PCs; **Table 1**). For monogenic obesity, *MC4R* variant carriers had ~2 kg/m² higher mean body mass index (BMI) than non-carriers in both datasets, however this difference reached significance only in UKB ($P = 0.063$ AMP-T2D-GENES, $P = 0.006$ UKB). Despite differences in the study populations and designs in AMP-T2D-GENES and UKB, the effect sizes of clinically significant variants on relevant biomarkers were remarkably consistent across the two studies for dyslipidemia and obesity gene sets, once the former was adjusted for lipid medication use (**Table 1, Supplementary Table 4**). MODY variant carriers had significantly increased odds of having diabetes compared to non-carriers in both studies ($OR > 7$, $P < 10^{-4}$; **Table 1, Supplementary Table 4**); differences in risk estimates between the two studies were likely influenced by ascertainment practices in AMP-T2D-GENES, as it was a T2D case-control study and several sub-studies intentionally excluded diabetes cases with clinical features suggestive of MODY¹⁸ (**Supplementary Table 2**).

We also performed the same effect size estimates noted above, but for the variants filtered out during our curation process. We reclassified 7% (21/276) of curated variants from review articles and from ClinVar (which had been designated as clinically significant by at least one submitting source) to “benign” or “likely benign.” Likewise, 27% (59/218) of the pLoF variants were downgraded by our manual review of sequence reads. Together, these ClinVar, review, and pLoF variants which were downgraded by our curation (77/463, 17%) had marked reduced effect sizes compared to variants we curated as clinically significant (**Supplementary Table 5**)²⁴⁻²⁷. These findings support our curation process and highlight the need for caution in relying on available variant designations without additional review.

Monogenic variant effect sizes are significantly larger than the top 1% of polygenic risk scores

We next directly compared the effect of monogenic variation to common genetic variation on the same corresponding biomarkers in UKB participants. We employed published polygenic scores capturing millions of common genetic variants across the genome, termed global extended polygenic scores (gePS)²⁸ (see **Methods**). Since the gePS predicts lifetime risk of developing a disease, and the population mean age in UKB was 58 years, it was possible that estimates by gePS would be underestimates not capturing individuals who would later in life develop a given condition. We therefore performed gePS analyses restricted to individuals age ≥ 60 year (mean age 65 years) so as to have a fairer comparison with monogenic conditions, which are typically diagnosed at a younger age.

Individuals with the top 1% of gePS had more extreme lipid levels or diabetes risk compared to those with average gePS (25-75%tiles) (**Supplementary Table 6**); however the carriers of clinically

significant monogenic variants for these same conditions had even more severe values compared to those top 1% respective gePS's ($P < 0.05$ for each condition, **Figure 2, Supplementary Table 6**). For obesity, the difference in BMI between *MC4R* monogenic variant carriers and the top 1% BMI gePS was not significant (**Figure 2**).

Monogenic metabolic conditions display highly variable penetrance estimates

While in aggregate clinically significant monogenic variants had marked effect sizes, individual-level trait values in carriers varied considerably (**Figure 3A**). In both datasets, penetrance estimates based on standard disease cut-offs were estimated to be less than 60% for all monogenic metabolic conditions except elevated HDL cholesterol and monogenic diabetes (**Figure 3B, Supplementary Table 2**). Penetrance estimates using composite gene sets for conditions ranged from no more than 20% for low LDL cholesterol conditions caused by *APOB* and *PCSK9* variants (20.0%, 95% CI 8.4-36.9% in AMP-T2D-GENES, 5.6%, 95% CI 1.8-12.5% in UKB) to greater than 80% for diabetes or prediabetes caused by MODY genes *GCK*, *HNF1A*, *HNF4A*, *HNF1B*, and *PDX1* (86.4%, 95% CI's 65.1-97.1% in AMP-T2D-GENES, 81.2%, 54.4-96% in UKB) (**Figure 3B, Supplementary Table 4**).

Genetic vs phenotypic ascertainment of MODY suggests broad phenotypic spectrum

We performed deeper phenotyping of MODY variant carriers in the two datasets to determine whether these genetically ascertained individuals manifested clinical features suggestive of MODY, as typically seen in phenotypically ascertained MODY cases. Monogenic diabetes, and particularly MODY (the most common form) can often be misdiagnosed as type 1 or type 2; however, MODY has subtle phenotypic differences from these other forms of diabetes and also, importantly, distinct gene-specific therapeutic strategies²⁹.

Focusing on the MODY genes most commonly offered in commercial panels available in the United States (*HNF1A*, *GCK*, *HNF4A*, *HNF1B*, and *PDX1*)³⁰, more than 80% of carriers of clinically significant variants had evidence of prediabetes or diabetes (81.2% (95% CI 54.4-96.0%) in UKB and 86.4% (95% CI 65.1-97.1%) in AMP-T2D-GENES) (**Supplementary Table 8, Supplementary Figure 2**). Notably clinical features classically associated with MODY (BMI ≤ 30 and triglycerides ≤ 150 ^{31,32}) were only observed in 50% (11/22) of MODY variant-carrying individuals in AMP-T2D-GENES and 75% (12/16) in UKB. Similarly, an expected young age of diagnosis (age ≤ 35 years), was only observed in 21% (3/14) of those with available data across both datasets (**Supplementary Table 8**). Thus, at least 63% of all MODY variant carriers did not have expected clinical features. Since participants in AMP-T2D-GENES were selected to be T2D cases or controls, and specific exclusion practices employed by several studies to remove possible monogenic diabetes cases (**Supplementary Table 2**)¹⁸, such ascertainment practices could certainly introduce bias away from classical MODY features. Nevertheless, when all MODY carriers were compared to others with diabetes in each study, they had significantly lower mean BMI and serum triglycerides (BMI: AMP-T2D-GENES: 26.6 vs 28.7 kg/m², $P=0.027$; UKB: 25.8 vs 31.7 kg/m², $P=0.004$; triglycerides: AMP-T2D-GENES: 136 vs 182 mg/dL, $P=0.032$; UKB: 97 vs 186 mg/dL, $P=0.004$; adjusted for age, sex, and 10 PCs). Thus, in aggregate, MODY variant carriers displayed expected clinical features, but on an individual level, genetically ascertained individuals revealed a broader spectrum of disease phenotype.

We also identified gene-specific findings for the two most common MODY genes, *GCK* and *HNF1A*. *GCK*-MODY is characterized by non-progressive asymptomatic mild hyperglycemia that is present from birth and may remain in the prediabetes state rather than progress to diabetes³³. In both datasets 100% (17/17) of carriers of clinically significant *GCK* variants developed diabetes or prediabetes (penetrance estimates of 100%, 95% CI: 59.0-100% in AMP-T2D-GENES and 69.2-100% in UKB). All those with glycated hemoglobin (HbA1c) values available (N=13) had levels consistent with *GCK*-MODY, ranging from 5.7 to 7.2% (HbA1c in *GCK*-MODY is typically 5.6-7.6%³⁴) (**Supplementary Table 8**). Penetrance estimates for diabetes in *HNF1A*-MODY from our two datasets (81% in AMP-T2D-GENES, 95% CI 48.2-97.7 and 40% in UKB, 95% CI 5.27-85.3 diagnosed with diabetes by 55 years) were lower than what has previously been reported in the literature (e.g. 97%, 95% CI 96–98 by 50 years³⁵) (**Supplementary Table 8**).

Phenotypic ascertainment strongly impacts estimates of expressivity

It is well-appreciated that phenotypic ascertainment of individuals can upwardly bias estimates of expressivity^{14,36}, and we sought to better define this impact by studying conditions of high and low LDL cholesterol levels, where we had information on phenotypic ascertainment within a specific AMP-T2D-GENES cohort. A set of 535 individuals selected for extreme LDL cholesterol (>98th or <2nd percentile), without knowledge of their monogenic condition carrier status, were sequenced as part of the Exome Sequencing Project (ESP) cohort in AMP-T2D-GENES³⁷ and not included in the prior analyses. Within this ascertained sample, we identified 18 carriers of clinically significant monogenic high LDL cholesterol variants in *APOB* and *LDLR* (mean LDL 329 mg/dL) and 15 carriers in low LDL cholesterol variants in *APOB* and *PCSK9* (mean LDL 49.2 mg/dL). As expected, compared to carriers of variants for the same LDL cholesterol conditions but not ascertained on LDL phenotype, the two ascertained groups had more extreme LDL cholesterol levels (mean LDL cholesterol values 198 mg/dL, $P=4 \times 10^{-4}$ and 77 mg/dL, $P=0.06$ respectively, **Figure 4, Supplementary Table 7**).

Five variants (High LDL: *LDLR* p.Glu101Lys, *LDLR* p.Asp266Glu, *LDLR* p.Gly592Glu, *APOB* p.Arg3527Gln; Low LDL: *PCSK9* p.Tyr142Ter) were carried by individuals both in the phenotypically ascertained group and in the rest of the AMP-T2D-GENES cohort. These variants showed the same pattern of significantly more extreme LDL cholesterol values in the phenotypically ascertained compared to genetically ascertained individuals ($P<0.05$; all analyses adjusted for age, sex, ancestry, and diabetes status; **Figure 4; Supplementary Table 7**). These marked differences in LDL cholesterol values between the phenotypic vs genetic ascertained carriers, even among those carrying exactly the same LDL cholesterol variant, could not be explained by the use of lipid-lowering medication, assay use, or biased selection of the LDL cholesterol values among those available (e.g. selection of maximum LDL cholesterol value ever for phenotypically ascertained participants)³⁷.

In fact, the mean absolute impact of phenotype ascertainment on serum LDL cholesterol levels among individuals with monogenic LDL-raising or lowering variants (27.8-131.0 mg/dL, **Supplementary Table 7, Figure 4**) was thus similar or greater than the mean impact of carrying these same variants compared to non-carriers (31.5-65.3 mg/dL, **Table 1, Figure 4**). Such a substantial effect from phenotypic ascertainment reflects the large variation in expressivity at the single variant level and underscores the importance of considering phenotypic ascertainment bias in monogenic risk prediction.

Polygenic risk may increase expressivity of monogenic variants

The substantial variability in phenotype expressivity that we observed across all monogenic conditions (**Figure 3A**) suggests that additional environmental and/or genetic factors contribute to expressivity beyond the given monogenic variant. We assessed whether common genetic variation alters expressivity in UKB participants carrying monogenic disease variants.

Among carriers of high HDL cholesterol, low LDL cholesterol, high triglycerides, and monogenic obesity variants, we found that a higher gePS for each condition was associated with a more severe phenotype (e.g., among carriers of monogenic high HDL cholesterol variants, having an increased HDL gePS was associated with even higher HDL cholesterol). However, these trends were only significant for high HDL cholesterol (gePS one SD: beta 17.52 mg/dL, $P=0.012$) and high triglycerides (gePS one SD: beta 80.57 mg/dL, $P=0.014$) (**Figure 5, Supplementary Table 9**). Notably, despite our large study size, power in this analysis was limited, and we estimate that at least 98 carriers of clinically significant variants for a given monogenic condition would be needed for 80% power to detect a correlation of 0.25 (the minimum noted for the above traits) between a given trait and gePS at significance level $\alpha=0.05$. Therefore, for a monogenic condition with prevalence of 1 in 10,000 individuals, a population-based study with sample size on the order of one million individuals would be required to categorically determine the impact of polygenic risk.

We also assessed the interaction between gePS and monogenic risk in both monogenic carriers and non-carriers in the UKB, and observed significant positive interactions for the same two conditions, high HDL cholesterol ($P=0.001$) and high triglyceride levels ($P=0.01$); however, given the complexities of interaction analyses, additional work will also be needed in larger cohorts before we can conclude that gePS contributes to phenotype expression differently in carriers and non-carriers³⁸.

Discussion

Until recently, the impact of clinically significant monogenic variants on predicting phenotype expression has been predominantly studied in individuals or families ascertained on phenotype¹³. Our analysis employed population-based studies to provide less biased estimates of penetrance and expressivity, and to quantify the impact of phenotypic ascertainment and polygenic risk. We were able to directly compare monogenic and polygenic risk for each condition, and also assess the additional contribution of polygenic risk to expressivity for carriers of monogenic variants.

We applied the current gold standard ACMG/AMP clinical variant classification criteria⁹ to ensure relevance to current clinical practice and demonstrated resultant improvement in risk stratification. Gene variant curation was blinded to participant phenotypes and assessed the pathogenicity of variants expected to cause multiple metabolic conditions in 77,184 exomes of adults (age ≥ 40 years) from the AMP-T2D-GENES consortium and the UK Biobank. Our current analysis adds to a growing set of studies aimed at re-evaluating penetrance estimates using population-based studies^{7,9,14-17,39}; however, to the best of our knowledge, this work represents the largest study utilizing clinical standard-of-care ACMG/AMP criteria to curate gene variants in order to establish the penetrance and expressivity of all these monogenic conditions.

Carriers of the highly curated clinically significant variants for MODY and monogenic dyslipidemias had significantly more extreme trait effect sizes compared to non-carriers (OR > 7 for diabetes risk, betas 16.5-130.0 mg/dL for dyslipidemias, P values $< 10^{-4}$, **Table 1**). Despite differences in study populations and designs, the effect estimates for rare monogenic variation for all conditions

aside from monogenic diabetes (which was subject to ascertainment bias in AMP-T2D-GENES) were remarkably consistent between the two studies, supporting the integrity of our variant curation. We also assessed the impact of common genetic variation with polygenic scores. There has recently been a great deal of interest around the potential clinical contribution of such scores, especially gePS, and particularly in comparison to monogenic variant risk⁶. We show here that with the exception of monogenic obesity, polygenic risk at the top 1% of the risk distribution is not equivalent to monogenic risk, consistent with recent observations,⁷ but in contrast with others.⁴⁰ There will likely be further development of polygenic scores with improved disease prediction in the coming years⁴¹; however, in their current state and for the conditions we studied, the risk conferred by polygenic scores on their own was still substantially less than clinically significant monogenic variants; the only exception to this was *MC4R* obesity variants, which are known to have low predictive value for obesity risk⁴².

We observed a wide range of expressivity among clinically significant monogenic variant carriers across all traits (**Figure 2A**), and consequently estimates of penetrance were below 60% for all conditions except elevated HDL cholesterol and monogenic diabetes. At the single gene-level, the penetrance of *MC4R* for obesity (BMI ≥ 30 kg/m²) was less than 55%, consistent with previous findings^{7,42,43} (**Figure 2B**), while the penetrance of GCK-MODY was 100% for diabetes or prediabetes (95% CI's 59-100% in AMP-T2D-GENES, 69-100% in UKB). The range of penetrance estimates across genes and conditions may relate to ability to measure the direct biomarker(s) impacted by a given gene, the extent to which there are redundant mechanisms available in a given pathway to overcome a genetic defect⁴⁴, and the extent to which additional factors, such as other genetic and environmental factors (e.g. diet), impact the trait¹². The finding of 100% penetrance for diabetes or prediabetes seen in the 17 carriers of GCK-MODY across both datasets is particularly intriguing. *GCK* encodes glucokinase, which acts as the cell's glucose sensor as it facilitates phosphorylation of glucose to glucose-6-phosphate in the pancreatic beta cell, which is the first and rate-limiting step in glucose metabolism⁴⁵. The complete penetrance we have observed may be due to the ability to directly measure glucose as a relevant biomarker, as well as the essential role of *GCK* in glucose homeostasis, with suspected non-redundancy in functioning as a glucose sensor⁴⁵.

We also characterized the impact of phenotypic ascertainment bias on expressivity of clinically significant variants, showing that in individuals with the same LDL cholesterol-raising or -lowering variants there were significant differences in biomarker levels depending on the mode of ascertainment (genetic vs phenotypic) (**Figure 3**) and that the magnitude of this difference on LDL cholesterol levels (29-129 mg/dL) was similar or greater than the mean effect size of such variants (31.5-65.3 mg/dL, **Table 1**). This substantial impact of ascertainment bias was seen at the individual variant level, consistent with other similar observations of LDL cholesterol levels in *LDLR* and *APOB* carriers in a different study population^{36,46} and *HNF4A* p.Arg114Trp in diabetes risk¹⁴ (*HNF4A* p.Arg114Trp was present in the present datasets, but filtered out due to its designation as a variant of uncertain significant (VUS), reflecting its known low penetrance). The extent of ascertainment bias that we and others have identified highlights an important genetic counseling consideration, particularly with respect to interpretations of genomic sequencing data with limited clinical context available: interpretation of the same test result will likely have different prognostic implications depending on whether the individual tested or family members carry the phenotype of interest (e.g. hyperlipidemia) vs if a variant is identified secondarily; a Bayesian framework that takes into account pre-test probability might therefore be useful⁴⁷. Additionally, the variable expressivity seen at the single-variant level in multiple instances further supports additional risk factor modulation from other genetic and environmental exposures.

With regard to additional genetic factors impacting expressivity, we assessed the impact of more common polygenic variation on carriers of monogenic variants and found significant contributions for both high HDL cholesterol and high triglyceride levels ($P < 0.05$). These results add to a growing body of research supporting a significant polygenic contribution to monogenic risk^{7,39,48,49}. As such analyses were restricted to carriers of monogenic variants, power was limited, and it will be important to investigate in even larger datasets. We estimate that for a monogenic condition with prevalence of 1 in 10,000 individuals, population-based analyses well-powered to capture the contribution of polygenic risk to individuals with the monogenic condition would require on the order of one million individuals.

One limitation of this study is that our selection of variants for curation did not include all possible missense variants, but rather was confined to those reported in ClinVar or subject area reviews. This approach was designed to streamline the variant curation process and restrict our analyses to highly-confident pathogenic variants, but also meant that we were unable to generate estimates of the prevalence of monogenic condition in the two datasets. As discussed previously, there is also the potential for residual bias within the datasets. In the case of AMP-T2D-GENES, ascertainment of participants could have impacted penetrance of monogenic diabetes and expressivity of the metabolic phenotypes (**Supplementary Table 2**). In the UKB, a healthy participant bias⁵⁰ might be expected to reduce estimates of penetrance. Furthermore, despite our large dataset of exomes, the likelihood of observing any specific rare pathogenic variant is still low; this raises the possibility of bias toward lower penetrance of clinically significant variants, since allele frequency is a major predictor of pathogenicity⁵¹, and rarer variants with potentially greater penetrance are less likely to be observed.

Strengths of this study include the large number of participants with both phenotype and exome data, and the strict variant curation methodology applied. Our analysis of 276 variants designated by ClinVar as pathogenic or likely pathogenic highlights the need for careful curation of variants in clinical practice, with 57% reclassified to “benign,” “likely benign,” or “variant of uncertain significance” with application of ACMG/AMP criteria (**Figure 1**). Of note, however, the ClinVar variants we curated included those submitted to the database before establishment of current standards for curation⁹. With time, we can expect that the ClinVar database will become a more reliable resource for ascertaining clinically significant variants, as more submitters utilize standardized curation practices and additionally as condition-specific standards and curation are provided by ClinGen Expert Panels, including the Monogenic Diabetes Expert Panel in which several of the co-authors participate⁵².

Our study emphasizes the critical need for careful interpretation of monogenic variation, highlighting the roles of variant curation, phenotypic ascertainment, and polygenic risk in the estimates of penetrance and expressivity. In the coming years, access to larger sequencing studies will allow assessment of increasingly rare variants; however, deep phenotyping of such datasets, for example information on medication use and age of disease onset, will to be needed in parallel to better define genetic risk estimates. Improved understanding of monogenic variant expressivity will also likely require broader incorporation of genetic variation across the allelic frequency spectrum and integration of environmental factors. Such advances will facilitate modeling of disease risk and ultimately guide individualized patient genetic counseling and management recommendations.

Acknowledgements

This work was supported by NIH/NIDDK U01 DK105554 to JCF. This research has been conducted using the UK Biobank Resource under application number 27892. MSU is supported by NIH/NIDDK

K23 DK114551. AODL was supported by NIH/NICHD K12 HD052896. MB is supported by NIH/NIDDK DK062370. JCF is also supported by NIH/NIDDK K24 DK110550.

Funding for GO ESP was provided by NHLBI grants RC2 HL-103010 (HeartGO), with exome sequencing was performed through NHLBI grants RC2 HL-102925 (BroadGO) and RC2 HL-102926 (SeattleGO). HeartGO components and their support include Atherosclerosis Risk in Communities (NHLBI contracts N01 HC-55015, N01 HC-55016, N01HC-55017, N01 HC-55018, N01 HC-55019, N01 HC-55020, and N01 HC-55021); Cardiovascular Health Study (NHLBI contracts HHSN268201200036C, HHSN268200800007C, N01HC55222, N01HC85079, N01HC85080, N01HC85081, N01HC85082, N01HC85083, and N01HC85086); and NHLBI grants U01HL080295, R01HL087652, R01HL105756, R01HL103612, and R01HL120393, with additional contribution from the National Institute of Neurological Disorders and Stroke. Additional support was provided through R01AG023629 from the National Institute on Aging. A full list of principal Cardiovascular Health Study investigators and institutions can be found at CHS-NHLBI.org; Coronary Artery Risk Development in Young Adults (NHLBI contracts N01-HC95095, N01-HC48047, N01-HC48048, N01-HC48049, and N01-HC48050); Framingham Heart Study (NHLBI contract N01-HC-25195 and grants NS17950, AG08122, and AG033193); Jackson Heart Study (NHLBI contracts N01 HC-95170, N01 HC-95171, and N01 HC-95172); Multi-Ethnic Study of Atherosclerosis (NHLBI contracts N01-HC-95159 through N01-HC-95169 and grant 024156).

Cardiovascular Health Study: This CHS research was supported by NHLBI contracts HHSN268201200036C, HHSN268200800007C, HHSN268201800001C, N01HC55222, N01HC85079, N01HC85080, N01HC85081, N01HC85082, N01HC85083, N01HC85086; and NHLBI grants U01HL080295, R01HL087652, R01HL105756, R01HL103612, R01HL120393, and U01HL130114 with additional contribution from the National Institute of Neurological Disorders and Stroke (NINDS). Additional support was provided through R01AG023629 from the National Institute on Aging (NIA). A full list of principal CHS investigators and institutions can be found at CHS-NHLBI.org. We gratefully acknowledge the Eunice Kennedy National Institute of Child Health and Human Development for support of the MODY variant curation through U24 HD093486 (to TIP). The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

The Jackson Heart Study (JHS) is supported and conducted in collaboration with Jackson State University (HHSN268201800013I), Tougaloo College (HHSN268201800014I), the Mississippi State Department of Health (HHSN268201800015I) and the University of Mississippi Medical Center (HHSN268201800010I, HHSN268201800011I and HHSN268201800012I) contracts from the National Heart, Lung, and Blood Institute (NHLBI) and the National Institute on Minority Health and Health Disparities (NIMHD). The authors also wish to thank the staffs and participants of the JHS.

Novo Nordisk Foundation Center for Basic Metabolic Research is an independent Research Center, based at the University of Copenhagen, Denmark and partially funded by an unconditional donation from the Novo Nordisk Foundation (www.cbmr.ku.dk) (Grant number NNF18CC0034900).

The LOLIPOP study is supported by the National Institute for Health Research (NIHR) Comprehensive Biomedical Research Centre Imperial College Healthcare NHS Trust, the NIHR Official Development Assistance (ODA, award 16/136/68), the European Union FP7 (EpiMigrant, 279143) and H2020 programs (iHealth-T2D, 643774). The views expressed are those of the author(s) and not necessarily those of the Imperial College Healthcare NHS Trust, the NHS, the NIHR or the Department of Health. We thank the participants and research staff who made the study possible. JC is supported by the Singapore Ministry of Health's National Medical Research Council under its Singapore Translational Research Investigator (STaR) Award (NMRC/STaR/0028/2017).

The TwinsUK study was funded by the Wellcome Trust and European Community's Seventh Framework Programme (FP7/2007-2013). The TwinsUK study also receives support from the National Institute for Health Research (NIHR)- funded BioResource, Clinical Research Facility and Biomedical Research Centre based at Guy's and St Thomas' NHS Foundation Trust in partnership with King's College London.

The views expressed in this article are those of the author(s) and not necessarily those of the NHS, the NIHR, or the Department of Health. MMcC has served on advisory panels for Pfizer, NovoNordisk and Zoe Global, has received honoraria from Merck, Pfizer, Novo Nordisk and Eli Lilly, and research funding from Abbvie, Astra Zeneca, Boehringer Ingelheim, Eli Lilly, Janssen, Merck, NovoNordisk, Pfizer, Roche, Sanofi Aventis, Servier, and Takeda. As of June 2019, MMcC is an employee of Genentech, and a holder of Roche stock. MMcC wishes to acknowledge NIDDK U01-DK105535 and Wellcome: 090532, 098381, 106130, 203141, 212259.

The San Antonio Mexican American Family Studies (SAMAFS) are supported by the following grants/institutes. The San Antonio Family Heart Study (SAFHS) and San Antonio Family Diabetes/Gallbladder Study (SAFDGS) were supported by U01DK085524, R01 HL0113323, P01 HL045222, R01 DK047482 and R01 DK053889. The Veterans Administration Genetic Epidemiology Study (VAGES) study was supported by a Veterans Administration Epidemiologic grant. The Family Investigation of Nephropathy and Diabetes - San Antonio (FIND-SA) study was supported by NIH grant U01DK57295. The SAMAFS research team acknowledges the contributions of late Dr. H. E. Abboud to the research activities of the SAMAFS.

The KARE cohort was supported by grants from Korea Centers for Disease Control and Prevention(4845–301, 4851–302, 4851–307) and intramural grants from the Korea National Institute of Health (2016-NI73001-00, 2019-NG-053-01).

RJFL is supported by the NIH (R01DK110113, R01DK107786, 1R01DK124097). NC is supported by a grant from the Canadian Institutes of Health Research (CIHR Fellowship). The Mount Sinai BioMe Biobank has been supported by The Andrea and Charles Bronfman Philanthropies and in part by Federal funds from the NHLBI and NHGRI (U01HG00638001; U01HG007417; X01HL134588).

The Framingham Heart Study (FHS) acknowledges the support of Contracts NO1-HC-25195, HHSN268201500001I and 75N92019D00031 from the National Heart, Lung and Blood Institute and grant supplement R01 HL092577-06S1 for this research. We also acknowledge the dedication of the

FHS study participants without whom this research would not be possible. Dr. Vasani is supported in part by the Evans Medical Foundation and the Jay and Louis Coffman Endowment from the Department of Medicine, Boston University School of Medicine.

Methods

Study populations and phenotype curation

AMP-T2D-GENES

The complete AMP-T2D-GENES cohort consists of 20,791 cases and 24,440 controls selected from multiple distinct multi-ancestry studies¹⁸. The present study includes a subset of 22,875 T2D or prediabetes and 15,743 controls from studies who consented for the data to be used in this analysis, which included Genetics of Type 2 Diabetes (GoT2D), the Exome Sequencing Project (ESP), Lubeck Foundation Centre for Applied Medical Genomics in Personalised Disease Prediction, Prevention and Care (LuCamp), Slim Initiative in Genomic Medicine for the Americas (SIGMA T2D), and T2D-GENES (Type 2 Diabetes Genetic Exploration by Next-generation sequencing in multi-Ethnic Samples). General study characteristics are provided in **Supplementary Tables 1** with more details, including exclusion criteria available in **Supplementary Table 2**, which is adapted from Flannick *et al.*, 2019¹⁸. All samples were approved for use by their home institution's institutional review board or ethics committee. Analysis of the data was approved by the Mass General Brigham (formerly Partners) institutional review board in Boston, Massachusetts were limited to those participants in each cohort with available DNA who consented to genetic studies.

Phenotype information related to diabetes status was collected by each case-control or cohort study, as previously described in Flannick *et al.* Additionally, we defined prediabetes as any individual with HbA1c $\geq 5.7\%$, fasting blood glucose ≥ 100 mg/dL, or oral glucose tolerance test (OGTT) 2 hour blood glucose ≥ 140 mg/dL. In individuals who were reported to be on lipid-lowering medication, serum LDL cholesterol and triglyceride levels were adjusted for statin use based on previous studies estimating the impact^{53,54}: we divided LDL by 0.7 and triglycerides by 0.85 as has been previously been implemented⁵⁵. Self-reported ancestry was used, as this was previously shown to correlate well with principal component analysis (PCA) defined ancestry and specific exceptions were dropped from analyses¹⁸. Analyses described below used a dataset restricted to individuals in the 'unrelated analysis set' (see Flannick *et al.* methods). To provide consistency with the UKB dataset, individuals younger than age 40 were also excluded. Individuals recruited to the Pakistan Genomic Resource cohort were excluded for all analyses involving lipid levels or BMI.

UK Biobank

UK Biobank (UKB) is a prospective cohort of approximately 500,000 recruited individuals from the general population aged 40–69 years in 2006–2010 from across the United Kingdom, with genotype, phenotype, and linked healthcare record data⁵⁶. All participants provided electronic informed consent at their initial visit. Analysis of the data was approved by the Mass General Brigham (formerly Partners) institutional review board in Boston, Massachusetts, and was performed under UK Biobank application 27892.

Direct LDL cholesterol (mmol/L), direct HDL cholesterol (mmol/L), triglyceride (mmol/L), BMI (kg/m²) (field codes: 30780, 30760, 30870, 21001) data were extracted for all individuals. Lipid measurements were converted from mmol/L to mg/dL. The mean for all visits was used in subsequent analyses. The 'Medication for cholesterol, blood pressure, diabetes, or take exogenous hormones' fields (6177 and 6153) was used to determine lipid-lowering medication, where an individual was considered to be on lipid-lowering medication if it was recorded at any of the visits. LDL and triglyceride values were adjusted for use of lipid-lowering medication, as described above.

Glycated hemoglobin (HbA1c; field code 30750) was taken as the maximum observed across visits. Since monogenic diabetes may be misdiagnosed as type 1 or type 2 diabetes, we used an inclusive definition of diabetes: possible and probable type 1 or type 2 diabetes was determined in a manner similar to previously described methods⁵⁷. We also considered individuals as having diabetes if they had ICD10 codes E10-E14 (fields: 41202), and recorded diabetes medication use (fields: 6177, 6153), diabetes ever diagnosed by a doctor (field: 2443), nurse interview codes indicating diabetes (fields: 1220 - any diabetes, 1222 - T1D, 1223 - T2D), or HbA1c $\geq 6.5\%$. Prediabetes was defined as

any individual with HbA1c \geq 5.7%. We also extracted data for the first recorded age of diabetes diagnosis (fields: 20009, 2976), age, and sex.

This dataset was filtered to only unrelated individuals with European ancestry to facilitate comparisons of biomarkers in analyses using polygenic risk scores. Filtering to unrelated individuals was done using the column 'used.in.pca.calculation' in the UKB genotype data sample QC document (ukb_sqc_v2.txt) as a proxy. This column indicates samples which UKB used in a principal component analysis (PCA), and this analysis was only performed on unrelated, high quality samples. To filter to European ancestry only, samples were first projected onto 1000 Genomes phase 3⁵⁸ PCA coordinate space. Then Aberrant R package⁵⁹ clustering was used to identify individuals falling within 1000 Genomes project EUR PC1 and PC2 limits ($\lambda=4.5$). Individuals that self-reported as non-European ethnicity were also filtered. There were 38,566 individuals remaining after all filtering and intersection with individuals that also have exome sequence data released in the first tranche (Category 170).

Generation of gene list

We sought to include genes that would be ordered in the United States in clinical practice to diagnose conditions of monogenic diabetes, lipodystrophy, obesity, and lipid disorders. We searched the Genetic Testing Registry (<https://www.ncbi.nlm.nih.gov/gtr/>) and Concert Genetics (<https://app.concertgenetics.com/>), last accessed March 14th, 2018, for lists of available commercial gene panels for clinical genetic testing for these diseases available in the United States. We filtered this list of genes to those with an autosomal dominant mode of inheritance, as determined by the Online Mendelian Inheritance in Man® (OMIM, <https://www.omim.org/>). For the genes in OMIM where mode of inheritance was not specified, the genes were researched in ClinVar (<https://www.ncbi.nlm.nih.gov/clinvar/>) and related literature. In total there were 26 autosomal dominant genes across the conditions. We further excluded any gene where there was no ClinVar submission (April 2019 ClinVar submission summary) of pathogenic or likely pathogenic for the phenotype of interest that also included clinical testing as a collection method, leaving 20 genes. We determined phenotype overlap by manual review of the "SubmittedPhenotypeInfo" and "ReportedPhenotypeInfo" fields in the submission summary where present and 'ExplanationOfInterpretation' or submitted PubMed articles when phenotype info was not reported in the other fields. For MODY, most commercial panels available in March, 2018 included *HNF1A*, *HNF4A*, *GCK*, *HNF1B*, and *PDX1*, with larger panels less widely available. We therefore separated the MODY genes into two categories: "MODY" including those five genes and "MODY extended" including eight additional genes.

Determination of genes with LoF mechanism

The pLoF curation was restricted to genes alleged to cause disease with a LoF mechanism based on reporting in ClinVar or a PubMed publication of an LoF variant in an individual with the phenotype of interest.

Two genes were determined to be related to both high LDL (familial hypercholesterolemia) and low LDL (familial hypobetalipoproteinemia): *APOB* and *PCSK9*. Gain-of-function missense mutations in both genes result in increased LDL levels, while LoF mutations cause lower LDL levels⁶⁰⁻⁶². Therefore, only missense ClinVar variants in *APOB* and *PCSK9* were assessed in the curation process for high LDL, and LoF variants were considered for low LDL.

Exome data variant filtering and annotation

All filtering and annotation described below was performed using Hail 0.2 (<https://hail.is>).

AMP-T2D-GENES

Exome sequencing and quality control were described previously¹⁸. We applied additional genotype filters to retain only high-quality genotypes: genotype quality \geq 20, depth \geq 10, and minor

allele balance > 0.25 for heterozygous genotypes. Variants were annotated using Ensembl's Variant Effect Predictor (VEP) v85⁶³ with the Loss-of-function Transcript Effect Estimator (LOFTEE) plugin⁶⁴. The dataset was then filtered to only variants with a consequence on any of the genes of interest. The filtered VCF was used in analyses described below that involve EPACTS.

We determined which variants in our dataset have been submitted to ClinVar by cross-referencing this filtered variant list with the ClinVar VCF (April 2019) (further curation described below). A list of predicted loss-of-function (pLoF) variants, including stop gained, frameshift or essential splice site (splice donor or splice acceptor), was generated by filtering to variants with a LOFTEE high-confidence (HC) annotation on any transcript. Finally, we used transcript expression-aware annotation⁶⁵ to add pext (proportion expression across transcripts) values for the worst consequence annotation to each variant for use in pLoF curation discussed below.

UK Biobank

UKB exome sequencing PLINK files were imported into Hail and all the same annotation described for AMP-T2D-GENES was added using appropriate files for genotype reference GRCh38 and VEP v95. In order to compare UKB variants to AMP-T2D-GENES variants we used Hail's liftover method to lift data from GRCh38 to GRCh37. Since the PLINK files do not contain genotype quality information that we can use for filtering low-quality genotypes, we downloaded the gVCFs for all variant carriers and determined which individuals genotypes were not high-quality (genotype quality \geq 20, depth \geq 10, and minor allele balance > 0.25 for heterozygous genotypes) and set each of these to missing in the VCF.

ClinVar variant curation

We identified individuals carrying variants in the genes of interest that had at least one 'pathogenic' or 'likely pathogenic' submission in ClinVar by a clinical testing lab for the relevant trait. To streamline variant curation we first generated a list of high confidence clinical genetic testing laboratories. Using the April 2019 release of the ClinVar submission summary, a lab was considered high confidence if it had submitted more than 15,000 variants to ClinVar and had updated its submission after 2017 when the most recent ACMG variant interpretation guidelines were published⁹. This resulted in eight labs: Invitae; GeneDx; Ambry Genetics; EGL Genetic Diagnostics; Eurofins Clinical Diagnostics; PreventionGenetics; Laboratory of Molecular Medicine, Partners Healthcare Personalized Medicine; Genetic Services Laboratory, University of Chicago; and Counsyl. Variants that were reported by any lab on this list since January 1st, 2017 were then accepted as having the pathogenicity reported by the lab.

These labs were further verified through manual curation. First, five variants from each lab that were also present in our study were chosen to be manually curated, so that the manual curation could be compared to the lab's analysis. Through this, we found no differences in curation results. Then, five variants from each lab were chosen at random through ClinVar – one Pathogenic, one Likely Pathogenic, one VUS, one Likely Benign, and one Benign. As PreventionGenetics only submitted Benign and Likely Benign to ClinVar, their variants were limited to those categories. These variants were then also manually curated, and the results were compared. The only difference in curation of the non-study variants involved University of Chicago, due to internal data initially not available to our study curator; however, the same conclusion was reached upon inclusion of this internal data, which was included in their reporting in ClinVar. During the manual phenotype curation (described below), we discovered Counsyl reported conflicting phenotypes for the same variant, so we opted to manually curate variants assessed by Counsyl.

The variants not analyzed by high confidence labs were analyzed separately using manual curation with the curator blinded to carrier phenotypes. The ClinGen Variant Curation Interface (<https://curation.clinicalgenome.org/>) was used to analyze the variants and assign evidence following the ACMG guidelines⁹ and recommendation for interpretation of LoF variants⁶⁶, with input from gene-

specific rules under development by the Monogenic Diabetes Expert Panel VCEP (<https://clinicalgenome.org/affiliation/50016/>) for the MODY variants. Databases and other resources such as ClinVar (<https://www.ncbi.nlm.nih.gov/clinvar/>), Human Gene Mutation Database (HGMD) (<https://digitalinsights.qiagen.com/products-overview/clinical-insights-portfolio/human-gene-mutation-database/>), gnomAD (<https://gnomad.broadinstitute.org/>), PubMed (<https://pubmed.ncbi.nlm.nih.gov/>), Google Scholar (<https://scholar.google.com/>), Alamut (<https://www.interactive-biosoftware.com/alamut-visual/>), and the UCSC browser (<https://genome.ucsc.edu/>) were utilized to collect evidence for curation purposes. The general guidelines were adjusted slightly for certain criteria such as control population frequency as shown in **Supplementary Table 10**. Since most AMP-T2D-GENES participants are included in gnomAD, AMP-T2D-GENES allele frequency decisions were made by subtracting the number of AMP-T2D-GENES carriers from the number of total gnomAD carriers to determine an adjusted gnomAD allele frequency which was compared to the cut-offs shown in **Supplementary Table 10**.

Three variants within *HNF1A* were excluded from further analysis because of poor genotyping quality at this site making it difficult to determine which individuals are actually carriers (GRCh37: 12-121432114-CG-C, 12-121432116-G-GC, 12-121432117-G-GC, GRCh38: chr12:120994311-CG-C, chr12:120994313-G-GC, chr12:120994314-G-GC). As all three are frameshifts, these variants were also excluded from the pLoF curation described below.

Variants in MODY genes were curated by a second set of reviewers at University of Maryland School of Medicine, the home institution of the ClinGen Monogenic Diabetes Expert Panel, to ensure accuracy. All variants were consistently classified as collectively pathogenic or likely pathogenic (**Supplementary Table 12**).

All variants curated for this project, along with their classification and supporting evidence, were submitted to ClinVar on January 30th, 2020.

High confidence Loss of Function variants

As described above, we used LOFTEE⁶⁴ to generate a list of high confidence pLoF variants, restricting to the set of genes we determined to have a LoF mechanism of pathogenicity. Each pLoF variant was assessed by manual review of reads by two independent reviewers. The reads were examined for poor quality, homopolymer artifacts, and multinucleotide variants (MNVs) causing a synonymous or missense variant instead of the reported stop codon. Where available, gnomAD data was examined to identify variants that were flagged as filtered by gnomAD's random forest variant quality control method. UCSC genome browser data was assessed to determine the conservation of the region, the location of the variant, and how many transcripts the variant was coding. If the variant was present in the last exon or last 50 base pairs of the penultimate exon, it was deemed not LoF due to a predicted lack of nonsense mediated decay. However, this was overruled if the variant was predicted to delete over 25% of the gene. The potential for a splice site rescue was assessed by examining +/- 21bp around the variant. Any inframe splice site within 6bp was considered an essential splice site rescue and possible inframe splice site rescues between 6 and 21bp were considered a rescue if validated by the alternative splice site prediction tool Alamut (v.2.11). We also used pext values obtained from the transcript expression-aware annotation⁶⁵ to indicate variants that fell in exons that have evidence of poor expression (specific cutoffs are detailed in **Supplementary Table 11**). Variants were classified into 5 categories, 'LoF', 'likely LoF', 'uncertain', 'likely not LoF', or 'not LoF' using the guidelines described in **Supplementary Table 11**. Any variant that had a discordant assessment between the two reviewers ('LoF' or 'likely LoF' by only one reviewer) was examined by a third reviewer to determine the final pLoF annotation.

Carrier vs non-carrier effect size analysis

We considered an individual to be a carrier of a clinically significant variant if they carry a ClinVar variant assessed as pathogenic or likely pathogenic or a pLoF variant passing manual curation

('LoF' or 'likely LoF' as described above). For AMP-T2D-GENES, as previously described¹⁸, we accounted for the diverse ancestry and different sequencing technologies by using a modified version of EFACTS (<http://genome.sph.umich.edu/wiki/EFACTS>) that sets specified variants to missing based on QC of sample subgroups (as described in Flannick et al, there are 25 subgroups that were determined by stratifying samples by cohort of origin, ancestry, and/or sequencing technology). As covariates in AMP-T2D-GENES analyses, we included sex, age, PCs 1-10, sample subgroup, and sequencing technology all as previously defined¹⁸. Analyses on UKB used covariates for sex, age, PCs 1-10 and the genotyping array.

For both AMP-T2D-GENES and UKB, we used VCFs produced after filtering variants as described above and performed the group b.burdenFirth for binary traits and q.burdenTest for continuous traits in EFACTS to compare carriers and non carriers for the following condition/phenotype pairs: high LDL cholesterol with LDL cholesterol (mg/dL); low LDL cholesterol with LDL cholesterol (mg/dL); high HDL cholesterol with HDL cholesterol (mg/dL); high triglycerides with triglycerides (mg/dL); monogenic obesity with BMI (kg/m²), MODY with diabetes status, and in diabetes cases only: HDL cholesterol, Triglycerides, and BMI.

Additionally, we included T2D or T2D with prediabetes as covariates in all tests on lipid measurements and BMI. Triglycerides and BMI were log transformed. All of these analyses were also performed per gene to ensure that we captured possible gene level differences in phenotype values.

Estimation of penetrance

Unlike diabetes, phenotypes used to assess the possibility that individuals have each lipid condition or obesity, are continuous. There are commonly used clinical guidelines for diagnosis of each of these conditions, so these cutoffs were used to dichotomize the phenotypes allowing us to determine individuals with and without each condition for use in estimating penetrance. The following clinical diagnosis cutoffs were used: High LDL cholesterol: LDL cholesterol ≥ 190 mg/dL, Low LDL cholesterol: LDL cholesterol ≤ 50 mg/dL, High HDL cholesterol: HDL cholesterol ≥ 60 mg/dL, High triglycerides: triglycerides ≥ 200 mg/dL⁶⁷, and Monogenic obesity: BMI ≥ 30 kg/m².

Penetrance estimates were calculated as the proportion of individuals carrying a clinically significant variant that also exhibit the expected condition. To determine the significance for all penetrance estimates we used the group Firth burden test in the modified version of EFACTS and the same covariates as described in 'Carrier vs non-carrier enrichment analysis'.

Calculation of global extended polygenic score (gePS)

Body mass index and type 2 diabetes

Global extended polygenic scores for T2D and BMI were previously calculated on UKB participants using LDpred^{6,43}. The variants and weights used in the calculation were downloaded (<http://www.broadcvdi.org/informational/data>). These weights were then applied to the UKB genotype data from the subset of individuals included in this study to calculate a gePS using Hail's equivalent to the --score method in PLINK (<https://hail.is/docs/0.2/guides/genetics.html?highlight=prs>). These values were then scaled and centered around zero with a standard deviation of one for downstream analysis. We confirmed that plots of T2D prevalence and BMI by respective polygenic scores converged at the same upper limits as previously published^{6,43}.

Lipid conditions

To estimate a gePS for each lipid phenotype, we filtered UK Biobank genotype data to only the individuals used in this study (unrelated, EUR ancestry, and exome sequenced) and excluded SNPs with an imputation INFO < 0.3 and allele frequency $< 1\%$. Summary statistics for lipid GWAS were downloaded from the European Network for Genetic and Genomic Epidemiology (ENGAGE) Consortium. This included LDL cholesterol, HDL cholesterol, and triglyceride GWAS summary stats from a meta-analysis of up to 62,166 individuals of European ancestry⁶⁸. We filtered to variants observed in HapMap3 (--only-hm3) and both the summary statistics and genotype data, and then

estimated SNP weights using the Bayesian computational method LDpred (version 1.0.6) which accounts for local LD patterns⁶⁹. SNP weight estimates were obtained using the infinitesimal (inf) model (assumes all genetic variants impact phenotype) with heritability estimates (TG: 0.1525, LDL: 0.1347, HDL: 0.1572) as previously calculated using LD Score regression⁷⁰ and displayed on LD Hub⁷¹. We then used PLINK version 1.9 (--score) to calculate polygenic scores using the SNP weights⁷². As in the BMI and T2D gePRS, the distribution was scaled to have a mean of zero and one standard deviation around the mean. Since there is a single gePS for LDL cholesterol, the scaled gePS was multiplied by -1 for figures and analyses comparing low LDL cholesterol carrier phenotype values to phenotypes aggregated by gePS deciles or quantiles.

Statistical Analysis

We used generalized linear models (GLM) to examine the gePS results in a few different ways. We compared the top 1% to the interquartile range (25-75%) of the gePS and to the clinically significant variant carriers (**Supplementary Table 6**). For both analyses we restricted the age in controls to ≥ 60 . Additionally, we determine the effect size of gePS on phenotypes in the subset of only clinically significant variant carriers and assessed the interaction of carrier status and gePS (**Supplementary Table 9**). In all GLMs age, sex and 10 PC's were included in the model as covariates. A linear regression was performed for all phenotypes except diabetes where a logistic regression was applied.

All plots were made using R version 3.5.2.

Data availability

Sequence data and phenotypes for this study are available via the database of Genotypes and Phenotypes (dbGAP) and/or the European Genome-phenome Archive, as indicated in Supplementary Table 2.

Code availability

Code used in analyses is will be made available prior to publication in GitHub.

References

1. Green, R.C., *et al.* CORRIGENDUM: ACMG recommendations for reporting of incidental findings in clinical exome and genome sequencing. *Genet. Med.* **19**, 606 (2017).
2. Directors, A.B.o. ACMG policy statement: updated recommendations regarding analysis and reporting of secondary findings in clinical genome-scale sequencing. *Genet. Med.* **17**, 68-69 (2015).
3. Green, R.C., *et al.* ACMG recommendations for reporting of incidental findings in clinical exome and genome sequencing. *Genet. Med.* **15**, 565-574 (2013).
4. Kalia, S.S., *et al.* Recommendations for reporting of secondary findings in clinical exome and genome sequencing, 2016 update (ACMG SF v2.0): a policy statement of the American College of Medical Genetics and Genomics. *Genet. Med.* **19**, 249-255 (2017).
5. Directors, A.B.o. The use of ACMG secondary findings recommendations for general population screening: a policy statement of the American College of Medical Genetics and Genomics (ACMG). *Genet. Med.* **21**, 1467-1468 (2019).
6. Khera, A.V., *et al.* Genome-wide polygenic scores for common diseases identify individuals with risk equivalent to monogenic mutations. *Nat. Genet.* **50**, 1219-1224 (2018).
7. Oetjens, M.T., Kelly, M.A., Sturm, A.C., Martin, C.L. & Ledbetter, D.H. Quantifying the polygenic contribution to variable expressivity in eleven rare genetic disorders. *Nat. Commun.* **10**, 4897 (2019).
8. Lewis, C.M. & Vassos, E. Polygenic risk scores: from research tools to clinical instruments. *Genome Med* **12**, 44 (2020).

9. Richards, S., *et al.* Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genet. Med.* **17**, 405-424 (2015).
10. Niehaus, A., *et al.* A survey assessing adoption of the ACMG-AMP guidelines for interpreting sequence variants and identification of areas for continued improvement. *Genet. Med.* **21**, 1699-1701 (2019).
11. Zlotogora, J. Penetrance and expressivity in the molecular age. *Genetics in Medicine* **5**, 347-352 (2003).
12. Cooper, D.N., Krawczak, M., Polychronakos, C., Tyler-Smith, C. & Kehrer-Sawatzki, H. Where genotype is not predictive of phenotype: towards an understanding of the molecular basis of reduced penetrance in human inherited disease. *Human Genetics* **132**, 1077-1130 (2013).
13. Turner, H. & Jackson, L. Evidence for penetrance in patients without a family history of disease: a systematic review. *Eur. J. Hum. Genet.* (2020).
14. Wright, C.F., *et al.* Assessing the Pathogenicity, Penetrance, and Expressivity of Putative Disease-Causing Variants in a Population Setting. *Am. J. Hum. Genet.* **104**, 275-286 (2019).
15. Natarajan, P., *et al.* Aggregate penetrance of genomic variants for actionable disorders in European and African Americans. *Sci. Transl. Med.* **8**, 364ra151 (2016).
16. Abul-Husn, N.S., *et al.* Genetic identification of familial hypercholesterolemia within a single U.S. health care system. *Science* **354**(2016).
17. Flannick, J., *et al.* Assessing the phenotypic effects in the general population of rare variants in genes for a dominant Mendelian form of diabetes. *Nat. Genet.* **45**, 1380-1385 (2013).
18. Flannick, J., *et al.* Exome sequencing of 20,791 cases of type 2 diabetes and 24,440 controls. *Nature* **570**, 71-76 (2019).
19. Bycroft, C., *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203-209 (2018).
20. Harrison, S.M., *et al.* Using ClinVar as a Resource to Support Variant Interpretation. *Current Protocols in Human Genetics*, 8.16.11-18.16.23 (2016).
21. Ellard, S. & Colclough, K. Mutations in the genes encoding the transcription factors hepatocyte nuclear factor 1 alpha (HNF1A) and 4 alpha (HNF4A) in maturity-onset diabetes of the young. *Human Mutation* **27**, 854-869 (2006).
22. Osbak, K.K., *et al.* Update on mutations in glucokinase (GCK), which cause maturity-onset diabetes of the young, permanent neonatal diabetes, and hyperinsulinemic hypoglycemia. *Hum. Mutat.* **30**, 1512-1526 (2009).
23. Colclough, K., Bellanne-Chantelot, C., Saint-Martin, C., Flanagan, S.E. & Ellard, S. Mutations in the genes encoding the transcription factors hepatocyte nuclear factor 1 alpha and 4 alpha in maturity-onset diabetes of the young and hyperinsulinemic hypoglycemia. *Hum Mutat* **34**, 669-685 (2013).
24. Rehm, H.L., *et al.* ClinGen — The Clinical Genome Resource. *New England Journal of Medicine* **372**, 2235-2242 (2015).
25. Yang, S., *et al.* Sources of discordance among germ-line variant classifications in ClinVar. *Genet. Med.* **19**, 1118-1126 (2017).
26. Harrison, S.M., *et al.* Scaling resolution of variant classification differences in ClinVar between 41 clinical laboratories through an outlier approach. *Hum. Mutat.* **39**, 1641-1649 (2018).
27. Campuzano, O., *et al.* Reanalysis and reclassification of rare genetic variants associated with inherited arrhythmogenic syndromes. *EBioMedicine* **54**, 102732 (2020).
28. Udler, M.S., McCarthy, M.I., Florez, J.C. & Mahajan, A. Genetic Risk Scores for Diabetes Diagnosis and Precision Medicine. *Endocrine Reviews* **40**, 1500-1520 (2019).
29. Hattersley, A.T., *et al.* ISPAD Clinical Practice Consensus Guidelines 2018: The diagnosis and management of monogenic diabetes in children and adolescents. *Pediatr Diabetes* **19 Suppl 27**, 47-63 (2018).
30. Home - Genetic Testing Registry (GTR) - NCBI.

31. Naylor, R., Knight Johnson, A. & del Gaudio, D. Maturity-Onset Diabetes of the Young Overview. in *GeneReviews* (eds. Adam, M.P., *et al.*) (University of Washington, Seattle, Seattle (WA), 2018).
32. Fajans, S.S., Bell, G.I. & Polonsky, K.S. Molecular mechanisms and clinical pathophysiology of maturity-onset diabetes of the young. *N. Engl. J. Med.* **345**, 971-980 (2001).
33. Chakera, A.J., *et al.* Recognition and Management of Individuals With Hyperglycemia Because of a Heterozygous Glucokinase Mutation. *Diabetes Care* **38**, 1383-1392 (2015).
34. Steele, A.M., *et al.* Use of HbA1c in the identification of patients with hyperglycaemia caused by a glucokinase mutation: observational case control studies. *PLoS One* **8**, e65326 (2013).
35. Patel, K.A., *et al.* Heterozygous RFX6 protein truncating variants are associated with MODY with reduced penetrance. *Nat Commun* **8**, 888 (2017).
36. Tybjaerg-Hansen, A., *et al.* Phenotype of heterozygotes for low-density lipoprotein receptor mutations identified in different background populations. *Arterioscler. Thromb. Vasc. Biol.* **25**, 211-215 (2005).
37. Lange, L.A., *et al.* Whole-exome sequencing identifies rare and low-frequency coding variants associated with LDL cholesterol. *Am. J. Hum. Genet.* **94**, 233-245 (2014).
38. Aschard, H., *et al.* Challenges and opportunities in genome-wide environmental interaction (GWEI) studies. *Hum. Genet.* **131**, 1591-1613 (2012).
39. Fahed, A.C., *et al.* Polygenic background modifies penetrance of monogenic variants for tier 1 genomic conditions. *Nat Commun* **11**, 3635 (2020).
40. Khera, A.V., *et al.* Genome-wide polygenic scores for common diseases identify individuals with risk equivalent to monogenic mutations. *Nat Genet* **50**, 1219-1224 (2018).
41. Chatterjee, N., *et al.* Projecting the performance of risk prediction based on polygenic analyses of genome-wide association studies. *Nature Genetics* **45**, 400-405 (2013).
42. Stutzmann, F., *et al.* Prevalence of melanocortin-4 receptor deficiency in Europeans and their age-dependent penetrance in multigenerational pedigrees. *Diabetes* **57**, 2511-2518 (2008).
43. Khera, A.V., *et al.* Polygenic Prediction of Weight and Obesity Trajectories from Birth to Adulthood. *Cell* **177**, 587-596.e589 (2019).
44. Narasimhan, V.M., *et al.* Health and population effects of rare gene knockouts in adult humans with related parents. *Science* **352**, 474-477 (2016).
45. Matschinsky, F.M. & Wilson, D.F. The Central Role of Glucokinase in Glucose Homeostasis: A Perspective 50 Years After Demonstrating the Presence of the Enzyme in Islets of Langerhans. *Front. Physiol.* **10**, 148 (2019).
46. Tybjaerg-Hansen, A., Steffensen, R., Meinertz, H., Schnohr, P. & Nordestgaard, B.G. Association of mutations in the apolipoprotein B gene with hypercholesterolemia and the risk of ischemic heart disease. *N Engl J Med* **338**, 1577-1584 (1998).
47. Sorscher, S. Ascertainment Bias and Estimating Penetrance. *JAMA Oncol* **4**, 587 (2018).
48. Paquette, M., *et al.* Polygenic risk score predicts prevalence of cardiovascular disease in patients with familial hypercholesterolemia. *J. Clin. Lipidol.* **11**, 725-732.e725 (2017).
49. Trinder, M., *et al.* Risk of Premature Atherosclerotic Disease in Patients With Monogenic Versus Polygenic Familial Hypercholesterolemia. *J. Am. Coll. Cardiol.* **74**, 512-522 (2019).
50. Fry, A., *et al.* Comparison of Sociodemographic and Health-Related Characteristics of UK Biobank Participants With Those of the General Population. *Am. J. Epidemiol.* **186**, 1026-1034 (2017).
51. Zuk, O., *et al.* Searching for missing heritability: designing rare variant association studies. *Proc. Natl. Acad. Sci. U. S. A.* **111**, E455-464 (2014).
52. Rivera-Muñoz, E.A., *et al.* ClinGen Variant Curation Expert Panel experiences and standardized processes for disease and gene-level specification of the ACMG/AMP guidelines for sequence variant interpretation. *Hum. Mutat.* **39**, 1614-1622 (2018).

53. Cholesterol Treatment Trialists, C., *et al.* Efficacy and safety of more intensive lowering of LDL cholesterol: a meta-analysis of data from 170,000 participants in 26 randomised trials. *Lancet* **376**, 1670-1681 (2010).
54. Zhao, Z., *et al.* Comparative efficacy and safety of lipid-lowering agents in patients with hypercholesterolemia: A frequentist network meta-analysis. *Medicine (Baltimore)* **98**, e14400 (2019).
55. Patel, A.P., *et al.* Association of Rare Pathogenic DNA Variants for Familial Hypercholesterolemia, Hereditary Breast and Ovarian Cancer Syndrome, and Lynch Syndrome With Disease Risk in Adults According to Family History. *JAMA Netw Open* **3**, e203959 (2020).
56. Sudlow, C., *et al.* UK biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLoS Med.* **12**, e1001779 (2015).
57. Eastwood, S.V., *et al.* Algorithms for the Capture and Adjudication of Prevalent and Incident Diabetes in UK Biobank. *PLoS One* **11**, e0162388 (2016).
58. Genomes Project, C., *et al.* A global reference for human genetic variation. *Nature* **526**, 68-74 (2015).
59. Bellenguez, C., *et al.* A robust clustering algorithm for identifying problematic samples in genome-wide association studies. *Bioinformatics* **28**, 134-135 (2012).
60. Sharifi, M., Futema, M., Nair, D. & Humphries, S.E. Genetic Architecture of Familial Hypercholesterolaemia. *Curr. Cardiol. Rep.* **19**, 44 (2017).
61. Peterson, A.S., Fong, L.G. & Young, S.G. PCSK9 function and physiology. *J. Lipid Res.* **49**, 1152-1156 (2008).
62. Whitfield, A.J., Barrett, P.H.R., van Bockxmeer, F.M. & Burnett, J.R. Lipid disorders and mutations in the APOB gene. *Clin. Chem.* **50**, 1725-1732 (2004).
63. McLaren, W., *et al.* The Ensembl Variant Effect Predictor. *Genome Biol.* **17**, 122 (2016).
64. Karczewski, K.J., *et al.* Variation across 141,456 human exomes and genomes reveals the spectrum of loss-of-function intolerance across human protein-coding genes. *bioRxiv* (2019).
65. Cummings, B.B., *et al.* Transcript expression-aware annotation improves rare variant discovery and interpretation. *bioRxiv* (2019).
66. Abou Tayoun, A.N., *et al.* Recommendations for interpreting the loss of function PVS1 ACMG/AMP variant criterion. *Hum Mutat* **39**, 1517-1524 (2018).
67. National Cholesterol Education Program Expert Panel on Detection, E. & Treatment of High Blood Cholesterol in, A. Third Report of the National Cholesterol Education Program (NCEP) Expert Panel on Detection, Evaluation, and Treatment of High Blood Cholesterol in Adults (Adult Treatment Panel III) final report. *Circulation* **106**, 3143-3421 (2002).
68. Surakka, I., *et al.* The impact of low-frequency and rare variants on lipid levels. *Nat. Genet.* **47**, 589-597 (2015).
69. Vilhjálmsson, B.J., *et al.* Modeling Linkage Disequilibrium Increases Accuracy of Polygenic Risk Scores. *Am. J. Hum. Genet.* **97**, 576-592 (2015).
70. Bulik-Sullivan, B.K., *et al.* LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.* **47**, 291-295 (2015).
71. Zheng, J., *et al.* LD Hub: a centralized database and web interface to perform LD score regression that maximizes the potential of summary level GWAS data for SNP heritability and genetic correlation analysis. *Bioinformatics* **33**, 272-279 (2017).
72. Purcell, S., *et al.* PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* **81**, 559-575 (2007).

Table 1: Impact of clinically significant variants on traits

Condition (proxy measure)	Gene	AMP-T2D-GENES (N=38,618)			UK Biobank (N=38,566)		
		N carrier	Beta (se)	P value*	N carrier	Beta (se)	P value*
High LDL (LDL mg/dL)	composite	55	56.0 (5.2)	3.9×10^{-24}	83	54.2 (3.9)	1.6×10^{-44}
	<i>APOB</i>	11	31.5 (12.0)	8.9×10^{-3}	26	52.2 (6.8)	2.2×10^{-14}
	<i>LDLR</i>	44	65.3 (6.3)	9.0×10^{-25}	57	55.1 (4.7)	1.1×10^{-31}
Low LDL (LDL mg/dL)	composite	35	-56.1 (7.1)	4.4×10^{-15}	90	-56.4 (3.7)	6.9×10^{-52}
	<i>APOB</i>	8	-79.8 (14.7)	5.9×10^{-8}	48	-74.5 (5.1)	6.7×10^{-48}
	<i>PCSK9</i>	27	-48.7 (8.2)	2.6×10^{-9}	42	-36.1 (5.4)	2.7×10^{-11}
High HDL (HDL mg/dL)	<i>CETP</i>	21	16.5 (3.0)	3.6×10^{-8}	20	16.8 (2.4)	2.3×10^{-12}
High triglycerides (TG mg/dL)	composite	20	130.0 (27.3)	2.8×10^{-6}	54	126.0 (12.2)	2.4×10^{-16}
	<i>APOA5</i>	15	122.4 (29.7)	2.6×10^{-5}	38	145.5 (13.6)	2.4×10^{-14}
	<i>LPL</i>	5	152.8 (54.6)	2.5×10^{-2}	16	79.3 (22.4)	9.4×10^{-4}
Monogenic obesity (BMI kg/m ²)	<i>MC4R</i>	28	1.5 (1.0)	6.3×10^{-2}	31	2.2 (0.8)	6.3×10^{-3}
Condition	gene	N carrier	OR	P value*	OR	P value*	
MODY (diabetes)	composite	22	7.8 (4.2-14.6)	6.5×10^{-5}	16	21 (12.5-35.2)	2.6×10^{-8}
	<i>GCK</i>	7	37.4 (6.3- 222.0)	1.3×10^{-3}	10	40.5 (20.3- 80.7)	3.1×10^{-8}
	<i>HNF1A</i>	11	4.8 (2.2-10.4)	1.7×10^{-2}	5	9.0 (3.51-22.9)	2.3×10^{-2}
MODY (T2D and prediabetes)	composite	22	4.8 (2.6-8.8)	2.5×10^{-3}	16	21.5 (11.5- 40.4)	9.1×10^{-9}
	<i>GCK</i>	7	17.8 (3.4-94.0)	8.2×10^{-3}	10	132.0 (28.7- 611.0)	1.4×10^{-9}
	<i>HNF1A</i>	11	3.1 (1.5-6.6)	8.9×10^{-2}	5	5.1 (2.0-12.9)	6.1×10^{-2}

Composite = individuals carrying variants in any of the genes analyzed for that condition. Note that MODY composite gene set included *GCK*, *HNF1A*, *HNF1B*, *HNF4A*, and *PDX1*.

*Comparison of variant carriers to non-carriers, adjusted for age, sex, 10 PCs

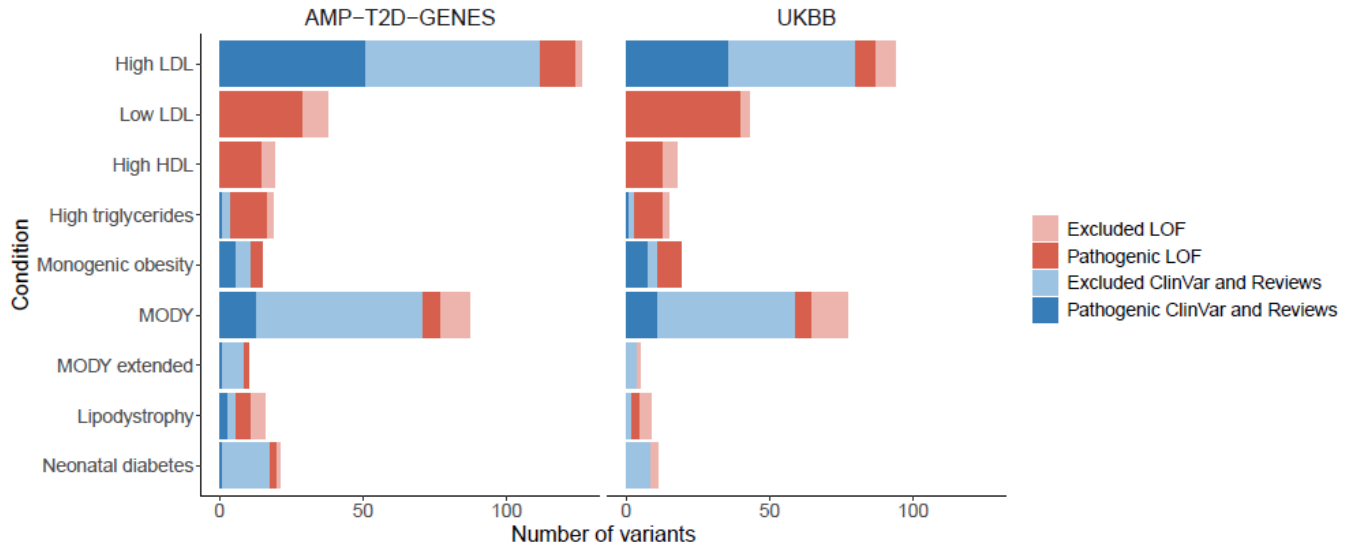


Figure 1. Curation of ClinVar and pLoF variants across the monogenic conditions.

Total number of curated ClinVar/Review (blue) and pLoF (red) variants with carriers in AMP-T2D-GENES (left panel) and UKB (right panel). Darker color shades indicate variants determined to be clinically significant (pathogenic, likely pathogenic, or pLoF) and lighter shades indicate variants excluded during curation from further analysis.

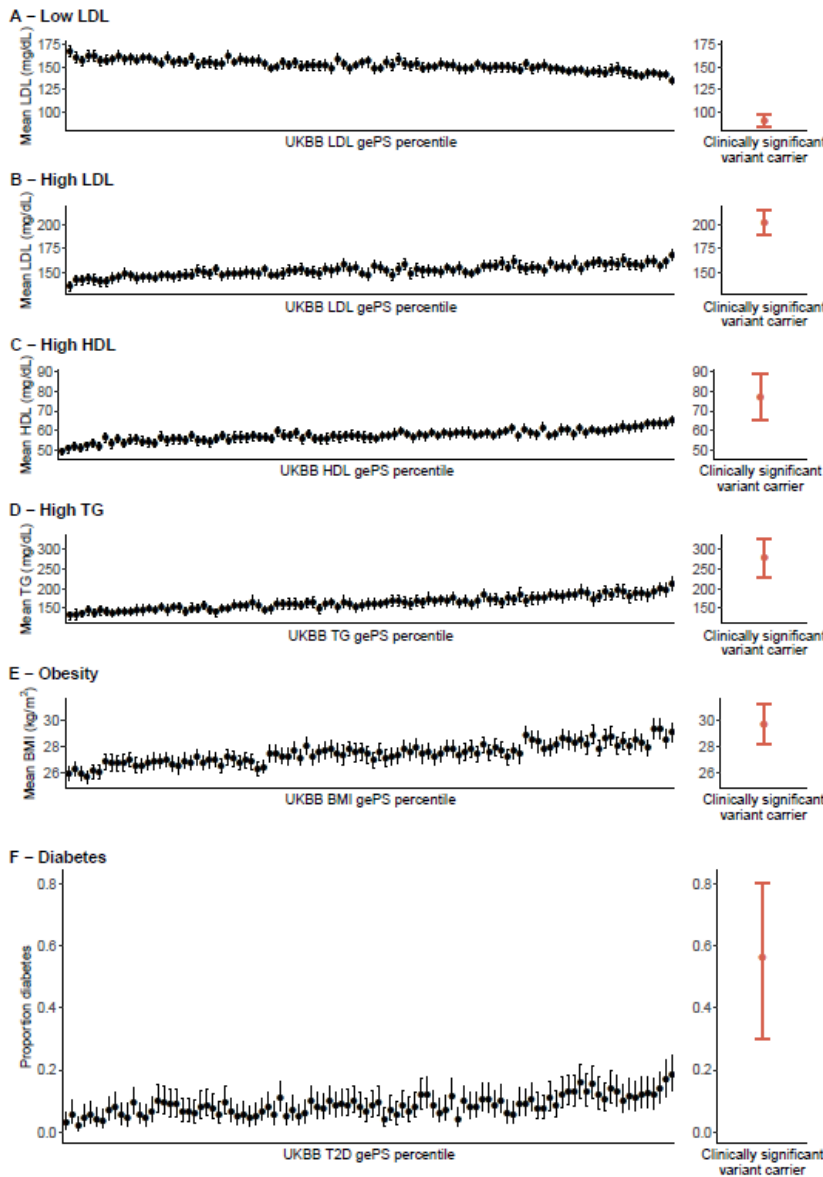


Figure 2. Carriers of rare clinically significant monogenic variants for lipid conditions and monogenic diabetes have more extreme effect size estimates than individuals with the top 1% of global extended polygenic scores (gePS).

In all plots, the left panels show the distribution of the phenotype in each percentile of the gePS for the relevant condition (black), and the right panel shows the phenotype distribution in carriers of rare clinically significant monogenic variants for the corresponding condition (red; low LDL cholesterol (*APOB*, *PCSK9*), high LDL cholesterol (*LDLR*, *APOB*), high HDL cholesterol (*CETP*), high triglycerides (*APOA5*, *LPL*), monogenic obesity (*MC4R*), and MODY (*GCK*, *HNF1A*, *PDX1*). **A-E**) Mean and 95% CI of each phenotype are indicated by the point and error bars respectively. The same gePS calculated for risk of increasing LDL levels was used for **A** and **B**; however, the inverse of this gePS was used for **B** to illustrate that higher gePS indicates risk of lower LDL cholesterol. **F**) The proportion of individuals with diabetes and 95% CI computed with the Clopper-Pearson method are shown as points and error bars respectively. Individuals in the gePS analysis were restricted to those age ≥ 60 years. LDL cholesterol and triglyceride values were adjusted for lipid-lowering medication use as per methods.

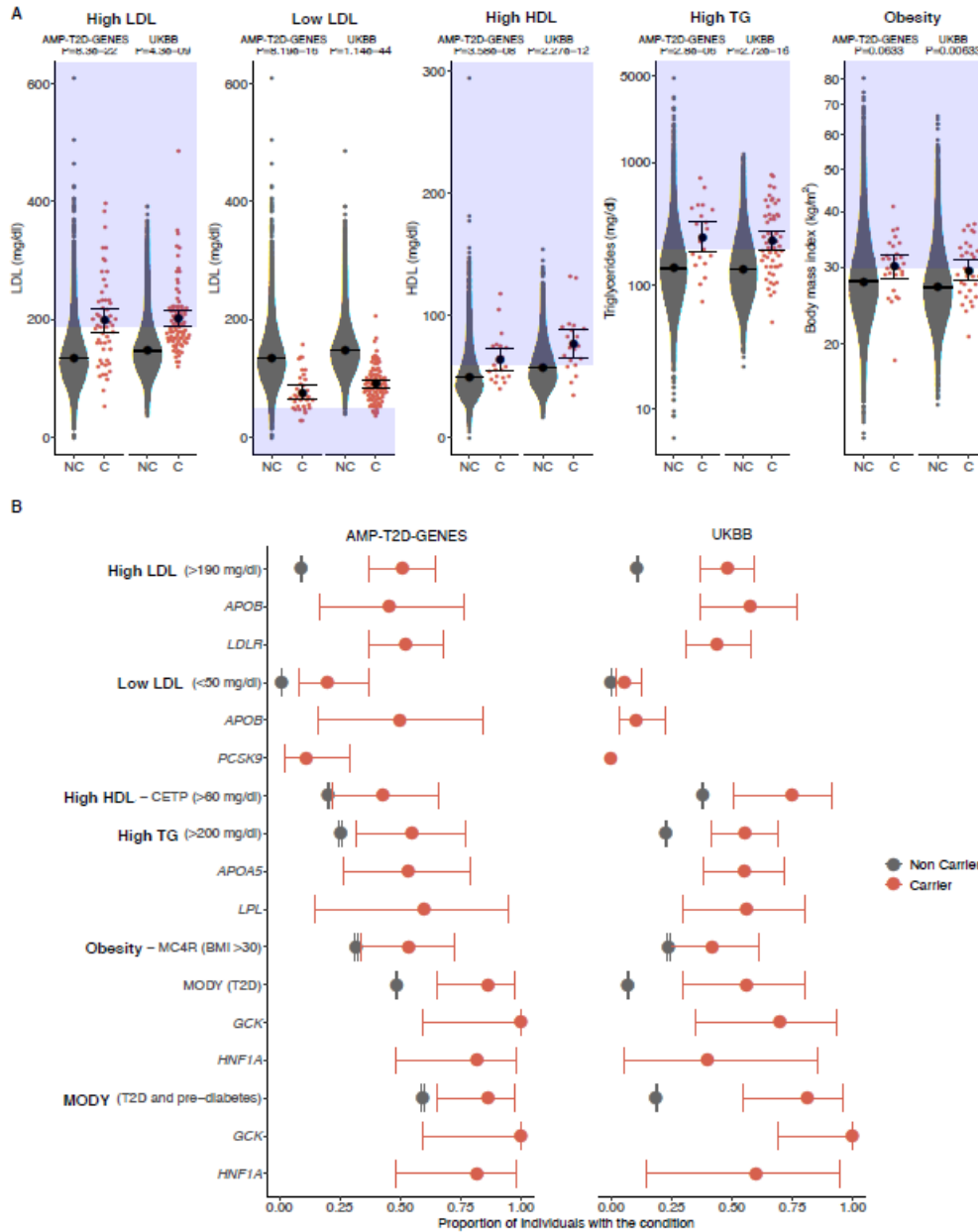


Figure 3. Phenotype distributions and penetrance estimates of clinically significant variant carriers.

In all plots, clinically significant variant carriers are shown in red and non-carriers are shown in grey. The left panel of each plot shows AMP-T2D-GENES participants (T2D case/control study) and the right panel shows UK Biobank participants (population-based study). **A**) Mean and 95% CI are represented by the black circle and black lines respectively. Relevant lipid levels (mg/dl) or body mass index (kg/m²) are shown for carriers (C) and non-carriers (NC) of clinically significant variants for the five monogenic conditions. The blue boxes indicate the phenotype values that meet a clinical threshold for diagnosis of each of the conditions, and *P* values were obtained by burden analysis in EPACTS (see **Methods**). **B**) Dots are the proportion of individuals that have the condition based on the clinical diagnosis threshold for each condition; for MODY, we show the proportion of individuals meeting T2D as well as T2D and prediabetes criteria (see **Methods**). Error bars reflect 95% CI computed with the Clopper-Pearson method.

It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

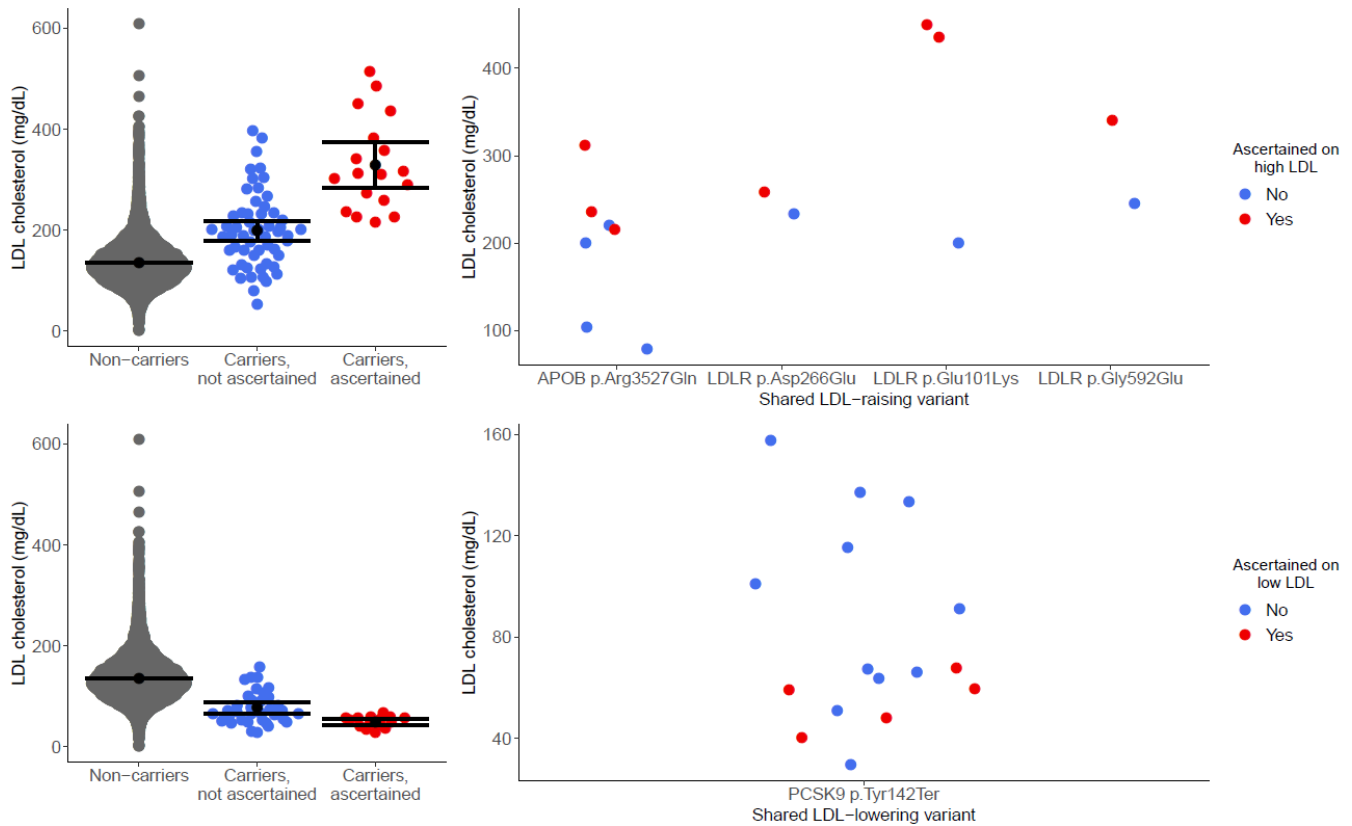


Figure 4. Ascertainment bias significantly impacts expressivity of clinically significant variants for LDL cholesterol conditions.

LDL cholesterol levels are shown for carriers and non-carriers of LDL cholesterol raising (top panels) or lowering (bottom panels) clinically significant variants. The variants carriers are stratified by whether they were identified in individuals phenotypically ascertained for extreme serum LDL cholesterol levels (Yes, Red) or in a separate unascertained population (No, Blue) (see **Methods**). The left panels show all clinically significant variant carriers. The right panels show carriers of the single variants that were present in both ascertained and unascertained individuals. LDL cholesterol values are adjusted for lipid-lowering medication use as per methods.

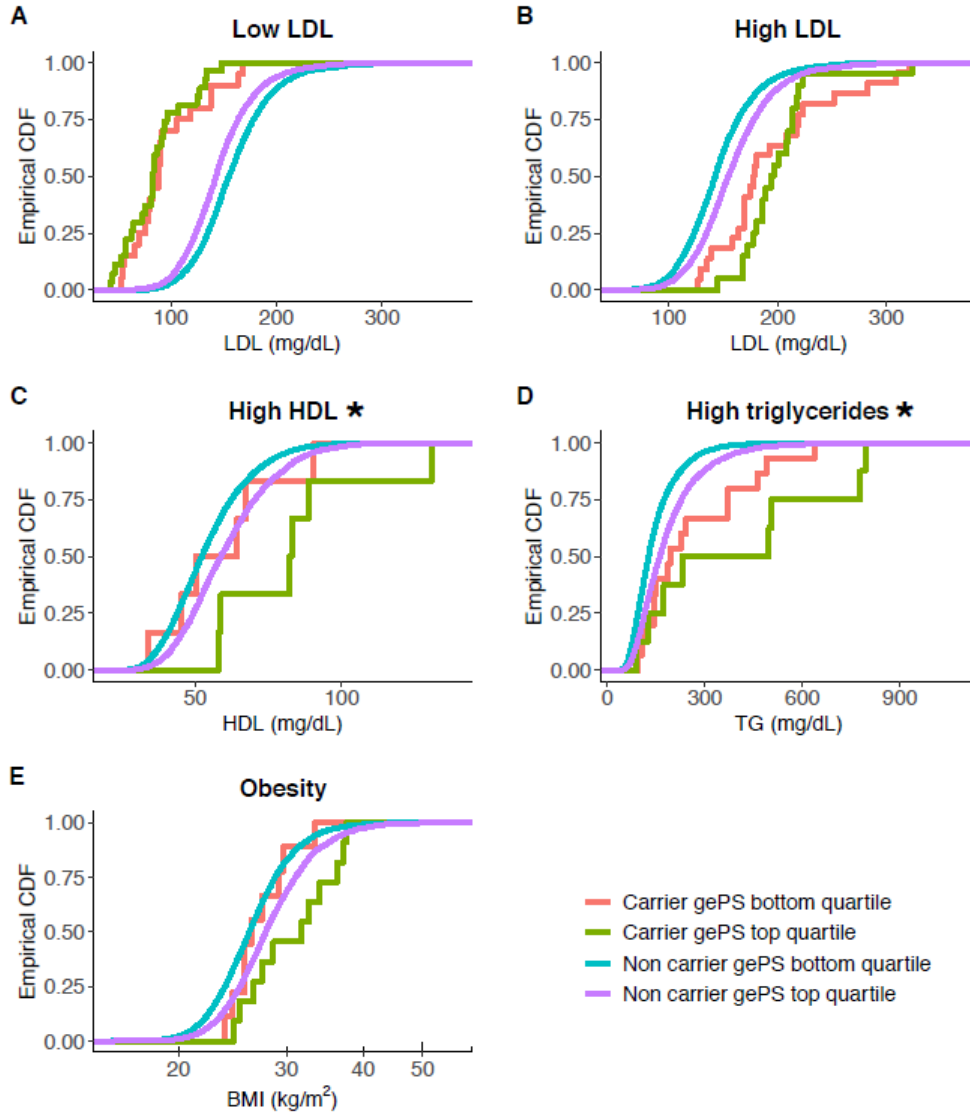


Figure 5. The combination of clinically significant monogenic variants and corresponding polygenic scores significantly improves prediction for high HDL cholesterol and high triglyceride conditions.

In all plots, an empirical cumulative distribution function (CDF) of each phenotype is shown for clinically significant variant carriers and non-carriers in the UKB for each monogenic condition stratified by bottom/top quartiles of the corresponding gePS. The monogenic conditions are low LDL cholesterol (*APOB*, *PCSK9*), high LDL cholesterol (*LDLR*, *APOB*), high HDL cholesterol (*CETP*), high triglycerides (*APOA5*, *LPL*), and monogenic obesity (*MC4R*). The same gePS calculated for risk of increasing LDL cholesterol levels was used for **A** and **B**, however the inverse of the gePS was used for **A** to illustrate that higher gePS indicates risk of lower LDL cholesterol. Asterisks indicate $P < 0.05$.