

Fine-mapping type 2 diabetes loci to single-variant resolution using high-density imputation and islet-specific epigenome maps

Anubha Mahajan^{1,2}, Daniel Taliun³, Matthias Thurner^{1,2}, Neil R Robertson^{1,2}, Jason M Torres¹, N William Rayner^{1,2,4}, Anthony J Payne¹, Valgerdur Steinthorsdottir⁵, Robert A Scott⁶, Niels Grarup⁷, James P Cook⁸, Ellen M Schmidt³, Matthias Wuttke⁹, Chloé Sarnowski¹⁰, Reedik Mägi¹¹, Jana Nano¹², Christian Gieger^{13,14}, Stella Trompet^{15,16}, Cécile Lecoeur¹⁷, Michael H Preuss¹⁸, Bram Peter Prins⁴, Xiuqing Guo¹⁹, Lawrence F Bielak²⁰, Jennifer E Below²¹, Donald W Bowden^{22,23,24}, John Campbell Chambers^{25,26,27,28,29}, Young Jin Kim³⁰, Maggie CY Ng^{22,23,24}, Lauren E Petty²¹, Xueling Sim³¹, Weihua Zhang^{25,26}, Amanda J Bennett², Jette Bork-Jensen⁷, Chad M Brummett³², Mickaël Canouil¹⁷, Kai-Uwe Eckardt³³, Krista Fischer¹¹, Sharon LR Kardia²⁰, Florian Kronenberg³⁴, Kristi Läll^{11,35}, Ching-Ti Liu¹⁰, Adam E Locke^{36,37}, Jian'an Luan⁶, Ioanna Ntalla³⁸, Vibe Nylander², Sebastian Schönherr³⁴, Claudia Schurmann¹⁸, Loïc Yengo¹⁷, Erwin P Bottinger¹⁸, Ivan Brandslund^{39,40}, Cramer Christensen⁴¹, George Dedoussis⁴², Jose C Florez^{43,44,45,46}, Ian Ford⁴⁷, Oscar H Franco¹², Timothy M Frayling⁴⁸, Vilmantas Giedraitis⁴⁹, Sophie Hacking⁴, Andrew T Hattersley⁵⁰, Christian Herder^{14,51}, M Arfan Ikram¹², Martin Ingelsson⁴⁹, Marit E Jørgensen^{52,53}, Torben Jørgensen^{54,55,56}, Jennifer Kriebel^{13,14}, Johanna Kuusisto⁵⁷, Symen Ligthart¹², Cecilia M Lindgren^{1,58,59}, Allan Linneberg^{54,60,61}, Valeriya Lyssenko^{62,63}, Vasiliki Mamakou⁶⁴, Thomas Meitinger^{65,66,67}, Karen L Mohlke⁶⁸, Andrew D Morris^{69,70}, Girish Nadkarni⁷¹, James S Pankow⁷², Annette Peters^{14,67,73}, Naveed Sattar⁷⁴, Alena Stančáková⁵⁷, Konstantin Strauch^{75,76}, Kent D Taylor¹⁹, Barbara Thorand^{14,73}, Gudmar Thorleifsson⁵, Unnur Thorsteinsdottir^{5,77}, Jaakko Tuomilehto^{78,79,80,81}, Daniel R Witte^{82,83}, Josée Dupuis^{10,84}, Patricia A Peyser²⁰, Eleftheria Zeggini⁴, Ruth J F Loos^{18,85}, Philippe Froguel^{17,86}, Erik Ingelsson^{87,88}, Lars Lind⁸⁹, Leif Groop^{62,90}, Markku Laakso⁵⁷, Francis S Collins⁹¹, J Wouter Jukema¹⁶, Colin N A Palmer⁹², Harald Grallert^{13,14,93,94}, Andres Metspalu¹¹, Abbas Dehghan^{12,25,29}, Anna Köttgen⁹, Goncalo R Abecasis³, James B Meigs^{43,46,95}, Jerome I Rotter^{19,96}, Jonathan Marchini^{1,97}, Oluf Pedersen⁷, Torben Hansen^{7,98}, Claudia Langenberg⁶,

Nicholas J Wareham⁶, Kari Stefansson^{5,77}, Anna L Gloyn^{1,2,99}, Andrew P Morris^{1,8,11,*}, Michael Boehnke^{3,*}, Mark I McCarthy^{1,2,99,*}

1. Wellcome Centre for Human Genetics, Nuffield Department of Medicine, University of Oxford, Oxford, OX3 7BN, UK.
2. Oxford Centre for Diabetes, Endocrinology and Metabolism, Radcliffe Department of Medicine, University of Oxford, Oxford, OX3 7LE, UK.
3. Department of Biostatistics and Center for Statistical Genetics, University of Michigan, Ann Arbor, Michigan, 48109, USA.
4. Department of Human Genetics, Wellcome Trust Sanger Institute, Hinxton, Cambridgeshire, CB10 1SA, UK.
5. deCODE Genetics, Amgen inc., Reykjavik, 101, Iceland.
6. MRC Epidemiology Unit, Institute of Metabolic Science, University of Cambridge, Cambridge, CB2 0QQ, UK.
7. The Novo Nordisk Foundation Center for Basic Metabolic Research, Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen, 2200, Denmark.
8. Department of Biostatistics, University of Liverpool, Liverpool, L69 3GL, UK.
9. Institute of Genetic Epidemiology, Department of Biometry, Epidemiology, and Medical Bioinformatics, Faculty of Medicine and Medical Center, University of Freiburg, Freiburg, 79106, Germany.
10. Department of Biostatistics, Boston University School of Public Health, Boston, Massachusetts, 02118, USA.
11. Estonian Genome Center, Institute of Genomics, University of Tartu, Tartu, 51010, Estonia.
12. Department of Epidemiology, Erasmus University Medical Center, Rotterdam, 3015CN, The Netherlands.
13. Research Unit of Molecular Epidemiology, Institute of Epidemiology, Helmholtz Zentrum München Research Center for Environmental Health, Neuherberg, 85764, Germany.
14. German Center for Diabetes Research (DZD), Neuherberg, 85764, Germany.
15. Section of Gerontology and Geriatrics, Department of Internal Medicine, Leiden University Medical Center, Leiden, 2300 RC, the Netherlands.

16. Department of Cardiology, Leiden University Medical Center, Leiden, 2300 RC, the Netherlands.
17. CNRS-UMR8199, Lille University, Lille Pasteur Institute, Lille, 59000, France.
18. The Charles Bronfman Institute for Personalized Medicine, The Icahn School of Medicine at Mount Sinai, New York, 10029, USA.
19. Department of Pediatrics, The Institute for Translational Genomics and Population Sciences, LABioMed at Harbor-UCLA Medical Center, Torrance, California, 90502, US.
20. Department of Epidemiology, School of Public Health, University of Michigan, Ann Arbor, Michigan, 48109, USA.
21. The Vanderbilt Genetics Institute, Vanderbilt University Medical Center, Nashville, TN 37232, USA.
22. Center for Genomics and Personalized Medicine Research, Wake Forest School of Medicine, Winston-Salem, North Carolina, 27157, USA.
23. Center for Diabetes Research, Wake Forest School of Medicine, Winston-Salem, North Carolina, 27157, USA.
24. Department of Biochemistry, Wake Forest School of Medicine, Winston-Salem, North Carolina, 27157, USA.
25. Department of Epidemiology and Biostatistics, Imperial College London, London, W2 1PG, UK.
26. Department of Cardiology, Ealing Hospital, London North West Healthcare NHS Trust, Middlesex, UB1 3HW, UK.
27. Imperial College Healthcare NHS Trust, Imperial College London, London, W12 0HS, UK.
28. Lee Kong Chian School of Medicine, Nanyang Technological University, Singapore, 308232, Singapore.
29. MRC-PHE Centre for Environment and Health, Imperial College London, London, W2 1PG, UK.
30. Division of Genome Research, Center for Genome Science, Korea National Institute of Health, Chungcheongbuk-do, Republic of Korea.
31. Saw Swee Hock School of Public Health, National University of Singapore, Singapore, 117549, Singapore.

32. Department of Anesthesiology, University of Michigan Medical School, Ann Arbor, Michigan, 48109, US.
33. Department of Nephrology and Medical Intensive Care, Charité, Universitätsmedizin Berlin, Berlin, 10117, Germany and German Chronic Kidney Disease study.
34. Division of Genetic Epidemiology, Department of Medical Genetics, Molecular and Clinical Pharmacology, Medical University of Innsbruck, Innsbruck, 6020, Austria.
35. Institute of Mathematics and Statistics, University of Tartu, J. Liivi 2, Tartu, 50409, Estonia.
36. McDonnell Genome Institute, Washington University School of Medicine, St. Louis, MO, USA.
37. Department of Medicine, Division of Genomics & Bioinformatics, Washington University School of Medicine, St. Louis, MO, USA.
38. William Harvey Research Institute, Barts and The London School of Medicine and Dentistry, Queen Mary University of London, London, UK.
39. Institute of Regional Health Research, University of Southern Denmark, Odense, 5000, Denmark.
40. Department of Clinical Biochemistry, Vejle Hospital, Vejle, 7100, Denmark.
41. Medical Department, Lillebælt Hospital Vejle, Vejle, Denmark.
42. Department of Nutrition and Dietetics, Harokopio University of Athens, Athens, 17671, Greece.
43. Department of Medicine, Harvard Medical School, Boston, Massachusetts, 02115, USA.
44. Diabetes Research Center (Diabetes Unit), Department of Medicine, Massachusetts General Hospital, Boston, Massachusetts, 02114, USA.
45. Center for Genomic Medicine, Massachusetts General Hospital, Boston, MA, 02114, USA.
46. Programs in Metabolism and Medical & Population Genetics, Broad Institute, Cambridge, MA, 02142, USA.
47. Robertson Centre for Biostatistics, University of Glasgow, Glasgow, UK.
48. Genetics of Complex Traits, University of Exeter Medical School, University of Exeter, Exeter, EX1 2LU, UK.

49. Department of Public Health and Caring Sciences, Geriatrics, Uppsala University, Uppsala, SE-751 85, Sweden.
50. University of Exeter Medical School, University of Exeter, Exeter, EX2 5DW, UK.
51. Institute for Clinical Diabetology, German Diabetes Center, Leibniz Center for Diabetes Research at Heinrich Heine University Düsseldorf, Düsseldorf, Germany.
52. Steno Diabetes Center Copenhagen, Gentofte, 2820, Denmark.
53. National Institute of Public Health, Southern Denmark University, Copenhagen, 1353, Denmark.
54. Research Centre for Prevention and Health, Capital Region of Denmark, Glostrup, 2600, Denmark.
55. Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen, Denmark.
56. Faculty of Medicine, Aalborg University, Aalborg, Denmark.
57. Institute of Clinical Medicine, Internal Medicine, University of Eastern Finland and Kuopio University Hospital, Kuopio, 70210, Finland.
58. Program in Medical and Population Genetics, Broad Institute, Cambridge, Massachusetts, 02142, USA.
59. Big Data Institute, Li Ka Shing Centre For Health Information and Discovery, University of Oxford, Oxford, OX37BN, UK.
60. Center for Clinical Research and Prevention, Bispebjerg and Frederiksberg Hospital, Frederiksberg, Denmark.
61. Department of Clinical Medicine, Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen, Denmark.
62. Department of Clinical Sciences, Diabetes and Endocrinology, Lund University Diabetes Centre, Malmö, 20502, Sweden.
63. Department of Clinical Science, KG Jebsen Center for Diabetes Research, University of Bergen, Bergen, Norway.
64. Dromokaiteio Psychiatric Hospital, National and Kapodistrian University of Athens, Athens, Greece.
65. Institute of Human Genetics, Technische Universität München, Munich, 81675, Germany.

66. Institute of Human Genetics, Helmholtz Zentrum München, German Research Center for Environmental Health, Neuherberg, 85764, Germany.
67. DZHK (German Centre for Cardiovascular Research), partner site Munich Heart Alliance, Munich, 81675, Germany.
68. Department of Genetics, University of North Carolina, Chapel Hill, North Carolina, 27599, USA.
69. Clinical Research Centre, Centre for Molecular Medicine, Ninewells Hospital and Medical School, Dundee, DD1 9SY, UK.
70. The Usher Institute to the Population Health Sciences and Informatics, University of Edinburgh, Edinburgh, EH16 4UX, UK.
71. Division of Nephrology, Department of Medicine, Icahn School of Medicine at Mount Sinai, New York, NY, 10069, USA.
72. Division of Epidemiology and Community Health, School of Public Health, University of Minnesota, Minneapolis, MN, 55454, US.
73. Institute of Epidemiology, Helmholtz Zentrum München, German Research Center for Environmental Health, Neuherberg, 85764, Germany.
74. Institute of Cardiovascular and Medical Sciences, University of Glasgow, Glasgow, G12 8TA, UK.
75. Institute of Genetic Epidemiology, Helmholtz Zentrum München, German Research Center for Environmental Health, Neuherberg, 85764, Germany.
76. Chair of Genetic Epidemiology, IBE, Faculty of Medicine, LMU, Munich, 81377, Germany.
77. Faculty of Medicine, University of Iceland, Reykjavik, 101, Iceland.
78. Department of Health, National Institute for Health and Welfare, Helsinki, 00271, Finland.
79. Dasman Diabetes Institute, Dasman, 15462, Kuwait.
80. Department of Neuroscience and Preventive Medicine, Danube-University Krems, Krems, 3500, Austria.
81. Diabetes Research Group, King Abdulaziz University, Jeddah, 21589, Saudi Arabia.
82. Department of Public Health, Aarhus University, Aarhus, Denmark.
83. Danish Diabetes Academy, Odense, Denmark.

84. National Heart, Lung, and Blood Institute's Framingham Heart Study, Framingham, Massachusetts, 01702, USA.
85. Mindich Child Health and Development Institute, The Icahn School of Medicine at Mount Sinai, New York, NY, 10029, USA.
86. Department of Genomics of Common Disease, School of Public Health, Imperial College London, London, W12 0NN, UK.
87. Department of Medicine, Division of Cardiovascular Medicine, Stanford University School of Medicine, Stanford, CA, 94305, US.
88. Department of Medical Sciences, Molecular Epidemiology and Science for Life Laboratory, Uppsala University, Uppsala, 75185, Sweden.
89. Department of Medical Sciences, Uppsala University, Uppsala, SE-751 85, Sweden.
90. Finnish Institute for Molecular Medicine (FIMM), University of Helsinki, Helsinki, Finland.
91. Genome Technology Branch, National Human Genome Research Institute, National Institutes of Health, Bethesda, Maryland, 20892, USA.
92. Pat Macpherson Centre for Pharmacogenetics and Pharmacogenomics, Ninewells Hospital and Medical School, University of Dundee, Dundee, DD1 9SY, UK.
93. Clinical Cooperation Group Type 2 Diabetes, Helmholtz Zentrum München, Ludwig-Maximilians University Munich, Germany.
94. Clinical Cooperation Group Nutrigenomics and Type 2 Diabetes, Helmholtz Zentrum München, Technical University Munich, Germany.
95. Division of General Internal Medicine, Massachusetts General Hospital, Boston, Massachusetts, 02114, USA.
96. Departments of Medicine, The Institute for Translational Genomics and Population Sciences, LABioMed at Harbor-UCLA Medical Center, Torrance, California, 90502, US.
97. Department of Statistics, University of Oxford, Oxford, OX1 3TG, UK.
98. Faculty of Health Sciences, University of Southern Denmark, Odense, 5000, Denmark.
99. Oxford NIHR Biomedical Research Centre, Oxford University Hospitals Trust, Oxford, OX3 7LE, UK.

*These authors contributed jointly to this work.

Correspondence to:

Anubha Mahajan (anubha@well.ox.ac.uk)

Michael Boehnke (boehnke@umich.edu)

Mark I McCarthy (mark.mccarthy@drl.ox.ac.uk)

We expanded GWAS discovery for type 2 diabetes (T2D) by combining data from 898,130 European-descent individuals (9% cases), following imputation to high-density reference panels. With these data, we: (a) extend the inventory of T2D-risk variants (243 loci, 135 newly-implicated in T2D-predisposition, comprising 403 distinct association signals); (b) enrich discovery of lower-frequency risk-alleles (80 index variants with minor allele frequency <5%, 14 with estimated allelic odds-ratio >2); (c) substantially improve fine-mapping of causal variants (at 51 signals, one variant accounted for >80% posterior probability of association (PPA)); (d) extend fine-mapping through integration of tissue-specific epigenomic information (islet regulatory annotations extend to 73 the number of variants with PPA>80%); (e) highlight validated therapeutic targets (18 genes with associations attributable to coding variants); and (f) demonstrate enhanced potential for clinical translation (genome-wide chip heritability explains 18% of T2D-risk; individuals in the extremes of a T2D polygenic risk score differ >9-fold in prevalence).

Array-based genome-wide association studies (GWAS) have identified ~140 loci influencing risk of type 2 diabetes (T2D)¹⁻³. Follow-up of these genetic discoveries has been compromised by the incomplete coverage of the most frequently-used genotyping arrays, imperfect performance of the reference panels available for imputation, extensive local linkage disequilibrium (LD), and inadequate sample sizes. These have combined to reduce power to detect low-frequency alleles with population-scale impact, to deliver clinically-relevant risk prediction, and to define molecular mechanisms involved in disease predisposition. Here, we address the limitations of previous studies by combining GWAS from ~900,000 Europeans with dense, high-quality imputation, producing the most comprehensive view to date of the genetic contribution to T2D with respect to locus discovery, causal variant resolution, and mechanistic insight.

RESULTS

Study overview. We combined data from 32 GWAS, including 74,124 T2D cases and 824,006 controls of European ancestry (effective sample size [N_{eff}] 231,436). This represents a 3.2-fold increase in N_{eff} from the largest previous genome-wide study of T2D-risk in Europeans¹. After harmonised quality control, 31 of the 32 GWAS were imputed using 64,976 whole-

genome sequenced haplotypes from the Haplotype Reference Consortium (HRC)⁴: the exception was the deCODE GWAS, imputed using a population-specific reference panel of 30,440 Icelandic haplotypes⁵ (**Methods, Supplementary Table 1**). We conducted T2D association analyses with and without adjustment for body-mass index (BMI).

Discovery of novel loci for T2D susceptibility. We tested for T2D association with ~27M variants passing quality control filters, ~21M having minor allele frequency (MAF)<5%. Our meta-analysis identified variants at 231 loci reaching genome-wide significance ($p < 5 \times 10^{-8}$) in the BMI-unadjusted analysis (N_{eff} 231,436) and 152 in the smaller (N_{eff} 157,401) BMI-adjusted analysis. Of the 243 loci identified across these two analyses, 135 mapped outside regions previously-implicated in T2D-risk (**Methods, Figure 1, Supplementary Table 2**).

Amongst samples not included in previous discovery efforts (42,734 cases, 497,261 controls), we replicated associations (directionally-consistent, $p < 0.05$) at 126 of 140 previously-reported T2D loci, including all 106 regions first discovered in European-only or trans-ethnic efforts^{3,6-8} and 20 initially reported in studies of non-European individuals^{9,10}. The 14 loci not replicated were all first-identified in non-European ancestry samples: at five, the reported lead variant had MAF<1% in Europeans.

Multiple association signals at T2D susceptibility loci. Across the 243 associated loci, we identified 160 additional signals at “locus-wide” significance ($p < 10^{-5}$; **Methods**), 110 within previously-reported T2D loci. Overall, we observed one signal at 151 loci, and two to ten signals at the remaining 92 (**Supplementary Table 2**), for a total of 403 distinct T2D-association signals.

We observed the first evidence for multiple signals at the *TCF7L2* locus. In addition to rs7903146, the largest-effect common variant signal for T2D in Europeans, we detected seven secondary signals, each represented by non-coding index variants ($0.5\% < \text{MAF} < 47.6\%$, $1.05 < \text{odds-ratio [OR]} < 1.36$).

In the ~1Mb telomeric region of chromosome 11 that encompasses the (previously-annotated) *INS-IGF2* and *KCNQ1* loci, we detected 15 distinct signals ($0.15\% < \text{MAF} < 42.8\%$,

1.03<OR<1.68). This multiplicity of signals in a region notable for complex imprinting effects, and several strong biological candidates (*INS*, *IGF2*, *KCNQ1*, *CDKN1C*), illustrates a previously-unrecognised degree of complexity in the risk-variant architecture at this locus.

The effects of BMI and sex. At most T2D-loci, there were only minimal differences in estimated T2D effect size between BMI-adjusted and unadjusted models (**Methods, Figure 2**). However, at index SNPs for 41 signals (mapping to 21 known and 16 novel loci), we observed significant differences in effect sizes between BMI-adjusted and unadjusted analyses ($p_{diff}<0.00012$, corrected for 403 variants; **Methods, Supplementary Table 3, Figure 2**). This effect-size heterogeneity followed two distinct patterns. At 26 signals, including index variants for signals at the *FTO*, *MC4R*, *TMEM18*, *SEC16B*, and *GNPDA2* loci, BMI-adjustment produced marked attenuation of associations detected in unadjusted analysis. These signals display positive correlations between BMI and T2D effect sizes, and represent T2D-risk effects primarily driven by adiposity. The other 15 signals were more strongly-associated in the BMI-adjusted analysis, and reflect a mixture of associations, some with a marked effect on insulin secretion (e.g. *TCF7L2*, *ARAP1*, *JAZF1*), and others that likely influence T2D-risk through reduced capacity for fat storage in peripheral adipose tissue¹¹ (e.g. *GRB14*, *PPARG*, *HMGGA1*, *ZNF664*).

In comparative analysis of T2D effects in males (41,846 cases, 383,767 controls) and females (30,053 T2D cases, 434,336 controls; **Methods**)¹², only one of the 403 T2D signals showed significant ($p_{diff}<0.00012$) differences in effect size (rs2925979 near *CMIP*, female OR=1.09, male OR=1.03, $p_{diff}=8.3\times 10^{-6}$; **Supplementary Figure 1, Supplementary Table 4**). We observed nominally-significant differences at several other loci, including *KLF14* (rs1562396, female OR=1.09, male OR=1.04, $p_{diff}=0.00048$) at which there is additional corroboration for sex-specific effects,^{13,14} indicating that additional examples of sex-differentiated signals are likely to be found in larger samples.

Fine-mapping variants driving T2D association signals. Previous efforts at fine-mapping causal variants within T2D loci have been hampered by both biological (extensive LD) and technical (diverse genotyping scaffolds, incomplete reference panels) factors. We sought to establish the extent to which the combination of increased sample size, enlarged reference

panel, and harmonised variant quality control would enhance fine-mapping resolution. We were able to undertake fine-mapping for 380 of the 403 distinct T2D association signals, following conditional decomposition of loci with multiple signals (**Methods**). For each, we constructed credible sets that collectively account for $\geq 99\%$ of the posterior probability of association (PPA; **Methods**)¹⁵. These credible sets included a median of 42 variants (range 1-3997; **Supplementary Figure 2**), and spanned a median of 116kb (range 1bp-995kb). At 51 signals, involving 44 loci (18 novel), the most strongly-associated variant accounted for $>80\%$ PPA (**Figure 3, Supplementary Table 5**). At 18, the credible set included a single variant (PPA $>99\%$).

We explored fine-mapping resolution at 83 distinct signals where detection in both studies allowed us to compare 99% credible sets from the HRC-based analysis with those constructed in a subset of these T2D GWAS imputed using the 1000 Genomes Project multi-ethnic reference panel¹ (26,676 T2D cases; 132,532 controls of European ancestry, N_{eff} 72,143). Although the former includes 2.3-fold more variants genome-wide than the latter, the HRC-imputed analysis resulted in smaller credible sets. The median number of variants at these 83 signals decreased from 59 to 10 and interval length from 60.3kb to 19.2kb. At 79 of 83 signals, HRC-based credible sets were either smaller than those generated from 1000G or unchanged (**Figure 4, Supplementary Table 6**).

This improved resolution likely reflects the combination of: (i) increased N_{eff} ; (ii) improved imputation quality, especially for lower-frequency variants⁴; and (iii) more effective, harmonised, quality control across contributing studies (**Methods**). To estimate the contribution to fine-mapping resolution attributable to the increase in N_{eff} (the other factors are more difficult to tease apart), we constructed 99% credible sets based on downscaling the HRC-imputation to a subset of 19 studies (31,387 cases; 326,742 controls, N_{eff} 92,960) that contributed to both 1000G and HRC-based analyses. Amongst 41 single signal loci with $p < 1 \times 10^{-5}$ in this downscaled meta-analysis, estimates of credible set size (median 66) and interval (median 196kb) indicate that the improvements in causal variant resolution derive mostly from increased sample size.

The HRC panel provides excellent coverage of all but very rare SNVs. However, one HRC limitation is the absence of indels: these constitute 4% of total variants in the phase 3 1000G reference panel¹⁶. We considered the 245,207 indels from the European subset of the 1000G panel which map within 500kb of index variants at the 380 fine-mapped signals: these account for 2.8% of variants across the 380Mb of sequence. Only 1% of these are in even moderate LD ($r^2 > 0.5$) with index variants for each T2D-association signal, indicating that indel omission is likely to have limited impact on our estimates of credible set size.

The contribution of lower-frequency variants. The limited yield of low-frequency and rare variant signals in previous T2D GWAS placed an upper bound on their individual and collective contribution to disease-risk¹⁷. The present analysis, with larger sample size and improved imputation, provides greater power in this regard, identifying 56 low-frequency and 24 rare T2D-associated variants across 60 loci (**Figure 5**). Six of these 80 signals mapped within known T2D loci, five reconfirming earlier observations (**Supplementary Table 2**).

Allelic OR for low-frequency and rare variants ranged from 1.08-8.05 (including 14 with estimated allelic OR > 2; at each, the minor allele conferred T2D-risk), compared with 1.03-1.37 for common variants (**Figure 5**). The 80 lower-frequency risk-variants cumulatively explained 1.1% of phenotypic variance in T2D, compared to 16.3% attributable to the 323 common-variant signals (**Methods**). Extrapolation beyond these discovered signals to estimate the full contribution of lower-frequency variants to T2D-risk is intrinsically difficult given the combination of effect-size overestimation and limited power to capture lower-frequency variants of lesser effect. Nonetheless, these data are consistent with recently-proposed models for the genetic architecture of T2D based on GWAS and sequencing data¹⁷. Notwithstanding, the identification of lower-frequency variants of modest to large effect can provide valuable biological inference, and we briefly describe some of these signals.

We observed a mix of common and low-frequency variant signals around *NEUROG3*, including T2D-risk attributable to the minor alleles at rs41277236 (p.Gly167Arg, MAF=4.3%, OR=1.09, $p=1.5 \times 10^{-6}$) and rs549498088 (non-coding, MAF=0.60%, OR=1.56, $p=4.7 \times 10^{-7}$). *NEUROG3* encodes the neurogenin-3 transcription factor implicated in pancreatic islet and

enteroendocrine cell development¹⁸. Rare homozygous, hypomorphic missense mutations in *NEUROG3* (*non-overlapping with those we detected*) are a cause of childhood-onset diabetes associated with severe malabsorptive diarrhea¹⁹. Age of T2D-diagnosis amongst carriers of these low-frequency T2D-risk alleles was, in UK Biobank, similar to non-carriers (rs41277236: 52.3 vs 52.7 years, $p=0.21$; rs549498088: 51.1 vs 52.7 years; $p=0.49$), consistent with a spectrum of functional impact that associates these variants with typical T2D. In UK Biobank, T2D-risk alleles at *NEUROG3* were associated with phenotypes recapitulating the gastrointestinal component of the neonatal syndrome (including “obstruction of bile duct” [OR=1.29; $p=0.023$], “gastrointestinal complications” [OR=1.79; $p=0.024$], and “functional digestive disorders” [OR=1.06; $p=0.027$]).

We detected two previously-unreported rare alleles with large ORs. The first was intronic to *DENND2C* (rs184660829, MAF=0.020%, OR=8.1, $p=2.5\times 10^{-8}$). In exploratory analyses within UK Biobank, the T2D-risk allele was associated with “lower gastrointestinal congenital anomalies” (OR=17.3 $p=0.00047$). The second mapped near *KIF2B* (rs569511541, MAF=0.020%, OR=7.6, $p=1.5\times 10^{-8}$) and was also associated with “congenital anomalies of endocrine gland” (OR=30.8; $p=0.00015$), “disease of pancreas” (OR=5.9; $p=0.0017$), and “hypokalemia” (OR=6.9; $p=0.0046$). Both sites are present in the Genome Aggregation Database (gnomAD)²⁰ and met quality control criteria in our data (average imputation quality >0.7; association signal visible in multiple studies), but their precise contribution to T2D-risk requires further validation.

Causal coding variants. We next considered the 51 signals (of 380) where fine-mapping strongly implicated (PPA>80%) a single causal variant. Eight of these were missense coding variants, six of which fell into established T2D-associated regions (**Supplementary Table 7**). With the exception of p.Cys130Arg at *APOE* (MAF=15.4%), all have been previously-implicated as causal for T2D: p.Ser539Trp in *PAM* (MAF=0.83%); p.Thr139Ile in *HNF4A* (MAF=3.5%); p.Asp1171Asn in *RREB1* (MAF=11.3%); p.Ala146Val in *HNF1A* (MAF=2.9%); and p.Pro446Leu in *GCKR* (MAF=39.3%)³. Coding variant associations at *PATJ* (p.Gly157Val; 9.5% MAF) and *CDKN1B* (p.Val109Gly; 23.5% MAF) are novel and highlight these genes as playing direct roles in T2D-risk. *PATJ* is highly expressed in brain²¹ and encodes Pals1-Associated Tight Junction component, a protein with multiple PDZ domains that mediate protein-

protein interactions. Associations for this variant indicate a central mechanism of action: the T2D-risk allele is associated with obesity in UK Biobank ($OR=1.11$; $p=3.8 \times 10^{-5}$) and the T2D-association signal is attenuated in BMI-adjusted analysis ($p_{diff}=9.3 \times 10^{-10}$). *CDKN1B* encodes a cyclin dependent kinase inhibitor: in mouse, deletion of this gene ameliorates hyperglycemia by increasing islet mass and maintaining compensatory hyperinsulinemia²². There were a further four signals (at *ANKH*, *POC5*, *NEUROG3*, and *ZNF771*) at which a single missense variant accounted for the majority (>50%) of the PPA (**Supplementary Table 7**).

Integration of regulatory annotations to support fine-mapping. Of the 51 variants with PPA>80%, 43 mapped to regulatory sequence: 12 of these were low-frequency or rare, including variants near *ANKH*, *CCND2*, and *WDR72*. To characterise the regulatory impact of these 51 variants, we overlaid them onto chromatin-state maps from T2D-relevant tissues (islets, liver, adipose, skeletal muscle²³⁻²⁵) and transcription factor binding sites^{23,24}. Twenty-eight mapped to islet enhancer or promoter elements; for 14, these chromatin states were islet-specific (**Supplementary Table 8, Supplementary Figure 3**). These data recapitulate previous findings implicating islet regulatory mechanisms at the *CDC123-CAMKD1* (rs11257655) and *MTNRB1* (rs10830963)²⁵⁻²⁷, and indicate that similar molecular mechanisms operate at signals for several other known T2D loci, including *IGF2BP2*, *ANK1*, *GLIS3*, *CDKN2B*, *KCNQ1*, *CCND2*, and *BCL2A*. Novel T2D signals near *ABCB10*, *FAM49A*, *LRFN2*, *CRHR2*, and *CASC11* also overlapped islet-specific enhancers or promoters. High-PPA (i.e. PPA>80%) variants at 13, 10, and 7 signals overlapped enhancers or promoters in adipose, skeletal muscle, and liver, respectively. All but four of these were also enhancers or promoters in islets: one signal (near *GLI2*) mapped to an adipose-specific enhancer, another (near *WDR72*) to a liver-specific enhancer, and two (near *PTGFRN* and *TSC22D2*) to enhancers in both adipose and skeletal muscle.

We next evaluated whether the integration of genome-wide regulatory annotation data could refine mapping resolution at those loci where genetic fine-mapping was less precise²⁵. We focused on regulatory annotations from human islets because: (a) most established T2D-risk variants are considered, given patterns of association to continuous metabolic traits, to act through primary effects on beta-cell function^{3,28,29}; (b) the strongest signal for regulatory enrichment at T2D association signals involves islet-specific regulatory

elements^{23,26}, a view supported by the annotation overlaps of the high-PPA variants described above, and by enrichment analyses we performed using epigenomic annotations from islets, fat, muscle, and liver²⁴ (**Supplementary Figure 4**); and (c) we had access to high-resolution epigenomic and chromatin state annotation maps for human islets combining available histone modification and transcription factor ChIP-seq, ATAC-seq and whole genome bisulphite sequencing²⁵.

Using the hierarchical modelling approach fGWAS³⁰, genome-wide, we observed strong (1.9-8.2 fold), significant (95% confidence not overlapping one) enrichment of T2D-associated variation with respect to multiple islet enhancer and promoter states, as well as coding sequence (with concomitant depletion of heterochromatin states; **Methods, Supplementary Figure 5**). We used the parameter estimates from the joint annotation model (which retained islet enhancers, promoters, and coding sequence, amongst other annotations; **Methods, Supplementary Figure 5**) as priors to redefine 99% credible sets for the 380 distinct T2D association signals amenable to fine-mapping. We circumvented the default assumption in fGWAS of a single casual variant per locus by conducting these analyses on conditionally-decomposed data (noting that this does still allow for the possibility that the association at each conditional signal is distributed across multiple variants on a risk haplotype; **Methods**).

As expected, this integrated fine-mapping analysis boosted PPA for variants overlapping enriched annotations (**Figure 6**). The median 99% credible set size declined from 42 to 32, credible intervals from 116kb to 100kb, and maximum variant PPA per signal climbed by a median of 21%. The number of signals at which the lead variant PPA exceeded 80% increased from 51 to 73, with dramatic improvements at some (e.g. at *GNG4* where the PPA for rs291367 rose from 24.0% to 84.2%; **Figure 3**).

These annotation-supported analyses highlighted seven additional loci (further to the 12 from genetic evidence alone) where the majority of the PPA was invested in a coding variant (**Supplementary Table 7**). Four were novel: *QSER1* (p.Arg1101Cys; MAF=4.3%), *SCD5* (p.Glu197Gln, MAF=33.8%), *IRS2* (p.Gly1057Asp, MAF=34.0%) and *MRPS30* (p.Glu128Gln=MAF 2.8%).

In our recent study of exome-array genotypes, we demonstrated that, for one-third of loci harbouring coding variant associations, a causal role could be excluded once information on local LD and annotation enrichment was incorporated³. For all 19 coding variant signals (at 18 loci) described in this study, present analyses (based on genome-wide data for both discovery and fine-mapping) were consistent with a causal role. These analyses therefore provide additional examples of human validated targets³¹. The value of these targets as leads to therapeutic development will ultimately depend not only on their impact on T2D phenotypes, but also on the consequences of perturbation on other traits, including coronary artery disease (CAD). Among the 19 T2D-associated coding variants, nine were also nominally associated ($p < 0.05$) with CAD³²: at three (*APOE*, *GCKR*, and *RREB1*) opposing effects on T2D- and CAD-predisposition render them less attractive targets (**Supplementary Table 7**).

Next, we concentrated on non-coding variant signals. In the annotation-informed analysis, we identified 15 additional signals (beyond the 43 non-coding signals described above) at which the lead variant PPA exceeded 80% (**Supplementary Table 8**). These signals overlap active islet regulatory sites including strong enhancers (e.g. at *TCF7L2*, *HNF4A*, *ANKH*, *RNF6*, *ZBED3*), active promoters (*EYA2*), weak enhancers (*ADSC2*, *ADCY5*, *CDKN2B*, *TBCE*), and weak promoters (*DGKB*). For many, orthogonal data (e.g. associations with continuous metabolic traits^{3,28,29}, *cis*-eQTL data³³; **Supplementary Table 8**) are consistent with a role in islet function. In contrast, at six signals, including three that are likely, on physiological grounds, to be acting, at least partly, through effects on islets, we saw reductions (10% to 76%) in the lead variant PPA after islet-annotation-informed fGWAS (**Supplementary Table 8**). This occurred when lead variants from the genetic fine-mapping overlapped annotations depleted in the genome-wide model. Examples include variants at primary *CDKAL1* and secondary *KCNQ1* and *INS-IGF2* signals, where the index-variant PPA decreased by 76% (rs7756992), 34% (rs2283164), and 22% (rs555759341), respectively. One explanation is that these represent T2D-association signals where the phenotypic impact on insulin secretion is mediated through long-term consequences of regulatory effects during islet development which are no longer reflected in regulatory annotations seen in mature islets.

At many of these fine-mapped regulatory loci, the integrated data provide novel insights into disease mechanisms: here, we highlight three. At *ST6GAL1*, rs3887925 achieves PPA=98.5% through genetic fine-mapping alone (99.3% in fGWAS), and overlaps enhancers active in islet, as well as liver, adipose, and skeletal muscle (**Supplementary Figure 6**). However, the T2D-risk allele at rs3887925 is associated with an increase in *ST6GAL1* *cis*-expression specific to islets³³ (**Methods, Supplementary Table 8**), consistent with evidence for reduced insulin secretion in risk-allele carriers during provocative testing³⁴. The candidate effector transcript, *ST6GAL1*, encodes β -galactoside α 2,6-sialyltransferase-1, a key enzyme responsible for the biosynthesis of α 2,6-linked sialic acid in N-linked glycans. Altered glycosylation has the potential to impact multiple processes, and global perturbation of *ST6GAL1* has broad effects including, in *St6gal1* knockout mice, increased body weight and visceral fat accumulation³⁵. However, no equivalent association between rs3887925 and anthropometric and lipid phenotypes is seen in human GWAS^{14,36,37}. This is consistent with T2D-predisposition attributable to rs3887925 being mediated through regulatory mechanisms restricted to the modulation of *ST6GAL1* expression in islets.

At *ANK1*, we observed three distinct association signals. The strongest causal variant attribution was for the primary signal at rs13262861 (PPA=97.3% on genetic data alone; 98.8% with fGWAS). This variant overlaps an islet promoter located 3' to *ANK1* and 5' to the transcription factor *NKX6.3* (**Supplementary Figure 7**). The T2D-risk allele at rs13262861 shows a directionally-consistent association with *in vivo* measures of reduced insulin secretion^{3,29,34}, and a *cis*-eQTL for reduced *NKX6-3* expression in human islets (**Supplementary Table 8**). Members of the NKX6 family (including NKX6.3) are implicated in islet development and function³⁸. A recent study highlighted the relationship between variants including rs515071 and rs508419 and the expression and splicing of *ANK1* in skeletal muscle³⁹. However, in our meta-analysis, variants influencing *ANK1* splicing have minimal impact on T2D-risk (PPA<1% in all three conditionally-decomposed signals [genetic fine-mapping only]). Collectively, these data indicate that the mechanism of T2D predisposition at this locus is more likely mediated through reduced islet expression of *NKX6.3* than altered muscle expression of *ANK1*.

At *TCF7L2*, patterns of overlap with epigenomic annotations across the eight distinct T2D-association signals offer explanations for the diverse metabolic consequences of *TCF7L2* perturbation in humans and animal models (**Supplementary Table 9**)⁴⁰. The primary signal at rs7903146, long-established as the largest common variant effect for T2D in Europeans, overlaps an islet enhancer (boosting PPA from 59.2% to 97.1% on fGWAS), multiple islet-relevant transcription factor binding sites, and islet open chromatin⁴¹, all features consistent with the islet phenotype (deficient insulin secretion) evident in non-diabetic individuals⁷ (**Supplementary Figure 8**). However, amongst the seven secondary signals, the picture is more mixed. Of the four secondary signals mapped to <10 credible set variants, only rs144155527 rises to moderate PPA (68%) following islet annotation-enriched fGWAS analysis. Other credible set variants map to adipose and liver enhancers, suggesting that their T2D-risk effects are mediated via modulation of *TCF7L2* expression in tissues relevant to insulin action.

Heritability estimates and polygenic risk score prediction. Using LD score regression⁴², and empirical estimates of population- and sample-level T2D prevalence, we estimated chip heritability (on the liability scale) for T2D at 18% (23% in females and 17% in males; **Supplementary Figure 9**), accounting for approximately half the median estimates of heritability derived from twin and family studies⁴³.

Identification of individuals at increased genetic risk for T2D could enhance screening strategies and allow targeted prevention. Previous attempts to deploy genetic data for disease prediction have shown limited utility^{44,45}. We used a revised BMI-unadjusted meta-analysis, generated from all samples other than UK Biobank, to develop genome-wide polygenic risk scores (PRS)⁴⁶ which we then applied to predict T2D status in the 18,197 cases and 423,697 controls from UK Biobank (Europeans only; **Methods**). Maximal discrimination (AUC C-statistic of 66%, equivalent to that derived from BMI, age, and sex in the same sample) was obtained from a PRS of 136,795 variants ($r^2 > 0.6$, $p < 0.076$; **Supplementary Figure 10**). Individuals in the top 2.5% of the PRS distribution were at 3.4-fold increased risk (prevalence=11.2%) compared to the median (prevalence=3.3%), and 9.4-fold compared to the bottom 2.5% (prevalence=1.2%). Low T2D prevalence rates in UK Biobank reflect the age-distribution of the cohort and preferential ascertainment of healthy individuals:

however, similar prevalence ratios were observed in the subset of individuals aged >55 years at recruitment (14.2% vs 1.6%). Applied to the general UK population, equivalent performance would equate to lifetime risks for T2D of ~59.7% and ~6.7% for individuals from those extremes, based on current UK general population prevalence rates for those >55 years of age⁴⁷.

Defining relationships with other traits. To characterise genetic relationships with other biomedical-relevant traits, we used LD score regression⁴² as implemented in LDHub⁴⁸. We tested 182 unique phenotypes, after excluding those with low heritability estimates and repeated measures. Eighty-five traits demonstrated a significant (Bonferroni corrected threshold $p < 0.00027$) genetic correlation with T2D (**Supplementary Table 10, Supplementary Figure 11**).

These highlighted several interesting genetic correlations, linking increased T2D-risk to sleeping behaviours (insomnia, excessive daytime sleeping), smoking (cigarettes smoked per day, ever versus never smoked), metabolites (glycoprotein acetyls, isoleucine, valine), depressive symptoms, urinary albumin-to-creatinine ratio, and urate. T2D-risk was negatively correlated with anorexia nervosa, intelligence, parent's age at death, lung function measures, education status/duration, age at menarche, and age of mother at first childbirth. Many of these relationships (including those related to intelligence, smoking behaviour, age at menarche, and education status) were primarily mediated by the shared impact of BMI/obesity on both T2D and the correlated phenotype (**Supplementary Figure 12**).

DISCUSSION

This study demonstrates how substantial increases in sample size coupled to more accurate and comprehensive imputation expand characterization of the genetic contribution to T2D-risk. The number of significantly-associated genomic regions has doubled, with a growing harvest of risk-alleles of lower-frequency, some with relatively large effects. At many of these signals, fine-mapping resolution has been substantially improved: we mapped 51 of 380 signals to single-variant resolution on genetic evidence alone, and demonstrated that

the integration of genomic annotations (here focused on the human islet epigenome) provides further specification of plausible causal variants. We highlight 18 genes as human validated targets based on causal coding variants and provide novel insights into the biological mechanisms operating at several fine-mapped regulatory signals. These findings represent mechanistic hypotheses that can now be targeted for large-scale empirical validation at both the level of the variants (e.g. through massively parallel reporter assays) and the candidate effector genes (e.g. through CRISPR screens in appropriate cellular models, and manipulation in *in vivo* models). The present study was limited to individuals of European ancestry: integration of these data with large-scale GWAS data from other major ancestral groups (as is being pursued by the DIAMANTE consortium) should provide an additional boost to locus discovery, and support further increases in causal variant resolution, most obviously at loci where extensive LD within Europeans limits resolution of fine-mapping.

DATA AVAILABILITY STATEMENT

Summary level data is available at the DIAGRAM consortium website <http://diagram-consortium.org/> and Accelerating Medicines Partnership T2D portal <http://www.type2diabetesgenetics.org/>.

ACKNOWLEDGMENTS

This work was supported primarily by NIDDK as part of the Accelerating Medicines Partnership-T2D, funded by U01DK105535, U01DK062370, and U01DK078616 grants. Part of this work was conducted using the UK Biobank resource under application number 9161. A full list of acknowledgments appears in the Supplementary Notes.

AUTHOR CONTRIBUTIONS

Project co-ordination. A.Mahajan, A.P.M., M.B., M.I.M. **Writing.** A.Mahajan, D.T., A.P.M., M.B., M.I.M. **Core analyses.** A.Mahajan, D.T., M.T., J.M.T., A.J.P., A.P.M., M.B., M.I.M. **DIAMANTE analysis group.** A.Mahajan, J.E.B., D.W.B., J.C.C., Y.J.K., M.C.N., L.E.P., X.S., W.Z., A.P.M., M.B., M.I.M. **Statistical analysis in individual studies.** A.Mahajan, D.T., N.R.R., N.W.R., V.S., R.A.S., N.G., J.P.C., E.M.S., M.W., C.Sarnowski, J.N., S.T., C.Lecoecur, M.H.P., B.P.P., X.G., L.F.B., J.B.-J., M.C., K.L., C.-T.L., A.E.L., J'a.L., C.Schurmann, L.Y., G.T., A.P.M. **Genotyping and phenotyping.** A.Mahajan, R.A.S., R.M., C.G., S.T., K.-U.E., K.F., S.L.R.K., F.K., I.N., C.M.B., C.Schurmann, E.P.B., I.B., C.C., G.D., I.f., V.G., M.I., M.E.J., S.L., A.L., V.L., V.M., A.D.M., G.N., N.S., A.S., D.R.W., S.S., E.P.B., S.H., C.H., J.Kriebel, T.M., A.P., B.T., D.A., A.K., G.R.A., C.Langenberg, N.J.W., A.P.M., M.B., M.I.M. **Islet annotations.** M.T., J.M.T., A.J.B., V.N., A.L.G., M.I.M. **Individual study design and principal investigators.** E.P.B., J.C.F., O.H.F., T.M.F., A.T.H., M.A.I., T.J., J.Kuusisto, C.M.L., K.L.M., J.S.P., K.Strauch, K.D.T., U.T., J.T., J.D., P.A.P., E.Z., R.J.F.L., P.F., E.I., L.L., L.G., M.L., F.S.C., J.W.J., C.N.A.P., H.G., A.Metspalu, D.A., A.K., G.R.A., J.B.M., J.I.R., J.M., O.P., T.H., C.Langenberg, N.J.W., K.Stefansson, A.P.M., M.B., M.I.M.

MATERIALS & CORRESPONDENCE

Correspondence and requests for materials should be addressed to anubha@well.ox.ac.uk and mark.mccarthy@drl.ox.ac.uk .

COMPETING FINANCIAL INTERESTS

Jose C Florez has received consulting honoraria from Merck and from Boehringer-Ingelheim. Oscar H Franco works in ErasmusAGE, a center for aging research across the life course funded by Nestlé Nutrition (Nestec Ltd.), Metagenics Inc., and AXA. Erik Ingelsson is a scientific advisor for Precision Wellness and Olink Proteomics for work unrelated to the present project. Abbas Dehghan has received consultancy and research support from Metagenics Inc. (outside the scope of submitted work). Timothy M Frayling has consulted for Boeringer Ingelheim and Sanofi-Aventis on the genetics of diabetes and has an MRC CASE studentship with GSK. Goncalo Abecasis is a consultant for 23andMe, Regeneron, Merck and Helix. Robert A Scott is an employee and shareholder in GlaxoSmithKline. Naveed Sattar is working with Boehringer-Ingelheim on a genetics project but no remuneration received. Mark I McCarthy has served on advisory panels for NovoNordisk and Pfizer, and received honoraria from NovoNordisk, Pfizer, Sanofi-Aventis and Eli Lilly. Companies named above had no role in the design and conduct of the study; collection, management, analysis, and interpretation of the data; and preparation, review or approval of the manuscript. Authors affiliated with deCODE (VS, GT, UT and KS) are employed by deCODE Genetics/Amgen, Inc.

REFERENCES

1. Scott, R.A. *et al.* An Expanded Genome-Wide Association Study of Type 2 Diabetes in Europeans. *Diabetes* **66**, 2888-2902 (2017).
2. Zhao, W. *et al.* Identification of new susceptibility loci for type 2 diabetes and shared etiological pathways with coronary heart disease. *Nat Genet* **49**, 1450-1457 (2017).
3. Mahajan, A. *et al.* Refining the accuracy of validated target identification through coding variant fine-mapping in type 2 diabetes. *Nature Genetics* **50**, 559-571 (2018).
4. McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nat Genet* **48**, 1279-83 (2016).
5. Jonsson, H. *et al.* Whole genome characterization of sequence diversity of 15,220 Icelanders. *Sci Data* **4**, 170115 (2017).
6. Flannick, J. & Florez, J.C. Type 2 diabetes: genetic data sharing to advance complex disease research. *Nat Rev Genet* **17**, 535-49 (2016).
7. Voight, B.F. *et al.* Twelve type 2 diabetes susceptibility loci identified through large-scale association analysis. *Nat Genet* **42**, 579-89 (2010).
8. Morris, A.P. *et al.* Large-scale association analysis provides insights into the genetic architecture and pathophysiology of type 2 diabetes. *Nat Genet* **44**, 981-90 (2012).
9. Kooner, J.S. *et al.* Genome-wide association study in individuals of South Asian ancestry identifies six new type 2 diabetes susceptibility loci. *Nat Genet* **43**, 984-9 (2011).
10. Cho, Y.S. *et al.* Meta-analysis of genome-wide association studies identifies eight new loci for type 2 diabetes in east Asians. *Nat Genet* **44**, 67-72 (2011).
11. Lotta, L.A. *et al.* Integrative genomic analysis implicates limited peripheral adipose storage capacity in the pathogenesis of human insulin resistance. *Nat Genet* **49**, 17-26 (2017).
12. Magi, R., Lindgren, C.M. & Morris, A.P. Meta-analysis of sex-specific genome-wide association studies. *Genet Epidemiol* **34**, 846-53 (2010).
13. Small, K.S. *et al.* Identification of an imprinted master trans regulator at the KLF14 locus related to multiple metabolic phenotypes. *Nat Genet* **43**, 561-4 (2011).

14. Teslovich, T.M. *et al.* Biological, clinical and population relevance of 95 loci for blood lipids. *Nature* **466**, 707-13 (2010).
15. Maller, J.B. *et al.* Bayesian refinement of association signals for 14 loci in 3 common diseases. *Nat Genet* **44**, 1294-301 (2012).
16. Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68-74 (2015).
17. Fuchsberger, C. *et al.* The genetic architecture of type 2 diabetes. *Nature* **536**, 41-7 (2016).
18. Gradwohl, G., Dierich, A., LeMeur, M. & Guillemot, F. neurogenin3 is required for the development of the four endocrine cell lineages of the pancreas. *Proc Natl Acad Sci U S A* **97**, 1607-11 (2000).
19. Rubio-Cabezas, O. *et al.* Permanent Neonatal Diabetes and Enteric Anendocrinosis Associated With Biallelic Mutations in NEUROG3. *Diabetes* **60**, 1349-53 (2011).
20. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285-91 (2016).
21. Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science* **348**, 648-60 (2015).
22. Uchida, T. *et al.* Deletion of Cdkn1b ameliorates hyperglycemia by maintaining compensatory hyperinsulinemia in diabetic mice. *Nat Med* **11**, 175-82 (2005).
23. Pasquali, L. *et al.* Pancreatic islet enhancer clusters enriched in type 2 diabetes risk-associated variants. *Nat Genet* **46**, 136-143 (2014).
24. Varshney, A. *et al.* Genetic regulatory signatures underlying islet gene expression and type 2 diabetes. *Proc Natl Acad Sci U S A* **114**, 2301-2306 (2017).
25. Thurner, M. *et al.* Integration of human pancreatic islet genomic data refines regulatory mechanisms at Type 2 Diabetes susceptibility loci. *Elife* **7**(2018).
26. Gaulton, K.J. *et al.* Genetic fine mapping and genomic annotation defines causal mechanisms at type 2 diabetes susceptibility loci. *Nat Genet* **47**, 1415-25 (2015).
27. Fogarty, M.P., Cannon, M.E., Vadlamudi, S., Gaulton, K.J. & Mohlke, K.L. Identification of a regulatory variant that binds FOXA1 and FOXA2 at the CDC123/CAMK1D type 2 diabetes GWAS locus. *PLoS Genet* **10**, e1004633 (2014).
28. Dimas, A.S. *et al.* Impact of type 2 diabetes susceptibility variants on quantitative glycemic traits reveals mechanistic heterogeneity. *Diabetes* **63**, 2158-71 (2014).

29. Wood, A.R. *et al.* A Genome-Wide Association Study of IVGTT-Based Measures of First-Phase Insulin Secretion Refines the Underlying Physiology of Type 2 Diabetes Variants. *Diabetes* **66**, 2296-2309 (2017).
30. Pickrell J , K. Joint Analysis of Functional Genomic Data and Genome-wide Association Studies of 18 Human Traits. *Am J Hum Genet* **94**, 559-73 (2014).
31. Plenge, R.M., Scolnick, E.M. & Altshuler, D. Validating therapeutic targets through human genetics. *Nat Rev Drug Discov* **12**, 581-94 (2013).
32. van der Harst, P. & Verweij, N. Identification of 64 Novel Genetic Loci Provides an Expanded View on the Genetic Architecture of Coronary Artery Disease. *Circulation Research* **122**, 433-443 (2018).
33. van de Bunt, M. *et al.* Transcript Expression Data from Human Islets Links Regulatory Signals from Genome-Wide Association Studies for Type 2 Diabetes and Glycemic Traits to Their Downstream Effectors. *PLoS Genet* **11**, e1005694 (2015).
34. Prokopenko, I. *et al.* A central role for GRB10 in regulation of islet function in man. *PLoS Genet* **10**, e1004235 (2014).
35. Kaburagi, T., Kizuka, Y., Kitazume, S. & Taniguchi, N. The Inhibitory Role of alpha2,6-Sialylation in Adipogenesis. *J Biol Chem* **292**, 2278-2286 (2017).
36. Locke, A.E. *et al.* Genetic studies of body mass index yield new insights for obesity biology. *Nature* **518**, 197-206 (2015).
37. Shungin, D. *et al.* New genetic loci link adipose and insulin biology to body fat distribution. *Nature* **518**, 187-96 (2015).
38. Lizio, M. *et al.* Mapping Mammalian Cell-type-specific Transcriptional Regulatory Networks Using KD-CAGE and CHIP-seq Data in the TC-YIK Cell Line. *Front Genet* **6**, 331 (2015).
39. Scott, L.J. *et al.* The genetic regulatory signature of type 2 diabetes in human skeletal muscle. *Nat Commun* **7**, 11764 (2016).
40. McCarthy, M.I., Rorsman, P. & Gloyn, A.L. TCF7L2 and diabetes: a tale of two tissues, and of two species. *Cell Metab* **17**, 157-9 (2013).
41. Gaulton, K.J. *et al.* A map of open chromatin in human pancreatic islets. *Nat Genet* **42**, 255-9 (2010).
42. Bulik-Sullivan, B. *et al.* An atlas of genetic correlations across human diseases and traits. *Nat Genet* **47**, 1236-41 (2015).

43. Meigs, J.B., Cupples, L.A. & Wilson, P.W. Parental transmission of type 2 diabetes: the Framingham Offspring Study. *Diabetes* **49**, 2201-7 (2000).
44. Meigs, J.B. *et al.* Genotype score in addition to common risk factors for prediction of type 2 diabetes. *N Engl J Med* **359**, 2208-19 (2008).
45. Weedon, M.N. *et al.* Combining information from common type 2 diabetes risk polymorphisms improves disease prediction. *PLoS Med* **3**, e374 (2006).
46. Euesden, J., Lewis, C.M. & O'Reilly, P.F. PRSice: Polygenic Risk Score software. *Bioinformatics* **31**, 1466-8 (2015).
47. UK government. Adult obesity and type 2 diabetes.
(https://www.gov.uk/government/uploads/system/uploads/attachment_data/file/338934/Adult_obesity_and_type_2_diabetes_.pdf).
48. 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).

FIGURE LEGENDS

Figure 1 | Manhattan plots of the sex-combined BMI-unadjusted and adjusted meta-analysis for T2D. **a**, Manhattan plot (top panel) of genome-wide association results for T2D without BMI adjustment from meta-analysis of up to 71,124 cases and 824,006 controls. The association p -value (on $-\log_{10}$ scale) for each SNP (y -axis) is plotted against the genomic position (NCBI Build 37; x -axis). Association signals that reached genome-wide significance ($p < 5 \times 10^{-8}$) are shown in purple if novel. **b**, Manhattan plot (bottom panel) of genome-wide association results for T2D with BMI adjustment from meta-analysis of up to 50,409 cases and 523,897 controls. Novel association signals that reached genome-wide significance ($p < 5 \times 10^{-8}$) only in the BMI-unadjusted analysis are shown in orange.

Figure 2 | Comparison of estimated T2D effect size between BMI-adjusted and unadjusted models. Z-score for each of the 403 distinct signals from BMI-unadjusted analysis (50,791 cases and 526,121 controls; x -axis) is plotted against the z-score from the BMI-adjusted analysis (50,402 cases and 523,888 controls; y -axis). Variants that display higher T2D effect size in BMI-adjusted analysis are shown in red and variants with higher T2D effects in BMI-unadjusted analysis are shown in blue. Diameter of the circle is proportional to $-\log_{10}$ heterogeneity p -value.

Figure 3 | Summary of fine-mapped associations. **a**, Distinct association signals. A single signal at 151 loci, and 2-10 signals at 92. **b**, Number of variants in genetic and functional 99% credible sets. Eighteen and 23 signals were mapped to a single variant in genetic and functional credible sets, respectively. **c**, Distribution of the posterior probability of association of the variants in credible sets.

Figure 4 | Comparison of fine-mapping resolution at 83 distinct signals. The number of variants included in the 99% credible set for each of the 83 distinct signals constructed using meta-analysis of GWAS data imputed using the 1000G multi-ethnic reference panel (26,676 T2D cases and 132,532 controls) (x -axis; logarithmic scale) is plotted against those (y -axis; logarithmic scale) derived using HRC-based imputation (74,124 T2D cases and 824,006 controls). Inset presents the same plot but with linear scales.

Figure 5 | The relationship between effect size and minor allele frequency. Conditional and joint analysis effect size (y -axis) and minor allele frequency (x -axis) for 403 conditionally independent SNPs. Previously-reported T2D associated variants are shown in green and novel variants are shown in purple. Stars and circles represent the “strongest regional lead at a locus” and “lead variants for secondary signals”, respectively.

Figure 6 | Comparison of posterior probability of association for each variant with and without incorporation enrichment information. Posterior probability of association from genetic credible sets (y -axis) and fGWAS analysis (x -axis) for each variant included in the 99% credible sets.

ONLINE METHODS

Ethics statement. All human research was approved by the relevant institutional review boards, and conducted according to the Declaration of Helsinki. All participants provided written informed consent.

Study-level analyses. We considered a total of 74,124 T2D cases and 824,006 controls from 32 GWAS undertaken in individuals of European ancestry (**Supplementary Table 1**), genotyped with a variety of genome-wide SNP arrays. Sample and variant quality control was performed within each study (**Supplementary Table 1**). To improve the quality of the genotype scaffold in each study, we developed a harmonised protocol in which variants were subsequently removed if: (i) allele frequencies differed from those for European ancestry haplotypes from the HRC reference panel⁴ by more than 20%; AT/GC variants had $MAF > 40\%$ because of potential undetected errors in strand alignment; or (iii) $MAF < 1\%$ because of difficulties in calling rare variants (with exception of BioMe, MGI, and UPCH, where no MAF exclusion was implemented as they were genotyped using newer GWAS arrays; **Supplementary Table 1**). Each scaffold, with exception of the deCODE GWAS, was then imputed up to the HRC reference panel⁴. The GWAS from deCODE was imputed up to a reference panel based on 30,440 Icelandic whole-genome sequences⁵, and only variants that were present on the HRC panel were considered for downstream analyses. Within each study, all variants were tested for association with T2D in a regression framework, with and without adjustment for BMI, in sex-combined and sex-specific analyses, under an additive model in the effect of the minor allele, with additional adjustment for study-specific covariates (**Supplementary Table 1**). To account for population structure and relatedness, association analyses were either adjusted for principal components (after excluding related individuals) or implemented in a mixed model with random effects for kinship from a genetic relationship matrix. For studies analysed using linear mixed models, implemented in EMMAX⁴⁹ or BOLT-LMM⁵⁰ (**Supplementary Table 1**), allelic effects and standard errors were converted to the log-odds scale to correct for case-control imbalance⁵¹. For each analysis, in each study, variants were removed from a study if: (i) minor allele count < 5 (in cases and controls combined); (ii) imputation quality $r^2\text{-hat} < 0.3$ (miniMAC) or $\text{proper-info} < 0.4$

(IMPUTE4); or (iii) standard error of the allelic log-OR>10. Association summary statistics for each analysis within each study were then corrected for residual structure by means of genomic control inflation factor⁵², calculated after excluding variants mapping to established T2D susceptibility loci (**Supplementary Table 1**).

Sex-combined meta-analysis. We aggregated association summary statistics from sex-combined analyses for each variant across studies, with and without adjustment for BMI, using fixed-effects meta-analysis with inverse-variance weighting of log-ORs, implemented in METAL⁵³. The BMI unadjusted meta-analysis was subsequently corrected for residual inflation (to account for structure between studies) by means of genomic control ($\lambda=1.013$)⁵², calculated after excluding variants mapping to established T2D susceptibility loci. No adjustment was required for the BMI adjusted meta-analysis ($\lambda=0.992$). From the meta-analysis, variants were extracted that passed quality control in at least two studies. Heterogeneity in allelic effect sizes between studies contributing to the meta-analysis was assessed by Cochran's *Q* statistic⁵⁴. We defined novel loci as mapping >500kb and conditionally independent from a previously reported lead GWAS SNP.

For the present study, we maintained the conventional genome-wide significance threshold of 5×10^{-8} , for compatibility with previous reports. We recognise that more comprehensive capture of lower-frequency variants in particular increases the effective number of tests, with some consequent increase in the false positive rate for signals just below this threshold. 162 of the 243 primary signals are significant at a more stringent threshold (5×10^{-9}) recently advocated for whole genome sequence data⁵⁵, and the major conclusions of the manuscript are unchanged if we select this more stringent (and given that our data lack the full coverage of WGS data, over-conservative) threshold. We make all summary level data results available so that readers can interpret the results themselves.

With our sample size ($N_{\text{eff}} 231,436$), assuming accurate imputation (imputation quality score >0.8), we had >80% power to detect T2D association (at $\alpha=5 \times 10^{-8}$) with variants of $\text{MAF} \geq 5\%$ and $\text{OR} \geq 1.10$, or $\text{MAF} \geq 0.1\%$ and $\text{OR} \geq 1.60$.

Sex-differentiated meta-analysis. The meta-analyses described above were repeated for males and females separately, with correction for population structure by genomic control as necessary: (i) male-specific BMI unadjusted $\lambda=1.029$; (ii) male-specific BMI adjusted $\lambda=1.001$; (iii) female-specific BMI unadjusted $\lambda=0.955$; and (iv) female-specific BMI adjusted $\lambda=0.932$. The male-specific meta-analysis consisted of up to 41,846 cases and 383,767 controls, whilst the female-specific meta-analysis consisted of up to 30,053 cases and 434,336 controls. The sex-specific meta-analyses were then combined to conduct a sex-differentiated test of association and a test of heterogeneity in allelic effects between males and females¹².

Assessment of effect of BMI adjustment. We compared the genetic effect sizes (beta coefficients) estimated from models with and without BMI adjustments using a matched meta-analysis conducted on the same subset of 28 studies:

$$\frac{\beta_{noBMI} - \beta_{BMI}}{\sqrt{SE(\beta_{noBMI})^2 + SE(\beta_{BMI})^2 - 2\rho \times SE(\beta_{noBMI}) \times SE(\beta_{BMI})}}$$

where β_{BMI} and β_{noBMI} are the estimated genetic effects from models with and without BMI adjustment, $SE(\beta)$ is the estimated standard error of the estimates, and $\rho = 0.89$, is the estimated correlation between β_{BMI} and β_{noBMI} across all variants¹.

Detection of distinct association signals. We used GCTA⁵⁶ to perform approximate conditional analyses to detect distinct association signals at each of the genome-wide significant risk loci for T2D (newly identified or confirmed, except at the major histocompatibility complex (MHC) region). GCTA performs conditional analysis using association summary statistics from GWAS meta-analysis and estimated LD from a sufficiently large reference study used in the meta-analysis. We used a reference sample of 6,000 (nearly) unrelated (pairwise relatedness <0.025) individuals of white British origin, randomly selected from the UK Biobank, to model patterns of LD between variants. The reference panel of genotypes consisted of the same 39 million variants from the HRC reference panel assessed in our GWAS, but with an additional quality control step to exclude SNPs with low imputation quality (proper-info <0.4) or deviation from Hardy-Weinberg

equilibrium ($p < 1 \times 10^{-6}$). For each locus, we first searched $\pm 500\text{kb}$ surrounding the lead SNP (using summary statistics from BMI unadjusted or adjusted analysis, as appropriate) to ensure potential long-range genetic influences were assessed. Within a region, conditionally independent variants that reached locus-wide significance ($p < 10^{-5}$) were considered as index SNPs for distinct association signals. If the minimum distance between any distinct signals from two separate loci was less than 500kb, we performed additional conditional analysis taking both regions (encompassing $\pm 500\text{kb}$ from both ends) and reassessed the independence of each signal.

Fine-mapping of distinct association signals with T2D susceptibility. We considered 380 of the 403 identified distinct signals, excluding 23 that were not amenable to fine-mapping: (i) 19 signals with $\text{MAF} < 0.25\%$; (ii) three signals where the index variant was rare and analysed in $< 50\%$ of the total effective sample size, defined as $N_e = 4 \times N_{\text{cases}} \times N_{\text{controls}} / (N_{\text{cases}} + N_{\text{controls}})$; and (iii) the one signal in the major histocompatibility complex because of the extended and complex structure of LD across the region, which complicates fine-mapping.

For each of the remaining distinct signals, we first defined a genomic region 500kb on either side of the index variant, considering only variants with $\text{MAF} > 0.25\%$ that were reported in at least 50% of the total effective sample size, thus removing those that were not well imputed in the majority of samples. We then adopted two approaches to compute 99% credible sets with 99% posterior probability of containing the causal variant: (i) using a (functionally unweighted) Bayesian approach, with the strength of evidence for association measured using the Bayes' factor in favour of association for each variant^{15,57}; and (ii) using (functionally weighted) fGWAS³⁰ that reweights the association measures by using information from functional genomics data.

(i) Genetic credible sets. For each distinct association signal, we first calculated an approximate Bayes factor⁵⁷ in favour of association on the basis of allelic effect sizes and standard errors from the meta-analysis (using BMI-unadjusted or adjusted meta-analysis, as appropriate). For loci with a single association signal, effect sizes and standard errors from unconditional meta-analysis were used. For loci with multiple distinct association signals,

these parameters were derived from the approximate conditional analysis adjusting for all other index variants in the region. Specifically, for the j th variant,

$$\Lambda_j = \sqrt{\frac{V_j}{V_j + \omega}} \exp \left[\frac{\omega \beta_j^2}{2V_j(V_j + \omega)} \right],$$

where β_j and V_j denote the estimated allelic effect (log-OR) and corresponding variance from the meta-analysis. The parameter ω denotes the prior variance in allelic effects, taken here to be 0.04⁵⁷.

We then calculated the posterior probability that the j th variant drives the association signal (PPA), given by

$$\pi_j = \frac{\Lambda_j}{\sum_k \Lambda_k}.$$

The 99% credible set¹⁵ for each locus was then constructed by: (i) ordering all variants in descending order of their PPA; and (ii) including ordered variants until the cumulative PPA reached 0.99. The number of variants and length of the genomic region covered by each 99% credible set was then calculated.

(ii) Functionally weighted credible sets. We first tested each of the 15 chromatin states in human islets and coding DNA sequence separately for enrichment using genome-wide data with the program fGWAS³⁰. Details on generation of the 15 chromatin states have been described elsewhere²⁵. The annotation with the most significant enrichment was retained and tested jointly with each remaining annotation. If the most significant two annotation model improved the model likelihood then the two annotations in the model were retained and the process continued until the model likelihood did not exceed the previous iteration. The resulting “full” model was iteratively pruned by dropping each annotation and assessing the cross-validated likelihood of the reduced model (i.e. an annotation was removed from the “full” model if dropping it increased the cross-validated likelihood). This process resulted in the “best joint model”.

By default, fGWAS partitions the genome into “blocks” of 5,000 SNPs and assumes no more than one causal variant per block. However, for direct comparison with the “genetic” credible sets and to account for multiple distinct association signals within a locus, we used a modified approach. For T2D associated regions with no evidence of more than one distinct signal, we delineated 1Mb windows comprising all SNPs within 500 kb of the index variant and partitioned the intervening regions into ~1Mb windows. These windows were manually input into fGWAS using the --bed command and a separate fGWAS analysis was performed using only the set of annotations remaining in the “best joint model”. The genome-wide enrichments were used as priors in a Bayesian fine-mapping analysis implemented in fGWAS to calculate posterior probabilities for each SNP in the designated windows. For the remaining regions with evidence of two or more distinct association signals, we used the results from the approximate conditional analyses described above and similarly performed a manually partitioned fGWAS analysis. We then constructed 99% credible sets as described above.

Association analyses with UK Biobank phenotypes. We performed targeted association analyses using genotype and phenotype data from electronic health records (EHRs) from the UK Biobank. Hierarchical phenotype codes from EHRs were curated by grouping International Classification of Disease, Ninth Revision (ICD-9) clinical/billing codes as previously described⁵⁸. Only phenotype codes with 20 or more cases and with minor allele count of 5 or greater in cases and controls were considered eligible for analysis. Logistic regression analyses were performed in individuals of European ancestry for relevant phenotype-genotype combination adjusting for six genetic ancestry principal components, array, and sex.

Estimating phenotypic variance explained by SNPs. We used UK Biobank samples (19,119 T2D cases and 423,698 controls) to calculate variance explained by genome-wide significant variants. We ran a model regressing T2D status on all independently associated rare and low-frequency variants, assuming an additive model (and adjusting for sex, age, array, and 6 principal components). A separate model was run to determine the variance captured by the independently associated common variants.

Co-localisation analysis. We used publicly available eQTL results from GTEx version 7 for adipose, liver and skeletal tissues. Islet eQTLs were called using published imputed genotype and aligned RNA-seq data (.vcf and .bam files) from human pancreatic islets of 118 individuals, downloaded from the European Genome-phenome Archive (access number EGAS00001001265). RNA extraction, sequencing, and mapping, as well as DNA extraction, genotyping, imputation, and variant filtering was performed as described previously³³. Gene-level reads were quantified using featureCounts version 1.5.0-p2⁵⁹, based on a patched version of GENCODE 19 published by the GTEx Consortium. Quantified gene level read counts for pancreatic islets were filtered in line with protocols used for GTEx version 7: only genes with at least 6 raw counts in 20% of the samples and TPM>0.1 in at least 20% of the samples were used for analysis. Gene-level counts for remaining genes were converted to counts per million and library sizes normalised using edgeR version 3.16.5⁶⁰, and resulting expression values were rank inverse normalised per-gene. Fifteen PEER factors⁶¹ were calculated, and *cis*-eQTLs were called using FastQTL version 2.0⁶² using a *cis* distance of 1Mb and PEER factors as covariates.

We performed co-localisation analysis using eCAVIAR version 2.0⁶³. Co-localisation was performed for each locus-tissue pair using genetic credible set variants from the locus that had: i) PPA>0.01, ii) correlation data from 1000G, and iii) available eQTL results from that tissue. Pairwise variant correlations between credible set SNPs were calculated with PLINK version 1.9⁶⁴ using the 1000 Genomes Project genotypes (phase 3, October 2014 release)¹⁶. Final co-localisation results were filtered to include only variant-gene pairs with significant eQTL effects, which were defined as associations with FDR<0.05 for islets, or published significant associations based on permuted *p*-values for GTEx. For a credible set variant, an eGene with co-localisation posterior probability >0.20 was considered as target gene.

Estimation of genetic variance explained. We used Linkage Disequilibrium Score Regression (LDSC)⁴² to estimate the proportion of variance explained by common genetic variants for T2D on the liability scale. As advised by the developers, these estimates were based on summary statistics (without any genomic control correction) of variants restricted to the subset of HapMap⁶⁵ variants after excluding the MHC region. Estimations were done for

both sex-combined and sex-specific (BMI-unadjusted) analyses, assuming population prevalence of 10%.

Polygenic risk score analyses. Polygenic risk score (PRSs) were created in UK Biobank samples using raw genotype data using the software PRSice⁴⁶ using the GWAS summary statistics of 4.6M common variants from the sex-combined BMI-unadjusted T2D meta-analysis excluding UK Biobank samples. PRSs were created using p -value thresholds ranging from 5×10^{-8} to 0.5 using LD pruning parameters of r^2 0.2 to 0.8 over 250 kb windows. We then tested each PRS for classification performance in UK Biobank.

Reporting summary. Information on experimental design is available in the accompanying Life Sciences Reporting Summary file.

49. Kang, H.M. *et al.* Variance component model to account for sample structure in genome-wide association studies. *Nat Genet* **42**, 348-54 (2010).
50. Loh, P.R. *et al.* Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat Genet* **47**, 284-90 (2015).
51. Cook, J.P., Mahajan, A. & Morris, A.P. Guidance for the utility of linear models in meta-analysis of genetic association studies of binary phenotypes. *Eur J Hum Genet* **25**, 240-245 (2017).
52. Devlin, B. & Roeder, K. Genomic control for association studies. *Biometrics* **55**, 997-1004 (1999).
53. Willer, C.J., Li, Y. & Abecasis, G.R. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* **26**, 2190-1 (2010).
54. Ioannidis, J.P., Patsopoulos, N.A. & Evangelou, E. Heterogeneity in meta-analyses of genome-wide association investigations. *PLoS One* **2**, e841 (2007).
55. Pulit, S.L., de With, S.A. & de Bakker, P.I. Resetting the bar: Statistical significance in whole-genome sequencing-based association studies of global populations. *Genet Epidemiol* **41**, 145-151 (2017).
56. Yang, J. *et al.* Conditional and joint multiple-SNP analysis of GWAS summary statistics identifies additional variants influencing complex traits. *Nat Genet* **44**, 369-75, s1-3 (2012).

57. Wakefield, J. A Bayesian measure of the probability of false discovery in genetic epidemiology studies. *Am J Hum Genet* **81**, 208-27 (2007).
58. Denny, J.C. *et al.* PheWAS: demonstrating the feasibility of a phenome-wide scan to discover gene-disease associations. *Bioinformatics* **26**, 1205-10 (2010).
59. Liao, Y., Smyth, G.K. & Shi, W. featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics* **30**, 923-30 (2014).
60. Robinson, M.D., McCarthy, D.J. & Smyth, G.K. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139-40 (2010).
61. Stegle, O., Parts, L., Piipari, M., Winn, J. & Durbin, R. Using probabilistic estimation of expression residuals (PEER) to obtain increased power and interpretability of gene expression analyses. *Nat Protoc* **7**, 500-7 (2012).
62. Ongen, H., Buil, A., Brown, A.A., Dermitzakis, E.T. & Delaneau, O. Fast and efficient QTL mapper for thousands of molecular phenotypes. *Bioinformatics* **32**, 1479-85 (2016).
63. Hormozdiari, F. *et al.* Colocalization of GWAS and eQTL Signals Detects Target Genes. *American Journal of Human Genetics* **99**, 1245-1260 (2016).
64. Purcell, S. *et al.* PLINK: A Tool Set for Whole-Genome Association and Population-Based Linkage Analyses. *American Journal of Human Genetics* **81**, 559-575 (2007).
65. Frazer, K.A. *et al.* A second generation human haplotype map of over 3.1 million SNPs. *Nature* **449**, 851-61 (2007).

URLs

UK Biobank, <http://www.ukbiobank.ac.uk/>; MACH, <http://csg.sph.umich.edu//abecasis/MaCH/>; SHAPEIT, https://mathgen.stats.ox.ac.uk/genetics_software/shapeit/shapeit.html; GTEx, <http://www.gtexportal.org/home/>; Locuszoom, <http://locuszoom.sph.umich.edu/locuszoom/>; 1000 Genomes Project, <http://www.1000genomes.org/>; HapMap project, <http://hapmap.ncbi.nlm.nih.gov/>; HRC, <http://www.haplotype-reference-consortium.org/>; GCTA, <http://cnsgenomics.com/software/gcta/>; LDSC, <https://github.com/bulik/ldsc/>; LD Hub,

<http://ldsc.broadinstitute.org/>; bedtools, <http://bedtools.readthedocs.io/en/latest/>;
DIAGRAM Consortium, <http://diagram-consortium.org/>; fGWAS
<https://github.com/joepickrell/fgwas>.