

Refining the accuracy of validated target identification through coding variant fine- mapping in type 2 diabetes

Anubha Mahajan¹, Jennifer Wessel², Sara M Willems³, Wei Zhao⁴, Neil R Robertson^{1,5},
Audrey Y Chu^{6,7}, Wei Gan¹, Hidetoshi Kitajima¹, Daniel Taliun⁸, N William Rayner^{1,5,9}, Xiuqing
Guo¹⁰, Yingchang Lu¹¹, Man Li^{12,13}, Richard A Jensen¹⁴, Yao Hu¹⁵, Shaofeng Huo¹⁵, Kurt K
Lohman¹⁶, Weihua Zhang^{17,18}, James P Cook¹⁹, Bram Peter Prins⁹, Jason Flannick^{20,21}, Niels
Grarup²², Vassily Vladimirovich Trubetskoy⁸, Jasmina Kravic²³, Young Jin Kim²⁴, Denis V
Rybin²⁵, Hanieh Yaghooskar²⁶, Martina Müller-Nurasyid^{27,28,29}, Karina Meidtner^{30,31}, Ruifang
Li-Gao^{32,33}, Tibor V Varga³⁴, Jonathan Marten³⁵, Jin Li³⁶, Albert Vernon Smith^{37,38}, Ping An³⁹,
Symen Ligthart⁴⁰, Stefan Gustafsson⁴¹, Giovanni Malerba⁴², Ayse Demirkan^{40,43}, Juan
Fernandez Tajés¹, Valgerdur Steinthorsdottir⁴⁴, Matthias Wuttke⁴⁵, Cécile Lecoeur⁴⁶, Michael
Preuss¹¹, Lawrence F Bielak⁴⁷, Marielisa Graff⁴⁸, Heather M Highland⁴⁹, Anne E Justice⁴⁸,
Dajiang J Liu⁵⁰, Eirini Marouli⁵¹, Gina Marie Peloso^{20,25}, Helen R Warren^{51,52}, ExomeBP
Consortium⁵³, MAGIC Consortium⁵³, GIANT consortium⁵³, Saima Afaq¹⁷, Shoaib Afzal^{54,55,56},
Emma Ahlqvist²³, Peter Almgren⁵⁷, Najaf Amin⁴⁰, Lia B Bang⁵⁸, Alain G Bertoni⁵⁹, Cristina
Bombieri⁴², Jette Bork-Jensen²², Ivan Brandslund^{60,61}, Jennifer A Brody¹⁴, Noël P Burt²⁰,
Mickaël Canouil⁴⁶, Yii-Der Ida Chen¹⁰, Yoon Shin Cho⁶², Cramer Christensen⁶³, Sophie V
Eastwood⁶⁴, Kai-Uwe Eckardt⁶⁵, Krista Fischer⁶⁶, Giovanni Gambaro⁶⁷, Vilmantas Giedraitis⁶⁸,
Megan L Grove⁶⁹, Hugoline G de Haan³³, Sophie Hackinger⁹, Yang Hai¹⁰, Sohee Han²⁴, Anne
Tybjaerg-Hansen^{55,56,70}, Marie-France Hivert^{71,72,73}, Bo Isomaa^{74,75}, Susanne Jäger^{30,31}, Marit E
Jørgensen^{76,77}, Torben Jørgensen^{56,78,79}, Annemari Käräjämäki^{80,81}, Bong-Jo Kim²⁴, Sung Soo
Kim²⁴, Heikki A Koistinen^{82,83,84,85}, Peter Kovacs⁸⁶, Jennifer Kriebel^{31,87}, Florian Kronenberg⁸⁸,
Kristi Läll^{66,89}, Leslie A Lange⁹⁰, Jung-Jin Lee⁴, Benjamin Lehne¹⁷, Huaixing Li¹⁵, Keng-Hung
Lin⁹¹, Allan Linneberg^{78,92,93}, Ching-Ti Liu²⁵, Jun Liu⁴⁰, Marie Loh^{17,94,95}, Reedik Mägi⁶⁶, Vasiliki
Mamakou⁹⁶, Roberta McKean-Cowdin⁹⁷, Girish Nadkarni⁹⁸, Matt Neville^{5,99}, Sune F
Nielsen^{54,55,56}, Ioanna Ntalla⁵¹, Patricia A Peyser¹⁰⁰, Wolfgang Rathmann^{31,101}, Kenneth
Rice¹⁰², Stephen S Rich¹⁰³, Line Rode^{54,55}, Olov Rolandsson¹⁰⁴, Sebastian Schönherr⁸⁸,

30 Elizabeth Selvin¹², Kerrin S Small¹⁰⁵, Alena Stančáková¹⁰⁶, Praveen Surendran¹⁰⁷, Kent D
 31 Taylor¹⁰, Tanya M Teslovich⁸, Barbara Thorand^{31,108}, Gudmar Thorleifsson⁴⁴, Adrienne Tin¹⁰⁹,
 32 Anke Tönjes¹¹⁰, Anette Varbo^{54,55,56,70}, Daniel R Witte^{111,112}, Andrew R Wood²⁶, Pranav
 33 Yajnik⁸, Jie Yao¹⁰, Loïc Yengo⁴⁶, Robin Young^{107,113}, Philippe Amouyel¹¹⁴, Heiner Boeing¹¹⁵,
 34 Eric Boerwinkle^{69,116}, Erwin P Bottinger¹¹, Rajiv Chowdhury¹¹⁷, Francis S Collins¹¹⁸, George
 35 Dedoussis¹¹⁹, Abbas Dehghan^{40,120}, Panos Deloukas^{51,121}, Marco M Ferrario¹²², Jean
 36 Ferrières^{123,124}, Jose C Florez^{71,125,126,127}, Philippe Frossard¹²⁸, Vilmundur Gudnason^{37,38},
 37 Tamara B Harris¹²⁹, Susan R Heckbert¹³⁰, Joanna M M Howson¹¹⁷, Martin Ingelsson⁶⁸, Sekar
 38 Kathiresan^{20,127,131,132}, Frank Kee¹³³, Johanna Kuusisto¹⁰⁶, Claudia Langenberg³, Lenore J
 39 Launer¹²⁹, Cecilia M Lindgren^{1,20,134}, Satu Männistö¹³⁵, Thomas Meitinger^{136,137}, Olle
 40 Melander⁵⁷, Karen L Mohlke¹³⁸, Marie Moitry^{139,140}, Andrew D Morris^{141,142}, Alison D
 41 Murray¹⁴³, Renée de Mutsert³³, Marju Orho-Melander¹⁴⁴, Katharine R Owen^{5,99}, Markus
 42 Perola^{135,145}, Annette Peters^{29,31,108}, Michael A Province³⁹, Asif Rasheed¹²⁸, Paul M Ridker^{7,127},
 43 Fernando Rivadineira^{40,146}, Frits R Rosendaal³³, Anders H Rosengren²³, Veikko Salomaa¹³⁵,
 44 Wayne H -H Sheu¹⁴⁷, Rob Sladek^{148,149,150}, Blair H Smith¹⁵¹, Konstantin Strauch^{27,152}, André G
 45 Uitterlinden^{40,146}, Rohit Varma¹⁵³, Cristen J Willer^{154,155,156}, Matthias Blüher^{86,110}, Adam S
 46 Butterworth^{107,157}, John Campbell Chambers^{17,18,158}, Daniel I Chasman^{7,127}, John
 47 Danesh^{107,157,159,160}, Cornelia van Duijn⁴⁰, Josée Dupuis^{6,25}, Oscar H Franco⁴⁰, Paul W
 48 Franks^{34,104,161}, Philippe Froguel^{46,162}, Harald Grallert^{31,87,163,164}, Leif Groop^{23,145}, Bok-Ghee
 49 Han²⁴, Torben Hansen^{22,165}, Andrew T Hattersley¹⁶⁶, Caroline Hayward³⁵, Erik Ingelsson^{41,167},
 50 Sharon LR Kardia¹⁶⁸, Fredrik Karpe^{5,99}, Jaspal Singh Kooner^{18,158,169}, Anna Köttgen⁴⁵, Kari
 51 Kuulasmaa¹³⁵, Markku Laakso¹⁰⁶, Xu Lin¹⁵, Lars Lind¹⁷⁰, Yongmei Liu⁵⁹, Ruth J F Loos^{11,171},
 52 Jonathan Marchini^{1,172}, Andres Metspalu⁶⁶, Dennis Mook-Kanamori^{33,173}, Børge G
 53 Nordestgaard^{54,55,56}, Colin N A Palmer¹⁷⁴, James S Pankow¹⁷⁵, Oluf Pedersen²², Bruce M
 54 Psaty^{176,177}, Rainer Rauramaa¹⁷⁸, Naveed Sattar¹⁷⁹, Matthias B Schulze^{30,31}, Nicole
 55 Soranzo^{9,157,180}, Timothy D Spector¹⁰⁵, Kari Stefansson^{38,44}, Michael Stumvoll¹⁸¹, Unnur
 56 Thorsteinsdottir^{38,44}, Tiinamaija Tuomi^{75,83,145,182}, Jaakko Tuomilehto^{82,183,184,185}, Nicholas J
 57 Wareham³, James G Wilson¹⁸⁶, Eleftheria Zeggini⁹, Robert A Scott³, Inês Barroso^{9,187},
 58 Timothy M Frayling²⁶, Mark O Goodarzi¹⁸⁸, James B Meigs¹⁸⁹, Michael Boehnke⁸, Danish
 59 Saleheen^{4,128,*}, Andrew P Morris^{1,19,66,*}, Jerome I Rotter^{190,*}, Mark I McCarthy^{1,5,99,*}
 60

- 61 1. Wellcome Trust Centre for Human Genetics, Nuffield Department of Medicine,
62 University of Oxford, Oxford, OX3 7BN, UK.
- 63 2. Departments of Epidemiology and Medicine, Diabetes Translational Research Center,
64 Indiana University, Indianapolis, IN, 46202-2872, USA.
- 65 3. MRC Epidemiology Unit, Institute of Metabolic Science, University of Cambridge,
66 Cambridge, CB2 0QQ, UK.
- 67 4. Department of Biostatistics and Epidemiology, University of Pennsylvania,
68 Philadelphia, Pennsylvania, 19104, USA.
- 69 5. Oxford Centre for Diabetes, Endocrinology and Metabolism, Radcliffe Department of
70 Medicine, University of Oxford, Oxford, OX3 7LE, UK.
- 71 6. National Heart, Lung, and Blood Institute's Framingham Heart Study, Framingham,
72 Massachusetts, 01702, USA.
- 73 7. Division of Preventive Medicine, Department of Medicine, Brigham and Women's
74 Hospital, Boston, MA, 02215, USA.
- 75 8. Department of Biostatistics and Center for Statistical Genetics, University of Michigan,
76 Ann Arbor, Michigan, 48109, USA.
- 77 9. Department of Human Genetics, Wellcome Trust Sanger Institute, Hinxton,
78 Cambridgeshire, CB10 1SA, UK.
- 79 10. Department of Pediatrics, The Institute for Translational Genomics and Population
80 Sciences, LABioMed at Harbor-UCLA Medical Center, Torrance, California, 90502, US.
- 81 11. The Charles Bronfman Institute for Personalized Medicine, The Icahn School of
82 Medicine at Mount Sinai, New York, 10029, USA.
- 83 12. Department of Epidemiology, Johns Hopkins Bloomberg School of Public Health,
84 Baltimore, Maryland, 21205, US.
- 85 13. Division of Nephrology and Hypertension, Department of Internal Medicine, University
86 of Utah School of Medicine, Salt Lake City, Utah, 84132, US.
- 87 14. Cardiovascular Health Research Unit, Department of Medicine, University of
88 Washington, Seattle, WA, 98101, USA.
- 89 15. Institute for Nutritional Sciences, Shanghai Institutes for Biological Sciences, Chinese
90 Academy of Sciences, University of the Chinese Academy of Sciences, Shanghai,
91 People's Republic of China.

- 92 16. Department of Biostatistical Sciences, Division of Public Health Sciences, Wake Forest
93 University Health Sciences, Winston Salem, North Carolina, 27157, USA.
- 94 17. Department of Epidemiology and Biostatistics, Imperial College London, London, W2
95 1PG, UK.
- 96 18. Department of Cardiology, Ealing Hospital, London North West Healthcare NHS Trust,
97 Middlesex, UB1 3HW, UK.
- 98 19. Department of Biostatistics, University of Liverpool, Liverpool, L69 3GA, UK.
- 99 20. Program in Medical and Population Genetics, Broad Institute, Cambridge,
100 Massachusetts, 02142, USA.
- 101 21. Department of Molecular Biology, Massachusetts General Hospital, Boston,
102 Massachusetts, 02114, USA.
- 103 22. The Novo Nordisk Foundation Center for Basic Metabolic Research, Faculty of Health
104 and Medical Sciences, University of Copenhagen, Copenhagen, 2200, Denmark.
- 105 23. Department of Clinical Sciences, Diabetes and Endocrinology, Lund University Diabetes
106 Centre, Malmö, 20502, Sweden.
- 107 24. Center for Genome Science, Korea National Institute of Health, Chungcheongbuk-do,
108 Republic of Korea.
- 109 25. Department of Biostatistics, Boston University School of Public Health, Boston,
110 Massachusetts, 02118, USA.
- 111 26. Genetics of Complex Traits, University of Exeter Medical School, University of Exeter,
112 Exeter, EX1 2LU, UK.
- 113 27. Institute of Genetic Epidemiology, Helmholtz Zentrum München, German Research
114 Center for Environmental Health, Neuherberg, 85764, Germany.
- 115 28. Department of Medicine I, University Hospital Grosshadern, Ludwig-Maximilians-
116 Universität, Munich, 81377, Germany.
- 117 29. DZHK (German Centre for Cardiovascular Research), partner site Munich Heart
118 Alliance, Munich, 81675, Germany.
- 119 30. Department of Molecular Epidemiology, German Institute of Human Nutrition
120 Potsdam-Rehbruecke (DIfE), Nuthetal, 14558, Germany.
- 121 31. German Center for Diabetes Research (DZD), Neuherberg, 85764, Germany.
- 122 32. Department of Clinical Epidemiology, Leiden, 2300 RC, The Netherlands.

- 123 33. Department of Clinical Epidemiology, Leiden University Medical Center, Leiden, 2300
124 RC, The Netherlands.
- 125 34. Department of Clinical Sciences, Lund University Diabetes Centre, Genetic and
126 Molecular Epidemiology Unit, Lund University, Malmö, SE-214 28, Sweden.
- 127 35. MRC Human Genetics Unit, Institute of Genetics and Molecular Medicine, University
128 of Edinburgh, Edinburgh, EH4 2XU, UK.
- 129 36. Division of Cardiovascular Medicine, Department of Medicine, Stanford University
130 School of Medicine, Palo Alto, CA, 94304, US.
- 131 37. Icelandic Heart Association, Kopavogur, 201, Iceland.
- 132 38. Faculty of Medicine, University of Iceland, Reykjavik, 101, Iceland.
- 133 39. Department of Genetics Division of Statistical Genomics, Washington University
134 School of Medicine, St. Louis, Missouri, 63110, USA.
- 135 40. Department of Epidemiology, Erasmus University Medical Center, Rotterdam, 3015CN,
136 The Netherlands.
- 137 41. Department of Medical Sciences, Molecular Epidemiology and Science for Life
138 Laboratory, Uppsala University, Uppsala, 75185, Sweden.
- 139 42. Section of Biology and Genetics, Department of Neurosciences, Biomedicine and
140 Movement sciences, University of Verona, Verona, 37134, Italy.
- 141 43. Department of Human Genetics, Leiden University Medical Center, Leiden,
142 Netherlands.
- 143 44. deCODE Genetics, Amgen inc., Reykjavik, 101, Iceland.
- 144 45. Institute of Genetic Epidemiology, Medical Center – University of Freiburg, Faculty of
145 Medicine, University of Freiburg, Freiburg, 79106, Germany.
- 146 46. CNRS-UMR8199, Lille University, Lille Pasteur Institute, Lille, 59000, France.
- 147 47. Department of Epidemiology, School of Public Health, University of Michigan, Ann
148 Arbor, Michigan, 48109, USA.
- 149 48. Department of Epidemiology, University of North Carolina, Chapel Hill, NC, 27514,
150 USA.
- 151 49. Human Genetics Center, The University of Texas Graduate School of Biomedical
152 Sciences at Houston, The University of Texas Health Science Center at Houston,
153 Houston, Texas, 77030, USA.

- 154 50. Department of Public Health Sciences, Institute of Personalized Medicine, Penn State
155 College of Medicine, Hershey, PA, USA.
- 156 51. William Harvey Research Institute, Barts and The London School of Medicine and
157 Dentistry, Queen Mary University of London, London, UK.
- 158 52. National Institute for Health Research, Barts Cardiovascular Biomedical Research Unit,
159 Queen Mary University of London, London, London, EC1M 6BQ, UK.
- 160 53. The members of this consortium and their affiliations are listed in the Supplementary
161 Note.
- 162 54. Department of Clinical Biochemistry, Herlev and Gentofte Hospital, Copenhagen
163 University Hospital, Herlev, 2730, Denmark.
- 164 55. The Copenhagen General Population Study, Herlev and Gentofte Hospital,
165 Copenhagen University Hospital, Copenhagen, DK-2730, Denmark.
- 166 56. Faculty of Health and Medical Sciences, University of Copenhagen, Copenhagen,
167 Denmark.
- 168 57. Department of Clinical Sciences, Hypertension and Cardiovascular Disease, Lund
169 University, Malmö, 20502, Sweden.
- 170 58. Department of Cardiology, Rigshospitalet, Copenhagen University Hospital,
171 Copenhagen, 2100, Denmark.
- 172 59. Department of Epidemiology & Prevention, Public Health Sciences, Wake Forest
173 University Health Sciences, Winston-Salem, NC, 27157-1063, USA.
- 174 60. Institute of Regional Health Research, University of Southern Denmark, Odense, 5000,
175 Denmark.
- 176 61. Department of Clinical Biochemistry, Vejle Hospital, Vejle, 7100, Denmark.
- 177 62. Department of Biomedical Science, Hallym University, Chuncheon, Republic of Korea.
- 178 63. Medical Department, Lillebælt Hospital Vejle, Vejle, Denmark.
- 179 64. Institute of Cardiovascular Science, University College London, London, WC1E 6BT.
- 180 65. Department of Nephrology and Medical Intensive Care Charité, University Medicine
181 Berlin, Berlin, 10117, Germany.
- 182 66. Estonian Genome Center, University of Tartu, Tartu, 51010, Estonia.
- 183 67. Università Cattolica del Sacro Cuore, Roma, 00168, Italy.
- 184 68. Department of Public Health and Caring Sciences, Geriatrics, Uppsala University,
185 Uppsala, SE-751 85, Sweden.

- 186 69. Human Genetics Center, Department of Epidemiology, Human Genetics, and
187 Environmental Sciences, School of Public Health, The University of Texas Health
188 Science Center at Houston, Houston, Texas, USA.
- 189 70. Department of Clinical Biochemistry, Rigshospitalet, Copenhagen University Hospital,
190 Copenhagen, 2100, Denmark.
- 191 71. Diabetes Research Center (Diabetes Unit), Department of Medicine, Massachusetts
192 General Hospital, Boston, Massachusetts, 02114, USA.
- 193 72. Department of Population Medicine, Harvard Pilgrim Health Care Institute, Harvard
194 Medical School, Boston, MA, 02215, USA.
- 195 73. Department of Medicine, Universite de Sherbrooke, Sherbrooke, QC, J1K 2R1, Canada.
- 196 74. Malmska Municipal Health Care Center and Hospital, Jakobstad, 68601, Finland.
- 197 75. Folkhälsan Research Centre, Helsinki, 00014, Finland.
- 198 76. Steno Diabetes Center Copenhagen, Gentofte, 2820, Denmark.
- 199 77. National Institute of Public Health, Southern Denmark University, Copenhagen, 1353,
200 Denmark.
- 201 78. Research Centre for Prevention and Health, Capital Region of Denmark, Glostrup,
202 2600, Denmark.
- 203 79. Faculty of Medicine, Aalborg University, Aalborg, Denmark.
- 204 80. Department of Primary Health Care, Vaasa Central Hospital, Vaasa, Finland.
- 205 81. Diabetes Center, Vaasa Health Care Center, Vaasa, Finland.
- 206 82. Department of Health, National Institute for Health and Welfare, Helsinki, 00271,
207 Finland.
- 208 83. Endocrinology, Abdominal Center, Helsinki University Hospital, Helsinki, Finland,
209 00029.
- 210 84. Minerva Foundation Institute for Medical Research, Helsinki, Finland.
- 211 85. Department of Medicine, University of Helsinki and Helsinki University Central
212 Hospital, Helsinki, Finland.
- 213 86. Integrated Research and Treatment (IFB) Center AdiposityDiseases, University of
214 Leipzig, Leipzig, 04103, Germany.
- 215 87. Research Unit of Molecular Epidemiology, Institute of Epidemiology II, Helmholtz
216 Zentrum München Research Center for Environmental Health, Neuherberg, 85764,
217 Germany.

- 218 88. Division of Genetic Epidemiology, Department of Medical Genetics, Molecular and
219 Clinical Pharmacology, Medical University of Innsbruck, Innsbruck, 6020, Austria.
- 220 89. Institute of Mathematical Statistics, University of Tartu, Tartu, Estonia.
- 221 90. Department of Medicine, Division of Bioinformatics and Personalized Medicine,
222 University of Colorado Denver, Aurora, CO, USA, 80045.
- 223 91. Department of Ophthalmology, Taichung Veterans General Hospital, Taichung, 40705,
224 Taiwan.
- 225 92. Department of Clinical Experimental Research, Rigshospitalet, Glostrup, Denmark.
- 226 93. Department of Clinical Medicine, Faculty of Health and Medical Sciences, University of
227 Copenhagen, Copenhagen, Denmark.
- 228 94. Institute of Health Sciences, University of Oulu, Oulu, 90014, Finland.
- 229 95. Translational Laboratory in Genetic Medicine (TLGM), Agency for Science, Technology
230 and Research (A*STAR), Singapore, 138648, Singapore.
- 231 96. Dromokaiteio Psychiatric Hospital, National and Kapodistrian University of Athens,
232 Athens, Greece.
- 233 97. Department of Preventive Medicine, Keck School of Medicine of the University of
234 Southern California, Los Angeles, California, 90007, US.
- 235 98. Division of Nephrology, Department of Medicine, Icahn School of Medicine at Mount
236 Sinai, New York, NY, 10069, USA.
- 237 99. Oxford NIHR Biomedical Research Centre, Oxford University Hospitals Trust, Oxford,
238 OX3 7LE, UK.
- 239 100. Department of Epidemiology, School of Public Health, University of Michigan, Ann
240 Arbor, Michigan, 48109, USA.
- 241 101. Institute for Biometrics and Epidemiology, German Diabetes Center, Leibniz Center
242 for Diabetes Research at Heinrich Heine University Düsseldorf, Düsseldorf, Germany.
- 243 102. Department of Biostatistics, University of Washington, Seattle, WA, 98195-7232, USA.
- 244 103. Center for Public Health Genomics, Department Public Health Sciences, University of
245 Virginia School of Medicine, Charlottesville, Virginia, 22908, US.
- 246 104. Department of Public Health and Clinical Medicine, Umeå University, Umeå, 90187,
247 Sweden.
- 248 105. Department of Twin Research and Genetic Epidemiology, King's College London,
249 London, SE1 7EH, UK.

- 250 106. Institute of Clinical Medicine, Internal Medicine, University of Eastern Finland and
251 Kuopio University Hospital, Kuopio, 70210, Finland.
- 252 107. MRC/BHF Cardiovascular Epidemiology Unit, Department of Public Health and Primary
253 Care, University of Cambridge, Cambridge, CB1 8RN, UK.
- 254 108. Institute of Epidemiology II, Helmholtz Zentrum München, German Research Center
255 for Environmental Health, Neuherberg, 85764, Germany.
- 256 109. Welch Center for Prevention, Epidemiology, and Clinical Research, Johns Hopkins
257 Bloomberg School of Public Health, Baltimore, Maryland, USA.
- 258 110. Department of Medicine, University of Leipzig, Leipzig, 04103, Germany.
- 259 111. Department of Public Health, Aarhus University, Aarhus, Denmark.
- 260 112. Danish Diabetes Academy, Odense, Denmark.
- 261 113. Robertson Centre for Biostatistics, University of Glasgow, Glasgow, UK.
- 262 114. Institut Pasteur de Lille, INSERM U1167, Université Lille Nord de France, Lille, F-59000,
263 France.
- 264 115. Department of Epidemiology, German Institute of Human Nutrition Potsdam-
265 Rehbruecke (DIfE), Nuthetal, 14558, Germany.
- 266 116. Human Genome Sequencing Center, Baylor College of Medicine, Houston, Texas,
267 77030, US.
- 268 117. Department of Public Health and Primary Care, University of Cambridge, Cambridge,
269 CB1 8RN, UK.
- 270 118. Genome Technology Branch, National Human Genome Research Institute, National
271 Institutes of Health, Bethesda, Maryland, 20892, USA.
- 272 119. Department of Nutrition and Dietetics, Harokopio University of Athens, Athens,
273 17671, Greece.
- 274 120. MRC-PHE Centre for Environment and Health, Imperial College London, London, W2
275 1PG, UK.
- 276 121. Princess Al-Jawhara Al-Brahim Centre of Excellence in Research of Hereditary
277 Disorders (PACER-HD), King Abdulaziz University, Jeddah, 21589, Saudi Arabia.
- 278 122. Research Centre on Epidemiology and Preventive Medicine (EPIMED), Department of
279 Medicine and Surgery, University of Insubria, Varese, 2100, Italy.
- 280 123. INSERM UMR 1027, Toulouse, 31000, France.

281 124. Department of Cardiology, Toulouse University School of Medicine, Rangueil Hospital,
282 Toulouse, 31059, France.

283 125. Center for Genomic Medicine, Massachusetts General Hospital, Boston, MA, 02114,
284 USA.

285 126. Programs in Metabolism and Medical & Population Genetics, Broad Institute,
286 Cambridge, MA, 02142, USA.

287 127. Department of Medicine, Harvard Medical School, Boston, Massachusetts, 02115,
288 USA.

289 128. Center for Non-Communicable Diseases, Karachi, Pakistan.

290 129. Laboratory of Epidemiology and Population Sciences, National Institute on Aging,
291 National Institutes of Health, Bethesda, MD, USA.

292 130. Department of Epidemiology, Cardiovascular Health Research Unit, University of
293 Washington, Seattle, WA, 98195, USA.

294 131. Center for Genomic Medicine, Massachusetts General Hospital, USA.

295 132. Cardiovascular Research Center, Massachusetts General Hospital, Boston, MA, USA.

296 133. UKCRC Centre of Excellence for Public Health (NI), Queens University of Belfast,
297 Northern Ireland, BT7 1NN, UK.

298 134. Big Data Institute, Li Ka Shing Centre For Health Information and Discovery, University
299 of Oxford, Oxford, OX37BN, UK.

300 135. National Institute for Health and Welfare, Helsinki, 00271, Finland.

301 136. Institute of Human Genetics, Technische Universität München, Munich, 81675,
302 Germany.

303 137. Institute of Human Genetics, Helmholtz Zentrum München, German Research Center
304 for Environmental Health, Neuherberg, 85764, Germany.

305 138. Department of Genetics, University of North Carolina, Chapel Hill, North Carolina,
306 27599, USA.

307 139. Department of Epidemiology and Public Health, University of Strasbourg, Strasbourg,
308 F-67085, France.

309 140. Department of Public Health, University Hospital of Strasbourg, Strasbourg, F-67081,
310 France.

311 141. Clinical Research Centre, Centre for Molecular Medicine, Ninewells Hospital and
312 Medical School, Dundee, DD1 9SY, UK.

- 313 142. The Usher Institute to the Population Health Sciences and Informatics, University of
314 Edinburgh, Edinburgh, EH16 4UX, UK.
- 315 143. Aberdeen Biomedical Imaging Centre, School of Medicine Medical Sciences and
316 Nutrition, University of Aberdeen, Aberdeen, AB25 2ZD, UK.
- 317 144. Department of Clinical Sciences, Diabetes and Cardiovascular Disease, Genetic
318 Epidemiology, Lund University, Malmö, 20502, Sweden.
- 319 145. Finnish Institute for Molecular Medicine (FIMM), University of Helsinki, Helsinki,
320 Finland.
- 321 146. Department of Internal Medicine, Erasmus University Medical Center, Rotterdam,
322 3015CN, The Netherlands.
- 323 147. Department of Internal Medicine, Taichung Veterans General Hospital, Taichung
324 Taiwan, National Yang-Ming University, School of Medicine, Taipei, Taiwan, National
325 Defense Medical Center, School of Medicine, Taipei, Taiwan, Taichung, 40705, Taiwan.
- 326 148. McGill University and Génome Québec Innovation Centre, Montreal, Quebec, H3A
327 0G1, Canada.
- 328 149. Department of Human Genetics, McGill University, Montreal, Quebec, H3A 1B1,
329 Canada.
- 330 150. Division of Endocrinology and Metabolism, Department of Medicine, McGill
331 University, Montreal, Quebec, H3A 1A1, Canada.
- 332 151. Division of Population Health Sciences, Ninewells Hospital and Medical School,
333 University of Dundee, Dundee, DD1 9SY, UK.
- 334 152. Institute of Medical Informatics, Biometry and Epidemiology, Chair of Genetic
335 Epidemiology, Ludwig-Maximilians-Universität, Munich, 80802, Germany.
- 336 153. USC Roski Eye Institute, Department of Ophthalmology, Keck School of Medicine of
337 the University of Southern California, Los Angeles, California, 90033, US.
- 338 154. Department of Internal Medicine, Division of Cardiovascular Medicine, University of
339 Michigan, Ann Arbor, Michigan, 48109, USA.
- 340 155. Department of Computational Medicine and Bioinformatics, University of Michigan,
341 Ann Arbor, Michigan, 48109, USA.
- 342 156. Department of Human Genetics, University of Michigan, Ann Arbor, Michigan, 48109,
343 USA.

344 157. NIHR Blood and Transplant Research Unit in Donor Health and Genomics, Department
 345 of Public Health and Primary Care, University of Cambridge, Cambridge, CB1 8RN, UK.
 346 158. Imperial College Healthcare NHS Trust, Imperial College London, London, W12 0HS,
 347 UK.
 348 159. Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton,
 349 Cambridge, CB10 1RQ.
 350 160. British Heart Foundation, Cambridge Centre of Excellence, Department of Medicine,
 351 University of Cambridge, Cambridge, CB2 0QQ, UK.
 352 161. Department of Nutrition, Harvard School of Public Health, Boston, Massachusetts,
 353 02115, USA.
 354 162. Department of Genomics of Common Disease, School of Public Health, Imperial
 355 College London, London, W12 0NN, UK.
 356 163. Clinical Cooperation Group Type 2 Diabetes, Helmholtz Zentrum München, Ludwig-
 357 Maximilians University Munich, Germany.
 358 164. Clinical Cooperation Group Nutrigenomics and Type 2 Diabetes, Helmholtz Zentrum
 359 München, Technical University Munich, Germany.
 360 165. Faculty of Health Sciences, University of Southern Denmark, Odense, 5000, Denmark.
 361 166. University of Exeter Medical School, University of Exeter, Exeter, EX2 5DW, UK.
 362 167. Department of Medicine, Division of Cardiovascular Medicine, Stanford University
 363 School of Medicine, Stanford, CA, 94305, US.
 364 168. Department of Epidemiology, School of Public Health, University of Michigan, Ann
 365 Arbor, Michigan, 48109, USA.
 366 169. National Heart and Lung Institute, Cardiovascular Sciences, Hammersmith Campus,
 367 Imperial College London, London, W12 0NN, UK.
 368 170. Department of Medical Sciences, Uppsala University, Uppsala, SE-751 85, Sweden.
 369 171. Mindich Child Health and Development Institute, The Icahn School of Medicine at
 370 Mount Sinai, New York, NY, 10029, USA.
 371 172. Department of Statistics, University of Oxford, Oxford, OX1 3TG, UK.
 372 173. Department of Public Health and Primary Care, Leiden University Medical Center,
 373 Leiden, 2300 RC, The Netherlands.
 374 174. Pat Macpherson Centre for Pharmacogenetics and Pharmacogenomics, Ninewells
 375 Hospital and Medical School, University of Dundee, Dundee, DD1 9SY, UK.

175. Division of Epidemiology and Community Health, School of Public Health, University of Minnesota, Minneapolis, MN, 55454, US.
176. Cardiovascular Health Research Unit, Departments of Medicine, Epidemiology and Health Services, University of Washington, Seattle, WA, 98101-1448, USA.
177. Kaiser Permanent Washington Health Research Institute, Seattle, WA, 98101, USA.
178. Foundation for Research in Health, Exercise and Nutrition, Kuopio Research Institute of Exercise Medicine, Kuopio, Finland.
179. Institute of Cardiovascular and Medical Sciences, University of Glasgow, Glasgow, G12 8TA, UK.
180. Department of Hematology, School of Clinical Medicine, University of Cambridge, Cambridge, CB2 0AH.
181. Divisions of Endocrinology and Nephrology, University Hospital Leipzig, Liebigstr. 18, Leipzig, 04103, Germany.
182. Research Programs Unit, Diabetes and Obesity, University of Helsinki, Helsinki, Finland.
183. Dasman Diabetes Institute, Dasman, 15462, Kuwait.
184. Department of Neuroscience and Preventive Medicine, Danube-University Krems, Krems, 3500, Austria.
185. Diabetes Research Group, King Abdulaziz University, Jeddah, 21589, Saudi Arabia.
186. Department of Physiology and Biophysics, University of Mississippi Medical Center, Jackson, Mississippi, 39216, USA.
187. Metabolic Research Laboratories, Wellcome Trust – MRC Institute of Metabolic Science, University of Cambridge, Cambridge, CB22 0QQ, UK.
188. Division of Endocrinology, Diabetes and Metabolism, Cedars-Sinai Medical Center, Los Angeles, CA, 90048.
189. General Medicine Division, Massachusetts General Hospital and Department of Medicine, Harvard Medical School, Boston, Massachusetts, 02114, USA.
190. Departments of Pediatrics and Medicine, The Institute for Translational Genomics and Population Sciences, LABioMed at Harbor-UCLA Medical Center, Torrance, California, 90502, US.

*These authors jointly directed this work.

408 **Correspondence to:**

409

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

411 Jerome I Rotter (jrotter@labiomed.org)

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

413

We aggregated coding variant data for 81,412 type 2 diabetes cases and 370,832 controls of diverse ancestry, identifying 40 coding variant association signals ($p < 2.2 \times 10^{-7}$): of these, 16 map outside known risk loci. We make two important observations. First, only five of these signals are driven by low-frequency variants: even for these, effect sizes are modest (odds ratio ≤ 1.29). Second, when we used large-scale genome-wide association data to fine-map the associated variants in their regional context, accounting for the global enrichment of complex trait associations in coding sequence, compelling evidence for coding variant causality was obtained for only 16 signals. At 13 others, the associated coding variants clearly represent “false leads” with potential to generate erroneous mechanistic inference. Coding variant associations offer a direct route to biological insight for complex diseases and identification of validated therapeutic targets: however, appropriate mechanistic inference requires careful specification of their causal contribution to disease predisposition.

Genome-wide association studies (GWAS) have identified thousands of association signals influencing multifactorial traits such as type 2 diabetes (T2D) and obesity¹⁻⁷. Most of these associations involve common variants that map to non-coding sequence, and identification of their cognate effector transcripts has proved challenging. Identification of coding variants causally implicated in trait predisposition offers a more direct route from association signal to biological inference.

The exome occupies 1.5% of overall genome sequence, but for many common diseases, coding variants make a disproportionate contribution to trait heritability^{8,9}. This enrichment indicates that coding variant association signals have an enhanced probability of being causal when compared to those involving an otherwise equivalent non-coding variant. This does not, however, guarantee that all coding variant associations are causal. Alleles driving common-variant (minor allele frequency [MAF] $\geq 5\%$) GWAS signals typically reside on extended risk haplotypes that, owing to linkage disequilibrium (LD), incorporate many common variants^{10,11}. Consequently, the presence of a coding allele on the risk haplotype does not constitute sufficient evidence that it represents the causal variant at the locus, or that the gene within which it lies is mediating the association signal. Since much coding variant discovery has proceeded through exome-specific analyses (either exome-array genotyping or exome sequencing), researchers have often been poorly-placed to position

coding variant associations in the context of regional genetic variation. It is unclear how often this may have led to incorrect assumptions regarding their causal role.

In our recent study of T2D predisposition¹², we surveyed the exomes of 34,809 T2D cases and 57,985 controls, of predominantly European descent, and identified 13 distinct coding variant associations reaching genome-wide significance. Twelve of these associations involved common variants, but the data hinted at a substantial pool of lower-frequency coding variants of moderate impact, potentially amenable to detection in larger samples. We also reported that, whilst many of these signals fell within common variant loci previously identified by GWAS, it was far from trivial to determine, using available data, whether those coding variants were causal or ‘hitchhiking’ on risk haplotypes.

Here, we report analyses that address these two issues. First, we extend the scope of our exome-array genotyping to include data from 81,412 T2D cases and 370,832 controls of diverse ancestry, substantially increasing power to detect coding variant associations across the allele-frequency spectrum. Second, to understand the extent to which identification of coding variant associations provides a reliable guide to causal mechanisms, we undertake high-resolution fine-mapping of identified coding variant association signals in 50,160 T2D cases and 465,272 controls of European ancestry with genome-wide genotyping data.

RESULTS

Discovery study overview. First, we set out to discover coding variant association signals by aggregating T2D association summary statistics in up to 452,244 individuals (effective sample size 228,825) across five ancestry groups, performing both European-specific (EUR) and trans-ethnic (TE) meta-analyses (**Supplementary Tables 1 and 2**). Analysis was restricted to the 247,470 variants represented on the exome-array. Genotypes were assembled from: (a) 58,425 cases and 188,032 controls genotyped with the exome-array; (b) 14,608 cases and 174,322 controls from UK Biobank and GERA (Resource for Genetic Epidemiology on Adult Health and Aging) genotyped with GWAS arrays enriched for exome content and/or coverage of low-frequency variation across ethnic groups^{13,14}; and (c) 8,379 cases and 8,478 controls with whole-exome sequence from GoT2D/T2D-GENES¹² and SIGMA¹⁵ studies. Overall, this represented a 3-fold increase in effective sample size over our previous study of T2D predisposition within coding sequence¹². To deconvolute the impact

of obesity on T2D-associated variants, association analyses were conducted with and without body mass index (BMI) adjustment.

We considered $p < 2.2 \times 10^{-7}$ as significant for protein truncating variants (PTVs) and moderate impact coding variants (including missense, in-frame indel and splice region variants) based on a weighted Bonferroni correction that accounts for the observed enrichment in complex trait association signals across sequence annotation¹⁶. This threshold matches those obtained through other approaches such as simple Bonferroni correction for the number of coding variants on the exome-array (**Methods**). Compared to our previous study¹², the expanded sample size substantially increased power to detect association for common variants of modest effect (e.g. from 14.4% to 97.9% for a variant with 20% MAF and odds ratio [OR]=1.05) and lower-frequency variants with larger effects (e.g. from 11.8% to 97.5% for a variant with 1% MAF and OR=1.20) assuming homogenous allelic effects across ancestry groups (**Methods**).

Insights into coding variant association signals underlying T2D susceptibility. We detected significant associations at 69 coding variants under an additive genetic model (either in BMI unadjusted or adjusted analysis), mapping to 38 loci (**Supplementary Fig. 1, Supplementary Table 3**). We observed minimal evidence of heterogeneity in allelic OR between ancestry groups (**Supplementary Table 3**), and no compelling evidence for non-additive allelic effects (**Supplementary Fig. 2, Supplementary Table 4**). Reciprocal conditional analyses (**Methods**) indicated that the 69 coding variants represented 40 distinct association signals (conditional $p < 2.2 \times 10^{-7}$) across the 38 loci, with two distinct signals each at *HNF1A* and *RREB1* (**Supplementary Table 5**). These 40 signals included the 13 associations reported in our earlier publication¹², each featuring more significant associations in this expanded meta-analysis (**Supplementary Table 6**). Twenty-five of the 40 signals were significant in both EUR and TE analyses. Of the other 15, three (*PLCB3*, *C17orf58*, and *ZHX3*) were significant in EUR, and all reached $p_{TE} < 6.8 \times 10^{-6}$ in the TE analysis: for *PLCB3* and *ZHX3*, risk allele frequencies were substantially lower outside European descent populations. Twelve loci (**Supplementary Table 3**) were significant in TE alone, but for these (except *PAX4* which is East Asian specific), the evidence for association was proportionate in the smaller EUR component ($p_{EUR} < 8.4 \times 10^{-5}$).

Sixteen of the 40 distinct association signals mapped outside regions previously implicated in T2D susceptibility (Methods, **Table 1**). These included missense variant signals in *POC5* (p.His36Arg, rs2307111, $p_{TE}=1.6\times10^{-15}$), *PNPLA3* (p.Ile148Met, rs738409, p_{TE} BMI-adjusted= 2.8×10^{-11}), and *ZZEF1* (p.Ile2014Val, rs781831, $p_{TE}=8.3\times10^{-11}$).

In addition to the 69 coding variant signals, we detected significant ($p<5\times10^{-8}$) and novel T2D-associations for 20 non-coding variants (at 15 loci) that were also assayed on the exome-array (**Supplementary Table 7**). Three of these (*POC5*, *LPL*, and *BPTF*) overlap with novel coding signals reported here.

Contribution of low-frequency and rare coding variation to T2D susceptibility. Despite increased power and good coverage of low-frequency variants on the exome-array¹², 35 of the 40 distinct coding variant association signals were common, with modest effects (allelic ORs 1.02-1.36) (**Supplementary Fig. 3, Supplementary Table 3**). The five signals attributable to lower-frequency variants were also of modest effect (allelic ORs 1.09-1.29) (**Supplementary Fig. 3**). Two of the lower-frequency variant signals were novel, and in both, the minor allele was protective against T2D: *FAM63A* p.Tyr95Asn (rs140386498, MAF=1.2%, OR= 0.82 [0.77-0.88], $p_{EUR}=5.8\times10^{-8}$) and *ANKH* p.Arg187Gln (rs146886108, MAF=0.4%, OR=0.78 [0.69-0.87], $p_{EUR}=2.0\times10^{-7}$). Both variants were very rare or monomorphic in non-European descent individuals.

In Fuchsberger et al.¹², we highlighted a set of 100 low-frequency coding variants with allelic ORs between 1.10 and 2.66, which despite relatively large estimates for liability-scale variance explained, had not reached significance. In this expanded analysis, only five of these variants, including the two novel associations at *FAM63A* p.Tyr95Asn and *ANKH* p.Arg187Gln, attained significance. More precise effect-size estimation in the larger sample size indicates that OR estimates in the earlier study were subject to a substantial upwards bias (**Supplementary Fig. 3**).

To detect additional rare variant association signals, we performed gene-based analyses (burden and SKAT¹⁷) using previously-defined “strict” and “broad” masks, filtered for annotation and MAF^{12,18} (**Methods**). We identified gene-based associations with T2D susceptibility ($p<2.5\times10^{-6}$, Bonferroni correction for 20,000 genes) for *FAM63A* (10 variants, combined MAF=1.9%, $p_{EUR}=3.1\times10^{-9}$) and *PAM* (17 variants, combined MAF=4.7%, $p_{TE}=8.2\times10^{-9}$). On conditional analysis (**Supplementary Table 8**), the gene-based signal at

FAM63A was entirely attributable to the low-frequency p.Tyr95Asn allele described earlier (conditional $p=0.26_{EUR}$). The gene-based signal for *PAM* was also driven by a single low-frequency variant (p.Asp563Gly; conditional $p_{TE}=0.15$). A second, previously-described, low-frequency variant, *PAM* p.Ser539Trp¹⁹, is not represented on the exome-array, and did not contribute to these analyses.

Fine-mapping of coding variant association signals with T2D susceptibility. These analyses identified 40 distinct coding variant associations with T2D, but this information is not sufficient to determine that these variants are causal for disease. To assess the role of these coding variants given regional genetic variation, we fine-mapped these association signals using a meta-analysis of 50,160 T2D cases and 465,272 controls (European-descent only; partially overlapping with the discovery samples), which we aggregated from 24 GWAS. Each component GWAS was imputed using appropriate high-density reference panels (for most, the Haplotype Reference Consortium²⁰; **Methods, Supplementary Table 9**). Before fine-mapping, distinct association signals were delineated using approximate conditional analyses (**Methods, Supplementary Table 5**). We included 37 of the 40 identified coding variants in this fine-mapping analysis, excluding three (those at the *MHC*, *PAX4*, and *ZHX3*) that were, for various reasons (see **Methods**), not amenable to fine-mapping in the GWAS data.

For each of these 37 signals, we first constructed “functionally-unweighted” credible variant sets, which collectively account for 99% of the posterior probability of association (PPA), based exclusively on the meta-analysis summary statistics²¹ (**Methods, Supplementary Table 10**). For each signal, we calculated the proportion of PPA attributable to coding variants (missense, in-frame indel, and splice region variants; **Figure 1, Supplementary Fig. 4 and 5**). There were only two signals at which coding variants accounted for $\geq 80\%$ of PPA: *HNF4A* p.Thr139Ile (rs1800961, PPA>0.999) and *RREB1* p.Asp1171Asn (rs9379084, PPA=0.920). However, at other signals, including those for *GCKR* p.Pro446Leu and *SLC30A8* p.Arg276Trp, for which robust empirical evidence has established a causal role^{22,23}, genetic support for coding variant causation was weak. This is because coding variants were typically in high LD ($r^2>0.9$) with large numbers of non-coding variants, such that the PPA was distributed across many sites with broadly equivalent evidence for association.

These functionally-unweighted sets are based on genetic fine-mapping data alone, and do not account for the disproportionate representation of coding variants amongst GWAS associations for complex traits^{8,9}. To accommodate this information, we extended the fine-mapping analyses by incorporating an “annotation-informed prior” model of causality. We derived priors from estimates of the enrichment of association signals by sequence annotation from analyses conducted by deCODE across 96 quantitative and 123 binary phenotypes¹⁶ (**Methods**). This model “boosts” the prior, and hence the posterior probabilities (we use ‘ aiPPA ’ to denote annotation-informed PPAs) of coding variants. It also takes account (in a tissue-non-specific manner) of the GWAS enrichment of variants within enhancer elements (as assayed through DNase I hypersensitivity) when compared to non-coding variants mapping elsewhere. The annotation-informed model generated smaller 99% credible sets across most signals, corresponding to fine-mapping at higher resolution (**Supplementary Table 10**). As expected, the contribution of coding variants was increased under the annotation-informed model. At these 37 association signals, we distinguished three broad patterns of causal relationships between coding variants and T2D risk.

Group 1: T2D association signal is driven by coding variants. At 16 of the 37 distinct signals, coding variation accounted for >80% of the aiPPA (**Fig. 1, Table 2, Supplementary Table 10**). This was attributable to a single coding variant at 12 signals and multiple coding variants at four. Reassuringly, group 1 signals confirmed coding variant causation for several loci (*GCKR*, *PAM*, *SLC30A8*, *KCNJ11-ABCC8*) at which functional studies have established strong mechanistic links to T2D pathogenesis (**Table 2**). T2D association signals at the 12 remaining signals (**Fig. 1, Supplementary Table 10**) had not previously been shown to be driven by coding variation, but our fine-mapping analyses pointed to causal coding variants with high aiPPA values: these included *HNF4A*, *RREB1* (p. Asp1171Asn), *ANKH*, *WSCD2*, *POC5*, *TM6SF2*, *HNF1A* (p. Ala146Val; p. Ile75Leu), *GIPR*, *LPL*, *PLCB3*, and *PNPLA3* (**Table 2**). At several of these, independent evidence corroborates the causal role of the genes harbouring the associated coding variants. For example, rare coding mutations at *HNF1A* and *HNF4A* are causal for monogenic, early-onset forms of diabetes²⁴; and at *TM6SF2* and *PNPLA3*, the associated coding variants are implicated in the development of non-alcoholic fatty liver disease (NAFLD)^{25,26}.

The use of priors to capture the enrichment of coding variants seems a reasonable model, genome-wide. However, at any given locus, strong priors (especially for PTVs) might elevate to apparent causality, variants that would have been excluded from a causal role on the basis of genetic fine-mapping alone. Comparison of the annotation-informed and functionally-unweighted credible sets for group 1 signals indicated that this scenario was unlikely. For 11 of the 16 (*GCKR*, *PAM*, *KCNJ11-ABCC8*, *HNF4A*, *RREB1* [p.Asp1171Asn], *ANKH*, *POC5*, *TM6SF2*, *HNF1A* [p.Ala146Val], *PLCB3*, *PNPLA3*), the coding variant had the highest PPA in the fine-mapping analysis (**Table 2**) even under the functionally-unweighted model. At *SLC30A8*, *WSCD2*, and *GIPR*, the coding variants had similar PPAs to the lead non-coding SNPs under the functionally-unweighted prior (**Table 2**). At these 14 signals therefore, coding variants have either greater or equivalent PPA to the best flanking non-coding SNPs under the functionally-unweighted model, but receive a boost in PPA after incorporating the annotation weights.

The situation is less clear at *LPL*. Here, fine-mapping resolution is poor under the functionally-unweighted prior, and the coding variant sits on an extended haplotype in strong LD with non-coding variants, some with higher PPA, such as rs74855321 (PPA=0.048) (compared to *LPL* p.Ser474* [rs328, PPA=0.023]). However, *LPL* p.Ser474* is annotated as a PTV, and benefits from a substantially-increased prior that boosts its annotation-informed ranking (**Table 2**). Ultimately, decisions regarding the causal role of any such variant must rest on the amalgamation of evidence from diverse sources including detailed functional evaluation of the coding variants, and of other variants with which they are in LD.

Group 2: T2D association signals are not attributable to coding variants. At 13 of the 37 distinct signals, coding variation accounted for <20% of the PPA, even after applying the annotation-informed prior model. These signals are likely to be driven by local non-coding variation and mediated through regulatory mechanisms. Five of these signals (*TPCN2*, *MLX*, *ZZEF1*, *C17orf58*, and *CEP68*) represent novel T2D-association signals identified in the exome-focused analysis. Given the exome-array discoveries, it would have been natural to consider the named genes at these, and other loci in this group, as candidates for mediation of their respective association signals. However, the fine-mapping analyses indicate that these coding variants do not provide useful mechanistic inference given low a_i PPA (**Fig. 1**, **Table 2**).

The coding variant association at the *CENTD2* (*ARAP1*) locus is a case-in-point. The association with the p.Gln802Glu variant in *ARAP1* (rs56200889, $p_{TE}=4.8 \times 10^{-8}$ but $aiPPA < 0.001$) is seen in the fine-mapping analysis to be secondary to a substantially stronger non-coding association signal involving a cluster of variants including rs11603334 ($p_{TE}=9.5 \times 10^{-18}$, $aiPPA=0.0692$) and rs1552224 ($p_{TE}=2.5 \times 10^{-17}$, $aiPPA=0.0941$). The identity of the effector transcript at this locus has been the subject of detailed investigation, and some early studies used islet expression data to promote *ARAP1*²⁷. However, a more recent study integrating human islet genomics and murine gene knockout data establishes *STARD10* as the gene mediating the GWAS signal, consistent with the reassignment of the *ARAP1* coding variant association as irrelevant to causal inference²⁸.

Whilst, at these loci, the coding variant associations represent “false leads”, this does not necessarily exclude the genes concerned from a causal role. At *WFS1* for example, coding variants too rare to be visible to the array-based analyses we performed, and statistically independent of the common p.Val333Ile variant we detected, cause an early-onset form of diabetes that renders *WFS1* the strongest local candidate for T2D predisposition.

Group 3: Fine-mapping data consistent with partial role for coding variants. At eight of the 37 distinct signals, the $aiPPA$ attributable to coding variation lay between 20% and 80%. At these signals, the evidence is consistent with “partial” contributions from coding variants, although the precise inference is likely to be locus-specific, dependent on subtle variations in LD, imputation accuracy, and the extent to which global priors accurately represent the functional impact of the specific variants concerned.

This group includes *PPARG* for which independent evidence corroborates the causal role of this specific effector transcript with respect to T2D-risk. *PPARG* encodes the target of antidiabetic thiazolidinedione drugs and harbours very rare coding variants causal for lipodystrophy and insulin resistance, conditions highly-relevant to T2D. The common variant association signal at this locus has generally been attributed to the p.Pro12Ala coding variant (rs1801282) although empirical evidence that this variant influences *PPARG* function is scant²⁹⁻³¹. In the functionally-unweighted analysis, p.Pro12Ala had an unimpressive PPA (0.0238); after including annotation-informed priors, the same variant emerged with the highest $aiPPA$ (0.410), although the 99% credible set included 19 non-coding variants,

spanning 67kb (**Supplementary Table 10**). These credible set variants included rs4684847 ($\text{aiPPA}=0.0089$), at which the T2D-associated allele has been reported to impact *PPARG2* expression and insulin sensitivity by altering binding of the homeobox transcription factor PRRX1³². These data are consistent with a model whereby regulatory variants contribute to altered PPARG activity in combination with, or potentially to the exclusion of, p.Pro12Ala. Future improvements in functional annotation for regulatory variants (gathered from relevant tissues and cell types) should provide increasingly granular priors that allow fine-tuned assignment of causality at loci such as this.

Functional impact of coding alleles. In other contexts, the functional impact of coding alleles is correlated with: (i) variant-specific features, including measures of conservation and predicted impact on protein structure; and (ii) gene-specific features such as extreme selective constraints as quantified by the intolerance to functional variation³³. To determine whether similar measures could capture information pertinent to T2D causation, we compared coding variants falling into the different fine-mapping groups for a variety of measures including MAF, Combined Annotation Dependent Depletion (CADD) score³⁴, and loss-of-function (LoF)-intolerance metric, pLI³³ (**Methods, Fig. 2**). Variants from group 1 had significantly higher CADD-scores than those in group 2 (Kolmogorov-Smirnov $p=0.0031$). Except for the variants at *KCNJ11-ABCC8* and *GCKR*, all group 1 coding variants considered likely to be driving T2D association signals had CADD-score ≥ 20 . On this basis, we predict that the East-Asian specific coding variant at *PAX4*, for which the fine-mapping data were not informative, is also likely causal for T2D.

T2D loci and physiological classification. The development of T2D involves dysfunction of multiple mechanisms. Systematic analysis of the physiological effects of known T2D-risk alleles has improved understanding of the mechanisms through which they exert their primary impact on disease risk³⁵. We obtained association summary statistics for diverse metabolic traits (and other outcomes) for 94 T2D-associated index variants. These 94 were restricted to sites represented on the exome-array and included the 40 coding signals plus 54 distinct non-coding signals (12 novel and 42 previously-reported non-coding GWAS lead SNPs). We applied clustering techniques (**Methods**) to generate multi-trait association patterns, allocating 71 of the 94 loci to one of three main physiological categories

(**Supplementary Figs. 6, Supplementary Table 11**). The first category, comprising nine T2D-risk loci with strong BMI and dyslipidemia associations, included three of the novel coding signals: *PNPLA3*, *POC5* and *BPTF*. The T2D associations at both *POC5* and *BPTF* were substantially attenuated (>2 -fold decrease in $-\log_{10}p$) after adjusting for BMI (**Table 1, Supplementary Table 3, Supplementary Fig. 7**), indicating that their impact on T2D-risk is likely mediated by a primary effect on adiposity. *PNPLA3* and *POC5* are established NAFLD²⁵ and BMI⁶ loci, respectively. The second category featured 39 loci at which multi-trait profiles indicated a primary effect on insulin secretion. This set included four of the novel coding variant signals (*ANKH*, *ZZEF1*, *TTL6*, *ZHX3*). The third category encompassed 23 loci with primary effects on insulin action, including signals at the *KIF9*, *PLCB3*, *CEP68*, *TPCN2*, *FAM63A*, and *PIM3* loci. For most variants in this category, the T2D-risk allele was associated with lower BMI, and T2D association signals were more pronounced after adjustment for BMI. At a subset of these loci, including *KIF9* and *PLCB3*, T2D-risk alleles were associated with higher waist-hip ratio and lower body fat percentage, indicating that the mechanism of action likely reflects limitations in storage capacity of peripheral adipose tissue³⁶.

DISCUSSION

The present study adds to mounting evidence constraining the contribution of lower-frequency variants to T2D-risk. Although the exome-array interrogates only a subset of the universe of coding variants, it captures the majority of low-frequency coding variants in European populations. The substantial increase in sample size in the present study over our previous effort¹² (effective sample sizes of 228,825 and 82,758, respectively), provides more robust evaluation of the effect size distribution in this low-frequency variant range, and indicates that previous analyses are likely, if anything, to have overestimated the contribution of low-frequency variants to T2D-risk.

The present study is less informative regarding rare variants. These are sparsely captured on the exome-array. In addition, the combination of greater regional diversity in rare allele distribution and the enormous sample sizes required to detect rare variant associations (likely to require meta-analysis of data from diverse populations) acts against their identification. Our complementary genome and exome sequence analyses have thus far failed to register strong evidence for a substantial rare variant component to T2D-risk¹².

It is therefore highly unlikely that rare variants missed in our analyses are causal for any of the common or low-frequency variant associations we have detected and fine-mapped. On the other hand, it is probable that rare coding alleles, with associations that are distinct from the common variant signals we have examined and detected only through sequence based analyses, will provide additional clues to the most likely effector transcripts at some of these signals (*WFS1* provides one such example).

Once a coding variant association is detected, it is natural to assume a causal connection between that variant, the gene in which it sits, and the phenotype of interest. Whilst such assignments may be robust for many rare protein-truncating alleles, we demonstrate that this implicit assumption is often inaccurate, particularly for associations attributable to common, missense variants. A third of the coding variant associations we detected were, when assessed in the context of regional LD, highly unlikely to be causal. At these loci, the genes within which they reside are consequently deprived of their implied connection to disease risk, and attention redirected towards nearby non-coding variants and their impact on regional gene expression. As a group, coding variants we assign as causal are predicted to have a more deleterious impact on gene function than those that we exonerate, but, as in other settings, coding annotation methods lack both sensitivity and specificity. It is worth emphasising that empirical evidence that the associated coding allele is “functional” (i.e. can be shown to influence cognate gene function in some experimental assay) provides limited reassurance that the coding variant is responsible for the T2D association, unless that specific perturbation of gene function can itself be plausibly linked to the disease phenotype.

Our fine-mapping analyses make use of the observation that coding variants are globally enriched across GWAS signals^{8,9,16} with greater prior probability of causality assigned to those with more severe impact on biological function. We assigned diminished priors to non-coding variants, with lowest support for those mapping outside of DNase I hypersensitive sites. The extent to which our findings corroborate previous assignments of causality (often substantiated by detailed, disease-appropriate functional assessment and other orthogonal evidence) suggests that even these sparse annotations provide valuable information to guide target validation. Nevertheless, there are inevitable limits to the extrapolation of these ‘broad-brush’ genome-wide enrichments to individual loci: improvements in functional annotation for both coding and regulatory variants, particularly

when gathered from trait-relevant tissues and cell types, should provide more granular, trait-specific priors to fine-tune assignment of causality within associated regions. These will motivate target validation efforts that benefit from synthesis of both coding and regulatory mechanisms of gene perturbation. It also needs to be acknowledged that, without whole genome sequencing data on sample sizes comparable to those we have examined here, imperfections arising from the imputation may confound fine-mapping precision at some loci, and that robust inference will inevitably depend on integration of diverse sources of genetic, genomic and functional data.

The term “smoking gun” has often been used to describe the potential of functional coding variants to provide causal inference with respect to pathogenetic mechanisms³⁷. This study provides a timely reminder that, even when a suspect with a smoking gun is found at the scene of a crime, it should not be assumed that they fired the fatal bullet.

ACKNOWLEDGMENTS

A full list of acknowledgments appears in the **Supplementary Information**. Part of this work was conducted using the UK Biobank Resource under Application Number 9161.

AUTHOR CONTRIBUTIONS

Project co-ordination. A.Mahajan, A.P.M., J.I.R., M.I.M.

Core analyses and writing. A.Mahajan, J.W., S.M.W, W.Zhao, N.R.R., A.Y.C., W.G., H.K., R.A.S., I.Barroso, T.M.F., M.O.G., J.B.M., M.Boehnke, D.S., A.P.M., J.I.R., M.I.M.

Statistical Analysis in individual studies. A.Mahajan, J.W., S.M.W., W.Zhao, N.R.R., A.Y.C., W.G., H.K., D.T., N.W.R., X.G., Y.Lu, M.Li, R.A.J., Y.Hu, S.Huo, K.K.L., W.Zhang, J.P.C., B.P., J.Flannick, N.G., V.V.T., J.Kravic, Y.J.K., D.V.R., H.Y., M.M.-N., K.M., R.L.-G., T.V.V., J.Marten, J.Li, A.V.S., P.An, S.L., S.G., G.M., A.Demirkan, J.F.T., V.Steinhorsdottir, M.W., C.Lecoeur, M.Preuss, L.F.B., P.Almgren, J.B.-J., J.A.B., M.Canouil, K.-U.E., H.G.d.H., Y.Hai, S.Han, S.J., F.Kronenberg, K.L., L.A.L., J.-J.L., H.L., C.-T.L., J.Liu, R.M., K.R., S.S., P.S., T.M.T., G.T., A.Tin, A.R.W., P.Y., J.Y., L.Y., R.Y., J.C.C., D.I.C., C.v.D., J.Dupuis, P.W.F., A.Köttgen, D.M.-K., N.Soranzo, R.A.S., A.P.M.

Genotyping. A.Mahajan, N.R.R., A.Y.C., Y.Lu, Y.Hu, S.Huo, B.P., N.G., R.L.-G., P.An, G.M., E.A., N.A., C.B., N.P.B., Y.-D.I.C., Y.S.C., M.L.G., H.G.d.H., S.Hackinger, S.J., B.-J.K., P.K., J.Kriebel, F.Kronenberg, H.L., S.S.R., K.D.T., E.B., E.P.B., P.D., J.C.F., S.R.H., C.Langenberg, M.A.P., F.R.,

796 A.G.U., J.C.C., D.I.C., P.W.F., B.-G.H., C.H., E.I., S.L.K., J.S.K., Y.Liu, R.J.F.L., N.Soranzo, N.J.W.,
797 R.A.S., T.M.F., A.P.M., J.I.R., M.I.M.

798 **Cross-trait lookups in unpublished data.** S.M.W., A.Y.C., Y.Lu, M.Li, M.G., H.M.H., A.E.J.,
799 D.J.L., E.M., G.M.P., H.R.W., S.K., C.J.W.

800 **Phenotyping.** Y.Lu, Y.Hu, S.Huo, P.An, S.L., A.Demirkan, S.Afaq, S.Afzal, L.B.B., A.G.B.,
801 I.Brandslund, C.C., S.V.E., G.G., V.Giedraitis, A.T.-H., M.-F.H., B.I., M.E.J., T.J., A.Käräjämäki,
802 S.S.K., H.A.K., P.K., F.Kronenberg, B.L., H.L., K.-H.L., A.L., J.Liu, M.Loh, V.M., R.M.-C., G.N.,
803 M.N., S.F.N., I.N., P.A.P., W.R., L.R., O.R., S.S., E.S., K.S.S., A.S., B.T., A.Tönjes, A.V., D.R.W.,
804 H.B., E.P.B., A.Deaghan, J.C.F., S.R.H., C.Langenberg, A.D.Morris, R.d.M., M.A.P., A.R., P.M.R.,
805 F.R.R., V.Salomaa, W.H.-H.S., R.V., J.C.C., J.Dupuis, O.H.F., H.G., B.-G.H., T.H., A.T.H., C.H.,
806 S.L.K., J.S.K., A.Köttgen, L.L., Y.Liu, R.J.F.L., C.N.A.P., J.S.P., O.P., B.M.P., M.B.S., N.J.W.,
807 T.M.F., M.O.G.

808 **Individual study design and principal investigators.** N.G., P.An, B.-J.K., P.Amouyel, H.B., E.B.,
809 E.P.B., R.C., F.S.C., G.D., A.Deaghan, P.D., M.M.F., J.Ferrières, J.C.F., P.Frossard, V.Gudnason,
810 T.B.H., S.R.H., J.M.M.H., M.I., F.Kee, J.Kuusisto, C.Langenberg, L.J.L., C.M.L., S.M., T.M., O.M.,
811 K.L.M., M.M., A.D.Morris, A.D.Murray, R.d.M., M.O.-M., K.R.O., M.Perola, A.P., M.A.P.,
812 P.M.R., F.R., F.R.R., A.H.R., V.Salomaa, W.H.-H.S., R.S., B.H.S., K.Strauch, A.G.U., R.V.,
813 M.Blüher, A.S.B., J.C.C., D.I.C., J.Danesh, C.v.D., O.H.F., P.W.F., P.Froguel, H.G., L.G., T.H.,
814 A.T.H., C.H., E.I., S.L.K., F.Karpe, J.S.K., A.Köttgen, K.K., M.Laakso, X.L., L.L., Y.Liu, R.J.F.L.,
815 J.Marchini, A.Metspalu, D.M.-K., B.G.N., C.N.A.P., J.S.P., O.P., B.M.P., R.R., N.Sattar, M.B.S.,
816 N.Soranzo, T.D.S., K.Stefansson, M.S., U.T., T.T., J.T., N.J.W., J.G.W., E.Z., I.Barroso, T.M.F.,
817 J.B.M., M.Boehnke, D.S., A.P.M., J.I.R., M.I.M.

818

819 **DISCLOSURES**

820 Jose C Florez has received consulting honoraria from Merck and from Boehringer-Ingelheim.
821 Daniel I Chasman received funding for exome chip genotyping in the WGHS from Amgen.
822 Oscar H Franco works in ErasmusAGE, a center for aging research across the life course
823 funded by Nestlé Nutrition (Nestec Ltd.), Metagenics Inc., and AXA. Nestlé Nutrition (Nestec
824 Ltd.), Metagenics Inc., and AXA had no role in the design and conduct of the study;
825 collection, management, analysis, and interpretation of the data; and preparation, review or
826 approval of the manuscript. Erik Ingelsson is an advisor and consultant for Precision
827 Wellness, Inc., and advisor for Cellink for work unrelated to the present project. Bruce M

828 Psaty serves on the DSMB for a clinical trial funded by the manufacturer (Zoll LifeCor) and
829 on the Steering Committee of the Yale Open Data Access Project funded by Johnson &
830 Johnson. Inês Barroso and spouse own stock in GlaxoSmithKline and Incyte Corporation.
831 Timothy Frayling has consulted for Boeringer IngelHeim and Sanofi on the genetics of
832 diabetes. Danish Saleheen has received support from Pfizer, Regeneron, Genentech and Eli
833 Lilly. Mark I McCarthy has served on advisory panels for NovoNordisk and Pfizer, and
834 received honoraria from NovoNordisk, Pfizer, Sanofi-Aventis and Eli Lilly.

REFERENCES

1. 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).
2. 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).
3. 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).
4. Mahajan, A. *et al.* Genome-wide trans-ancestry meta-analysis provides insight into the genetic architecture of type 2 diabetes susceptibility. *Nat Genet* **46**, 234-44 (2014).
5. Ng, M.C. *et al.* Meta-analysis of genome-wide association studies in African Americans provides insights into the genetic architecture of type 2 diabetes. *PLoS Genet* **10**, e1004517 (2014).
6. Locke, A.E. *et al.* Genetic studies of body mass index yield new insights for obesity biology. *Nature* **518**, 197-206 (2015).
7. Shungin, D. *et al.* New genetic loci link adipose and insulin biology to body fat distribution. *Nature* **518**, 187-96 (2015).
8. Gusev, A. *et al.* Partitioning heritability of regulatory and cell-type-specific variants across 11 common diseases. *Am J Hum Genet* **95**, 535-52 (2014).
9. Walter, K. *et al.* The UK10K project identifies rare variants in health and disease. *Nature* **526**, 82-90 (2015).
10. 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).
11. Horikoshi, M. *et al.* Transancestral fine-mapping of four type 2 diabetes susceptibility loci highlights potential causal regulatory mechanisms. *Hum Mol Genet* **25**, 2070-2081 (2016).
12. Fuchsberger, C. *et al.* The genetic architecture of type 2 diabetes. *Nature* **536**, 41-7 (2016).

- 865 13. Sudlow, C. *et al.* UK biobank: an open access resource for identifying the causes of a
866 wide range of complex diseases of middle and old age. *PLoS Med* **12**, e1001779
867 (2015).
- 868 14. Cook, J.P. & Morris, A.P. Multi-ethnic genome-wide association study identifies novel
869 locus for type 2 diabetes susceptibility. *Eur J Hum Genet* **24**, 1175-80 (2016).
- 870 15. Estrada, K. *et al.* Association of a low-frequency variant in HNF1A with type 2
871 diabetes in a Latino population. *JAMA* **311**, 2305-14 (2014).
- 872 16. Sveinbjornsson, G. *et al.* Weighting sequence variants based on their annotation
873 increases power of whole-genome association studies. *Nat Genet* **48**, 314-7 (2016).
- 874 17. Liu, D.J. *et al.* Meta-analysis of gene-level tests for rare variant association. *Nat*
875 *Genet* **46**, 200-4 (2014).
- 876 18. Purcell, S.M. *et al.* A polygenic burden of rare disruptive mutations in schizophrenia.
877 *Nature* **506**, 185-90 (2014).
- 878 19. Steinthorsdottir, V. *et al.* Identification of low-frequency and rare sequence variants
879 associated with elevated or reduced risk of type 2 diabetes. *Nat Genet* **46**, 294-8
880 (2014).
- 881 20. McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation.
882 *Nat Genet* **48**, 1279-83 (2016).
- 883 21. Maller, J.B. *et al.* Bayesian refinement of association signals for 14 loci in 3 common
884 diseases. *Nat Genet* **44**, 1294-301 (2012).
- 885 22. Flannick, J. *et al.* Loss-of-function mutations in SLC30A8 protect against type 2
886 diabetes. *Nat Genet* **46**, 357-63 (2014).
- 887 23. Beer, N.L. *et al.* The P446L variant in GCKR associated with fasting plasma glucose
888 and triglyceride levels exerts its effect through increased glucokinase activity in liver.
889 *Hum Mol Genet* **18**, 4081-8 (2009).
- 890 24. Murphy, R., Ellard, S. & Hattersley, A.T. Clinical implications of a molecular genetic
891 classification of monogenic beta-cell diabetes. *Nat Clin Pract Endocrinol Metab* **4**,
892 200-13 (2008).
- 893 25. Romeo, S. *et al.* Genetic variation in PNPLA3 confers susceptibility to nonalcoholic
894 fatty liver disease. *Nat Genet* **40**, 1461-5 (2008).
- 895 26. Kozlitina, J. *et al.* Exome-wide association study identifies a TM6SF2 variant that
896 confers susceptibility to nonalcoholic fatty liver disease. *Nat Genet* **46**, 352-6 (2014).

27. Kulzer, J.R. *et al.* A common functional regulatory variant at a type 2 diabetes locus upregulates ARAP1 expression in the pancreatic beta cell. *Am J Hum Genet* **94**, 186-97 (2014).
28. Carrat, G.R. *et al.* Decreased STARD10 expression is associated with defective insulin secretion in humans and mice. *Am J Hum Genet* **100**, 238-256 (2017).
29. Deeb, S.S. *et al.* A Pro12Ala substitution in PPARgamma2 associated with decreased receptor activity, lower body mass index and improved insulin sensitivity. *Nat Genet* **20**, 284-7 (1998).
30. Majithia, A.R. *et al.* Rare variants in PPARG with decreased activity in adipocyte differentiation are associated with increased risk of type 2 diabetes. *Proc Natl Acad Sci U S A* **111**, 13127-32 (2014).
31. Majithia, A.R. *et al.* Prospective functional classification of all possible missense variants in PPARG. *Nat Genet* **48**, 1570-1575 (2016).
32. Claussnitzer, M. *et al.* Leveraging cross-species transcription factor binding site patterns: from diabetes risk loci to disease mechanisms. *Cell* **156**, 343-58 (2014).
33. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285-91 (2016).
34. Kircher, M. *et al.* A general framework for estimating the relative pathogenicity of human genetic variants. *Nat Genet* **46**, 310-5 (2014).
35. Dimas, A.S. *et al.* Impact of type 2 diabetes susceptibility variants on quantitative glycemic traits reveals mechanistic heterogeneity. *Diabetes* **63**, 2158-71 (2014).
36. 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).
37. Altshuler, D. & Daly, M. Guilt beyond a reasonable doubt. *Nat Genet* **39**, 813-5 (2007).

FIGURE LEGENDS

Figure 1 | Posterior probabilities for coding variants across loci with annotation-informed

priors. Fine-mapping of 37 distinct association signals was performed using European ancestry GWAS meta-analysis including 50,160 T2D cases and 465,272 controls. For each signal, we constructed a credible set of variants accounting for 99% of the posterior probability of driving the association, incorporating an “annotation informed” prior model of causality which “boosts” the posterior probability of driving the association signal that is attributed to coding variants. Each bar represents a signal with the total probability attributed to the coding variants within the 99% credible set plotted on the y-axis. When the probability (bar) is split across multiple coding variants (at least 0.05 probability attributed to a variant) at a particular locus, these are indicated by blue, pink, yellow, and green colours. The combined probability of the remaining coding variants is highlighted in grey.

RREB1(a): RREB1 p. Asp1171Asn; RREB1(b): RREB1 p.Ser1499Tyr; HNF1A(a): HNF1A p.Ala146Val; HNF1A(b): HNF1A p.Ile75Leu; PPIP5K2† : PPIP5K2 p.Ser1207Gly; MTMR3†: MTMR3 p.Asn960Ser; IL17REL†: IL17REL p.Gly70Arg; NBEAL2†: NBEAL2 p.Arg511Gly, KIF9†: KIF9 p.Arg638Trp.

Figure 2 | Plot of measures of variant-specific and gene-specific features of distinct coding signals to assess the functional impact of coding alleles.

Each point represents a coding variant with the minor allele frequency plotted on the x-axis and the Combined Annotation Dependent Depletion score (CADD-score) plotted on the y-axis. Size of each point varies with the measure of intolerance of the gene to loss of function variants (pLI) and the colour represents the fine-mapping group each variant is assigned to. Group 1: signal is driven by coding variant. Group 2: signal attributable to non-coding variants. Group 3: consistent with partial role for coding variants. Group 4: Unclassified category; includes *PAX4*, *ZHX3*, and signal at *TCF19* within the MHC region where we did not perform fine-mapping. Inset: plot shows the distribution of CADD-score between different groups. The plot is a combination of violin plots and box plots; width of each violin indicates frequency at the corresponding CADD-score and box plots show the median and the 25% and 75% quantiles. *P* value indicates significance from two-sample Kolmogorov-Smirnov test.

958 **Table 1 | Summary of discovery and fine-mapping analyses of the 40 index coding variants associated with T2D ($p < 2.2 \times 10^{-7}$).**

Discovery meta-analysis using exome-array component: 81,412 T2D cases and 370,832 controls from diverse ancestries															Fine-mapping meta-analysis using GWAS: 50,160 T2D cases and 465,272 controls from European ancestry					
Locus	Index variant	rs ID	Chr	Pos	Alleles	RAF	BMI unadjusted				BMI adjusted				RAF	OR	L95	U95	p-value	Group
					R/O		OR	L95	U95	p-value	OR	L95	U95	p-value						
Previously reported T2D associated loci																				
MACF1	MACF1 p.Met1424Val	rs2296172	1	39,835,817	G/A	0.193	1.06	1.05	1.08	6.7x10 ⁻¹⁶	1.04	1.03	1.06	5.9x10 ⁻⁸	0.22	1.08	1.06	1.1	1.6x10 ⁻¹⁵	3
GCKR	GCKR p.Pro446Leu	rs1260326	2	27,730,940	C/T	0.630	1.06	1.05	1.08	5.3x10 ⁻²⁵	1.06	1.04	1.07	3.2x10 ⁻¹⁸	0.607	1.05	1.04	1.07	9.1x10 ⁻¹⁰	1
THADA	THADA p.Cys845Tyr	rs35720761	2	43,519,977	C/T	0.895	1.08	1.05	1.1	4.6x10 ⁻¹⁵	1.07	1.05	1.10	8.3x10 ⁻¹⁶	0.881	1.1	1.07	1.12	3.4x10 ⁻¹²	2
GRB14	COBL1 p.Asn901Asp	rs7607980	2	165,551,201	T/C	0.879	1.08	1.06	1.11	8.6x10 ⁻²⁰	1.09	1.07	1.12	5.0x10 ⁻²³	0.871	1.08	1.06	1.11	3.6x10 ⁻¹⁰	2
PPARG	PPARG p.Pro12Ala	rs1801282	3	12,393,125	C/G	0.887	1.09	1.07	1.11	1.4x10 ⁻¹⁷	1.10	1.07	1.12	2.7x10 ⁻¹⁹	0.876	1.12	1.09	1.14	3.7x10 ⁻¹⁷	3
IGF2BP2	SEN2 p.Thr291Lys	rs6762208	3	185,331,165	A/C	0.367	1.03	1.01	1.04	1.6x10 ⁻⁶	1.03	1.02	1.05	3.0x10 ⁻⁸	0.339	1.02	1.01	1.04	0.01	2
WFS1	WFS1 p.Val333Ile	rs1801212	4	6,302,519	A/G	0.748	1.07	1.06	1.09	1.1x10 ⁻²⁴	1.07	1.05	1.08	7.1x10 ⁻²¹	0.703	1.07	1.05	1.09	4.1x10 ⁻¹³	2
PAM-PP1P5K2	PAM p.Asp336Gly	rs35658696	5	102,338,811	G/A	0.045	1.13	1.10	1.17	1.2x10 ⁻¹⁶	1.13	1.09	1.17	7.4x10 ⁻¹⁵	0.051	1.17	1.13	1.22	2.5x10 ⁻¹⁷	1
RREB1	RREB1 p.Asp1171Asn	rs9379084	6	7,231,843	G/A	0.884	1.08	1.06	1.11	1.1x10 ⁻¹³	1.10	1.07	1.13	1.5x10 ⁻¹⁷	0.888	1.09	1.06	1.12	1.1x10 ⁻⁹	1
	RREB1 p.Ser1499Tyr	rs35742417	6	7,247,344	C/A	0.836	1.04	1.03	1.06	5.5x10 ⁻⁸	1.04	1.02	1.06	2.2x10 ⁻⁷	0.817	1.04	1.02	1.07	0.00012	2
MHC	TCF19 p.Met131Val	rs2073721	6	31,129,616	G/A	0.749	1.04	1.02	1.05	1.6x10 ⁻¹⁰	1.04	1.02	1.05	2.3x10 ⁻⁹	N/A	N/A	N/A	N/A	N/A	N/A
PAX4	PAX4 p.Arg190His	rs2233580	7	127,253,550	T/C	0.029	1.36	1.25	1.48	1.8x10 ⁻¹²	1.38	1.26	1.51	4.2x10 ⁻¹³	0	N/A	N/A	N/A	N/A	N/A
SLC30A8	SLC30A8 p.Arg276Trp	rs13266634	8	118,184,783	C/T	0.691	1.09	1.08	1.11	1.9x10 ⁻⁴⁷	1.09	1.08	1.11	1.3x10 ⁻⁴⁷	0.683	1.12	1.1	1.14	8.2x10 ⁻³⁶	1
GPSM1	GPSM1 p.Ser391Leu	rs60980157	9	139,235,415	C/T	0.771	1.06	1.05	1.08	3.2x10 ⁻¹⁶	1.06	1.05	1.08	6.6x10 ⁻¹⁶	0.756	1.06	1.04	1.09	8.3x10 ⁻⁸	3
KCNJ11-ABCC8	KCNJ11 p.Lys29Glu	rs5219	11	17,409,572	T/C	0.364	1.06	1.05	1.07	5.7x10 ⁻²²	1.07	1.05	1.08	1.5x10 ⁻²²	0.381	1.07	1.05	1.09	8.1x10 ⁻¹⁶	1
CENTD2	ARAP1 p.Gln802Glu	rs56200889	11	72,408,055	G/C	0.733	1.04	1.02	1.05	4.8x10 ⁻⁸	1.05	1.03	1.06	5.2x10 ⁻¹⁰	0.727	1.05	1.03	1.07	2.3x10 ⁻⁸	2
KLHDC5	MRPS35 p.Gly43Arg	rs1127787	12	27,867,727	G/A	0.850	1.06	1.04	1.08	1.4x10 ⁻¹¹	1.05	1.03	1.07	1.5x10 ⁻⁸	0.842	1.06	1.04	1.09	2.2x10 ⁻⁷	2
HNF1A	HNF1A p.Ile75Leu	rs1169288	12	121,416,650	C/A	0.323	1.04	1.03	1.06	1.1x10 ⁻¹¹	1.04	1.02	1.06	1.9x10 ⁻¹⁰	0.33	1.05	1.04	1.07	4.6x10 ⁻⁹	1
	HNF1A p.Ala146Val	rs1800574	12	121,416,864	T/C	0.029	1.11	1.06	1.15	6.1x10 ⁻⁸	1.10	1.06	1.15	1.3x10 ⁻⁷	0.03	1.16	1.1	1.21	5.0x10 ⁻⁹	1
MPHOSPH9	SBNO1 p.Ser729Asn	rs1060105	12	123,806,219	C/T	0.815	1.04	1.02	1.06	5.7x10 ⁻⁷	1.04	1.02	1.06	1.1x10 ⁻⁷	0.787	1.04	1.02	1.06	3.6x10 ⁻⁵	2
CILP2	TM6SF2 p.Glu167Lys	rs58542926	19	19,379,549	T/C	0.076	1.07	1.05	1.10	4.8x10 ⁻¹²	1.09	1.06	1.11	3.4x10 ⁻¹⁵	0.076	1.09	1.05	1.12	2.0x10 ⁻⁷	1
GIPR	GIPR p.Glu318Gln	rs1800437	19	46,181,392	C/G	0.200	1.03	1.02	1.05	7.1x10 ⁻⁵	1.06	1.04	1.07	6.8x10 ⁻¹²	0.213	1.09	1.06	1.12	4.6x10 ⁻⁹	1
HNF4A	HNF4A p.Thr139Ile	rs1800961	20	43,042,364	T/C	0.032	1.09	1.05	1.13	2.6x10 ⁻⁸	1.10	1.06	1.14	5.0x10 ⁻⁸	0.037	1.17	1.12	1.22	1.4x10 ⁻¹²	1
MTMR3-ASCC2	ASCC2 p.Asp407His	rs28265	22	30,200,761	C/G	0.925	1.09	1.06	1.11	2.1x10 ⁻¹²	1.09	1.07	1.12	4.4x10 ⁻¹⁴	0.916	1.1	1.07	1.14	9.6x10 ⁻¹¹	3
Novel T2D associated loci																				
FAM63A	FAM63A p.Tyr95Asn	rs140386498	1	150,972,959	A/T	0.988	1.21	1.14	1.28	7.5x10 ⁻⁸	1.19	1.12	1.26	6.7x10 ⁻⁷	0.986	1.15	1.06	1.25	0.00047	3
CEP68	CEP68 p.Gly74Ser	rs7572857	2	65,296,798	G/A	0.846	1.05	1.04	1.07	8.3x10 ⁻⁹	1.05	1.03	1.07	6.6x10 ⁻⁷	0.830	1.06	1.03	1.08	6.6x10 ⁻⁷	2
KIF9	KIF9 p.Arg638Trp	rs2276853	3	47,282,303	A/G	0.588	1.02	1.01	1.04	8.0x10 ⁻⁵	1.03	1.02	1.05	5.3x10 ⁻⁸	0.602	1.04	1.02	1.05	2.6x10 ⁻⁵	3
ANKH	ANKH p.Arg187Gln	rs146886108	5	14,751,305	C/T	0.996	1.29	1.16	1.45	1.4x10 ⁻⁷	1.27	1.13	1.41	3.5x10 ⁻⁷	0.995	1.51	1.29	1.77	3.5x10 ⁻⁷	1
POC5	POC5 p.His36Arg	rs2307111	5	75,003,678	T/C	0.562	1.05	1.04	1.07	1.6x10 ⁻¹⁵	1.03	1.01	1.04	2.1x10 ⁻⁵	0.606	1.06	1.05	1.08	1.1x10 ⁻¹²	1
LPL	LPL p.Ser474*	rs328	8	19,819,724	C/G	0.903	1.05	1.03	1.08	6.8x10 ⁻⁹	1.05	1.03	1.07	2.3x10 ⁻⁷	0.901	1.08	1.05	1.11	7.1x10 ⁻⁸	1
PLCB3 [†]	PLCB3 p.Ser778Leu	rs35169799	11	64,031,241	T/C	0.071	1.05	1.02	1.08	1.3x10 ⁻⁵	1.06	1.03	1.09	1.8x10 ⁻⁷	0.065	1.07	1.04	1.11	3.8x10 ⁻⁵	1
TPCN2	TPCN2 p.Val219Ile	rs72928978	11	68,831,364	G/A	0.890	1.05	1.02	1.07	5.2x10 ⁻⁷	1.05	1.03	1.07	1.8x10 ⁻⁸	0.847	1.03	1.00	1.05	0.042	2
WSCD2	WSCD2 p.Thr113Ile	rs3764002	12	108,618,630	C/T	0.719	1.03	1.02	1.05	3.3x10 ⁻⁸	1.03	1.02	1.05	1.2x10 ⁻⁷	0.736	1.05	1.03	1.07	8.1x10 ⁻⁷	1
ZZEF1	ZZEF1 p.Ile402Val	rs781831	17	3,947,644	C/T	0.422	1.04	1.03	1.05	8.3x10 ⁻¹¹	1.03	1.02	1.05	1.8x10 ⁻⁷	0.407	1.04	1.02	1.05	2.1x10 ⁻⁵	2
MLX	MLX p.Gln139Arg	rs665268	17	40,722,029	G/A	0.294	1.04	1.02	1.05	2.0x10 ⁻⁸	1.03	1.02	1.04	1.1x10 ⁻⁵	0.280	1.04	1.02	1.06	5.2x10 ⁻⁶	2

<i>TTL6</i>	<i>TTL6</i> p.Glu712Asp	rs2032844	17	46,847,364	C/A	0.754	1.04	1.02	1.06	1.2x10 ⁻⁷	1.03	1.01	1.04	0.00098	0.750	1.04	1.02	1.06	9.5x10 ⁻⁵	3
<i>C17orf58</i> [†]	<i>C17orf58</i> p.Ile92Val	rs9891146	17	65,988,049	T/C	0.277	1.04	1.02	1.06	1.3x10 ⁻⁷	1.02	1.00	1.04	0.00058	0.269	1.05	1.03	1.07	1.7x10 ⁻⁷	2
<i>ZHX3</i> [†]	<i>ZHX3</i> p.Asn310Ser	rs17265513	20	39,832,628	C/T	0.211	1.05	1.03	1.07	9.2x10 ⁻⁸	1.04	1.02	1.05	2.9x10 ⁻⁶	0.208	1.02	1.00	1.04	0.068	N/A
<i>PNPLA3</i>	<i>PNPLA3</i> p.Ile148Met	rs738409	22	44,324,727	G/C	0.239	1.04	1.03	1.05	2.1x10 ⁻¹⁰	1.05	1.03	1.06	2.8x10 ⁻¹¹	0.230	1.05	1.03	1.07	5.8x10 ⁻⁶	1
<i>PIM3</i>	<i>PIM3</i> p.Val300Ala	rs4077129	22	50,356,693	T/C	0.276	1.04	1.02	1.05	1.9x10 ⁻⁷	1.04	1.02	1.06	3.5x10 ⁻⁸	0.280	1.04	1.02	1.06	8.7x10 ⁻⁵	3

959

960 Chr: chromosome. Pos: Position build 37. RAF: risk allele frequency. R: risk allele. O: other allele. BMI: body mass index. OR: odds ratio. L95: lower 95% confidence interval.

961 U95: upper 95% confidence interval. GWAS: genome wide association studies.[†]Summary statistics from European ancestry specific meta-analyses of 48,286 cases and

962 250,671 controls. Fine-mapping group 1: signal is driven by coding variant, group 2: signal attributable to non-coding variants, and group 3: consistent with partial role for

963 coding variants. *p*-values are based on the meta-analyses of discovery stage and fine-mapping studies as appropriate.

964

965

966

967

968

969

970 **Table 2 | Posterior probabilities for coding variants within 99% credible set across loci**
971 **with annotation-informed and functionally-unweighted prior based on fine-mapping**
972 **analysis performed using 50,160 T2D cases and 465,272 controls of European ancestry.**
973

Locus	Variant	rs ID	Chr	Position	Posterior probability		Cumulative posterior probability attributed to coding variants	
					PPA	aiPPA	PPA	aiPPA
MACF1	MACF1 p.Ile39Val	rs16826069	1	39,797,055	0.012	0.240	0.032	0.628
	MACF1 p.Met1424Val	rs2296172	1	39,835,817	0.011	0.224		
	MACF1 p.Lys1625Asn	rs41270807	1	39,801,815	0.008	0.163		
FAM63A	FAM63A p.Tyr95Asn	rs140386498	1	150,972,959	0.005	0.129	0.012	0.303
GCKR	GCKR p. Pro 446Leu	rs1260326	2	27,730,940	0.773	0.995	0.773	0.995
THADA	THADA p.Cys845Tyr	rs35720761	2	43,519,977	<0.001	0.011	0.003	0.120
	THADA p.Thr897Ala	rs7578597	2	43,732,823	0.003	0.107		
CEP68	CEP68 p.Gly74Ser	rs7572857	2	65,296,798	<0.001	0.004	<0.001	0.004
GRB14	COBLL1 p.Asn901Asp	rs7607980	2	165,551,201	0.006	0.160	0.006	0.160
PPARG	PPARG p.Pro12Ala	rs1801282	3	12,393,125	0.023	0.410	0.024	0.410
KIF9	SETD2 p.Pro1962Lys	rs4082155	3	47,125,385	0.008	0.171	0.018	0.384
	NBEAL2 p.Arg511Gly	rs11720139	3	47,036,756	0.005	0.097		
	KIF9 p.Arg638Trp	rs2276853	3	47,282,303	0.003	0.059		
IGF2BP2	SEN2 p.Thr291Lys	rs6762208	3	185,331,165	<0.001	<0.001	<0.001	<0.001
WFS1	WFS1 p.Val333Ile	rs1801212	4	6,302,519	<0.001	0.001	<0.001	0.004
ANKH	ANKH p.Arg187Gln	rs146886108	5	14,751,305	0.459	0.972	0.447	0.972
POC5	POC5 p.His36Arg	rs2307111	5	75,003,678	0.697	0.954	0.702	0.986
PAM-PIIP5K2	PAM p.Asp336Gly	rs35658696	5	102,338,811	0.288	0.885	0.309	0.947
	PIIP5K2 p.Ser1207Gly	rs36046591	5	102,537,285	0.020	0.063		
RREB1 p.Asp1171Asn	RREB1 p.Asp1171Asn	rs9379084	6	7,231,843	0.920	0.997	0.920	0.997
RREB1 p.Ser1499Tyr	RREB1 p.Ser1499Tyr	rs35742417	6	7,247,344	<0.001	0.013	0.005	0.111
LPL	LPL p.Ser474*	rs328	8	19,819,724	0.023	0.832	0.023	0.832
SLC30A8	SLC30A8 p.Arg276Trp	rs13266634	8	118,184,783	0.295	0.823	0.295	0.823
GPSM1	GPSM1 p.Ser391Leu	rs60980157	9	139,235,415	0.031	0.557	0.031	0.557
KCNU11-ABCC8	KCNU11 p.Val250Ile	rs5215	11	17,408,630	0.208	0.412	0.481	0.951
	KCNU11 p.Lys29Glu	rs5219	11	17,409,572	0.190	0.376		
	ABCC8 p.Ala1369Ser	rs757110	11	17,418,477	0.083	0.163		
PLCB3	PLCB3 p.Ser778Leu	rs35169799	11	64,031,241	0.113	0.720	0.130	0.830
TPCN2	TPCN2 p.Val219Ile	rs72928978	11	68,831,364	<0.001	0.004	0.006	0.140
CENTD2	ARAP1 p.Gln802Glu	rs56200889	11	72,408,055	<0.001	<0.001	<0.001	<0.001
KLHDC5	MRPS35 p.Gly43Arg	rs1127787	12	27,867,727	<0.001	<0.001	<0.001	<0.001
WSCD2	WSCD2 p.Thr113Ile	rs3764002	12	108,618,630	0.281	0.955	0.282	0.958
HNF1A p.Ile75Leu	HNF1A Gly226Ala	rs56348580	12	121,432,117	0.358	0.894	0.358	0.894
	HNF1A p.Ile75Leu	rs1169288	12	121,416,650	<0.001	<0.001		
HNF1A p.Ala146Val	HNF1A p.Ala146Val	rs1800574	12	121,416,864	0.269	0.867	0.280	0.902
MPHOSPH9	SBN01 p.Ser729Asn	rs1060105	12	123,806,219	0.002	0.054	0.002	0.057
ZZEF1	ZZEF1 p.Ile402Val	rs781831	17	3,947,644	<0.001	0.001	<0.001	0.018
MLX	MLX p.Gln139Arg	rs665268	17	40,722,029	0.002	0.038	0.002	0.039
TTLL6	TTLL6 p.Glu712Asp	rs2032844	17	46,847,364	<0.001	<0.001	0.016	0.305
	CALCOCO2 p.Pro347Ala	rs10278	17	46,939,658	0.0100	0.187		
	SNF8 p.Arg155His	rs57901004	17	47,011,897	0.005	0.092		
C17orf58	C17orf58 p.Ile92Val	rs9891146	17	65,988,049	<0.001	0.009	<0.001	0.009
CILP2	TM6SF2 p.Glu167Lys	rs58542926	19	19,379,549	0.211	0.732	0.263	0.913
	TM6SF2 p.Leu156Pro	rs187429064	19	19,380,513	0.049	0.172		
GIPR	GIPR p.Glu318Gln	rs1800437	19	46,181,392	0.169	0.901	0.169	0.901
ZHX3	ZHX3 p.Asn310Ser	rs17265513	20	39,832,628	<0.001	0.003	0.003	0.110
HNF4A	HNF4A p.Thr139Ile	rs1800961	20	43,042,364	1.000	1.000	1.00	1.000
MTMR3-ASCC2	ASCC2 p.Asp407His	rs28265	22	30,200,761	0.011	0.192	0.028	0.481
	ASCC2 p.Pro423Ser	rs36571	22	30,200,713	0.007	0.116		
	ASCC2 p.Val123Ile	rs11549795	22	30,221,120	0.006	0.107		
	MTMR3 p.Asn960Ser	rs41278853	22	30,416,527	0.004	0.065		
PNPLA3	PNPLA3 p.Ile148Met	rs738409	22	44,324,727	0.112	0.691	0.130	0.806
	PARVB p.Trp37Arg	rs1007863	22	44,395,451	0.017	0.103		
PIM3	IL17REL p.Leu333Pro	rs5771069	22	50,435,480	0.041	0.419	0.047	0.475
	IL17REL p.Gly70Arg	rs9617090	22	50,439,194	0.005	0.054		
	PIM3 p.Val300Ala	rs4077129	22	50,356,693	<0.001	0.002		

974
975 19th chromosome. Pos: Position build 37. PPA: functionally-unweighted prior; aiPPA: annotation informed prior. Index
976 coding variants are highlighted in bold.

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.

Derivation of significance thresholds. We considered five categories of annotation¹⁶ of variants on the exome array in order of decreasing effect on biological function: (1) PTVs (stop-gain and stop-loss, frameshift indel, donor and acceptor splice-site, and initiator codon variants, $n_1=8,388$); (2) moderate-impact variants (missense, in-frame indel, and splice region variants, $n_2=216,114$); (3) low-impact variants (synonymous, 3' and 5' UTR, and upstream and downstream variants, $n_3=8,829$); (4) other variants mapping to DNase I hypersensitive sites (DHS) in any of 217 cell types⁸ (DHS, $n_4=3,561$); and (5) other variants not mapping to DHS ($n_5=10,578$). To account for the greater prior probability of causality for variants with greater effect on biological function, we determined a weighted Bonferroni-corrected significance threshold on the basis of reported enrichment¹⁶, denoted w_i , in each annotation category, i : $w_1=165$; $w_2=33$; $w_3=3$; $w_4=1.5$; $w_5=0.5$. For coding variants (annotation categories 1 and 2):

$$\alpha = \frac{0.05 \sum_{i=1}^2 n_i w_i}{(\sum_{i=1}^2 n_i)(\sum_{i=1}^5 n_i w_i)} = 2.21 \times 10^{-7}.$$

We note that this threshold is similar to a simple Bonferroni correction for the total number of coding variants on the array, which would yield:

$$\alpha = \frac{0.05}{224502} = 2.23 \times 10^{-7}.$$

For non-coding variants (annotation categories 3, 4 and 5) the weighted Bonferroni-corrected significance threshold is:

$$\alpha = \frac{0.05 \sum_{i=3}^5 n_i w_i}{(\sum_{i=3}^5 n_i)(\sum_{i=1}^5 n_i w_i)} = 9.45 \times 10^{-9}.$$

DISCOVERY: Exome-array study-level analyses. Within each study, genotype calling and quality control were undertaken according to protocols developed by the UK Exome Chip Consortium or the CHARGE central calling effort³⁸ (**Supplementary Table 1**). Within each study, variants were then excluded for the following reasons: (i) not mapping to autosomes or X chromosome; (ii) multi-allelic and/or insertion-deletion; (iii) monomorphic; (iv) call rate <99%; or (v) exact $p < 10^{-4}$ for deviation from Hardy-Weinberg equilibrium (autosomes only).

We tested association of T2D with each variant in a linear mixed model, implemented in RareMetalWorker¹⁷, using a genetic relationship matrix (GRM) to account for population structure and relatedness. For participants from family-based studies, known relationships were incorporated directly in the GRM. For founders and participants from population-based studies, the GRM was constructed from pair-wise identity by descent (IBD) estimates based on LD pruned ($r^2 < 0.05$) autosomal variants with MAF $\geq 1\%$ (across cases and controls combined), after exclusion of those in high LD and complex regions^{39,40}, and those mapping to established T2D loci. We considered additive, dominant, and recessive models for the effect of the minor allele, adjusted for age and sex (where appropriate) and additional study-specific covariates (**Supplementary Table 2**). Analyses were also performed with and without adjustment for BMI (where available Supplementary Table 2).

For single-variant association analyses, variants with minor allele count ≤ 10 in cases and controls combined were excluded. Association summary statistics for each analysis were corrected for residual inflation by means of genomic control⁴¹, calculated after excluding variants mapping to established T2D susceptibility loci. For gene-based analyses, we made no variant exclusions on the basis of minor allele count.

DISCOVERY: Exome-sequence analyses. We used summary statistics of T2D association from analyses conducted on 8,321 T2D cases and 8,421 controls across different ancestries, all genotyped using exome sequencing. Details of samples included, sequencing, and quality control are described elsewhere^{12,15} (<http://www.type2diabetesgenetics.org/>). Samples were subdivided into 15 sub-groups according to ancestry and study of origin. Each sub-group was analysed independently, with sub-group specific principal components and genetic relatedness matrices. Association tests were performed with both a linear mixed model, as implemented in EMMAX⁴², using covariates for sequencing batch, and the Firth

test, using covariates for principal components and sequencing batch. Related samples were excluded from the Firth analysis but maintained in the linear mixed model analysis. Variants were then filtered from each sub-group analysis, according to call rate, differential case-control missing-ness, or deviation from Hardy-Weinberg equilibrium (as computed separately for each sub-group). Association statistics were then combined via a fixed-effects inverse-variance weighted meta-analysis, at both the level of ancestry as well as across all samples. P-values were taken from the linear mixed model analysis, while effect sizes estimates were taken from the Firth analysis. Analyses were performed with and without adjustment for BMI. From exome sequence summary statistics, we extracted variants passing quality control and present on the exome array.

DISCOVERY: GWAS analyses. The UK Biobank is a large detailed prospective study of more than 500,000 participants aged 40-69 years when recruited in 2006-2010¹³. Prevalent T2D status was defined using self-reported medical history and medication in UK Biobank participants⁴³. Participants were genotyped with the UK Biobank Axiom Array or UK BiLEVE Axiom Array, and quality control and population structure analyses were performed centrally at UK Biobank. We defined a subset of “white European” ancestry samples (n=120,286) as those who both self-identified as white British and were confirmed as ancestrally “Caucasian” from the first two axes of genetic variation from principal components analysis. Imputation was also performed centrally at UK Biobank for the autosomes only, up to a merged reference panel from the 1000 Genomes Project (multi-ethnic, phase 3, October 2014 release)⁴⁴ and the UK10K Project⁹. We used SNPTESTv2.5⁴⁵ to test for association of T2D with each SNP in a logistic regression framework under an additive model, and after adjustment for age, sex, six axes of genetic variation, and genotyping array as covariates. Analyses were performed with and without adjustment for BMI, after removing related individuals.

GERA is a large multi-ethnic population-based cohort, created for investigating the genetic and environmental basis of age-related diseases [dbGaP phs000674.p1]. T2D status is based on ICD-9 codes in linked electronic medical health records, with all other participants defined as controls. Participants have previously been genotyped using one of four custom arrays, which have been designed to maximise coverage of common and low-frequency variants in non-Hispanic white, East Asian, African American, and Latino

ethnicities^{46,47}. Methods for quality control have been described previously¹⁴. Each of the four genotyping arrays were imputed separately, up to the 1000 Genomes Project reference panel (autosomes, phase 3, October 2014 release; X chromosome, phase 1, March 2012 release) using IMPUTEv2.3^{48,49}. We used SNPTESTv2.5⁴⁵ to test for association of T2D with each SNP in a logistic regression framework under an additive model, and after adjustment for sex and nine axes of genetic variation from principal components analysis as covariates. BMI was not available for adjustment in GERA.

For UK Biobank and GERA, we extracted variants passing standard imputation quality control thresholds (IMPUTE info \geq 0.4)⁵⁰ and present on the exome array. Association summary statistics under an additive model were corrected for residual inflation by means of genomic control⁴¹, calculated after excluding variants mapping to established T2D susceptibility loci: GERA (λ =1.097 for BMI unadjusted analysis) and UK Biobank (λ =1.043 for BMI unadjusted analysis, λ =1.056 for BMI adjusted analysis).

DISCOVERY: Single-variant meta-analysis. We aggregated association summary statistics under an additive model across studies, with and without adjustment for BMI, using METAL⁵¹: (i) effective sample size weighting of Z-scores to obtain p -values; and (ii) inverse variance weighting of log-odds ratios. For exome-array studies, allelic effect sizes and standard errors obtained from the RareMetalWorker linear mixed model were converted to the log-odds scale prior to meta-analysis to correct for case-control imbalance⁵².

The European-specific meta-analyses aggregated association summary statistics from a total of 48,286 cases and 250,671 controls from: (i) 33 exome-array studies of European ancestry; (ii) exome-array sequence from individuals of European ancestry; and (iii) GWAS from UK Biobank. Note that non-coding variants represented on the exome array were not available in exome sequence. The European-specific meta-analyses were corrected for residual inflation by means of genomic control⁴¹, calculated after excluding variants mapping to established T2D susceptibility loci: λ =1.091 for BMI unadjusted analysis and λ =1.080 for BMI adjusted analysis.

The trans-ethnic meta-analyses aggregated association summary statistics from a total of 81,412 cases and 370,832 controls across all studies (51 exome array studies, exome sequence, and GWAS from UK Biobank and GERA), irrespective of ancestry. Note that non-coding variants represented on the exome array were not available in exome sequence. The

trans-ethnic meta-analyses were corrected for residual inflation by means of genomic control⁴¹, calculated after excluding variants mapping to established T2D susceptibility loci: $\lambda=1.073$ for BMI unadjusted analysis and $\lambda=1.068$ for BMI adjusted analysis. Heterogeneity in allelic effect sizes between exome-array studies contributing to the trans-ethnic meta-analysis was assessed by Cochran's Q statistic⁵³.

DISCOVERY: Detection of distinct association signals. Conditional analyses were undertaken to detect association signals by inclusion of index variants and/or tags for previously reported non-coding GWAS lead SNPs as covariates in the regression model at the study level. Within each exome-array study, approximate conditional analyses were undertaken under a linear mixed model using RareMetal¹⁷, which uses score statistics and the variance-covariance matrix from the RareMetalWorker single-variant analysis to estimate the correlation in effect size estimates between variants due to LD. Study-level allelic effect sizes and standard errors obtained from the approximate conditional analyses were converted to the log-odds scale to correct for case-control imbalance⁵². Within each GWAS, exact conditional analyses were performed under a logistic regression model using SNPTTESTv2.5⁴⁵. GWAS variants passing standard imputation quality control thresholds (IMPUTE info ≥ 0.4)⁵⁰ and present on the exome array were extracted for meta-analysis.

Association summary statistics were aggregated across studies, with and without adjustment for BMI, using METAL⁵¹: (i) effective sample size weighting of Z-scores to obtain p -values; and (ii) inverse variance weighting of log-odds ratios.

We defined novel loci as mapping >500kb from a previously reported lead GWAS SNP. We performed conditional analyses where a novel signal mapped close to a known GWAS locus, and the lead GWAS SNP at that locus is present (or tagged) on the exome array (**Supplementary Table 5**).

DISCOVERY: Non-additive association models. For exome-array studies only, we aggregated association summary statistics under recessive and dominant models across studies, with and without adjustment for BMI, using METAL⁵¹: (i) effective sample size weighting of Z-scores to obtain p -values; and (ii) inverse variance weighting of log-odds ratios. Allelic effect sizes and standard errors obtained from the RareMetalWorker linear mixed model were converted to the log-odds scale prior to meta-analysis to correct for case-control

imbalance⁵². The European-specific meta-analyses aggregated association summary statistics from a total of 41,066 cases and 136,024 controls from 33 exome-array studies of European ancestry. The European-specific meta-analyses were corrected for residual inflation by means of genomic control⁴¹, calculated after excluding variants mapping to established T2D susceptibility loci: $\lambda=1.076$ and $\lambda=1.083$ for BMI unadjusted analysis, under recessive and dominant models respectively, and $\lambda=1.081$ and $\lambda=1.062$ for BMI adjusted analysis, under recessive and dominant models respectively. The trans-ethnic meta-analyses aggregated association summary statistics from a total of 58,425 cases and 188,032 controls across all exome-array studies, irrespective of ancestry. The trans-ethnic meta-analyses were corrected for residual inflation by means of genomic control⁴¹, calculated after excluding variants mapping to established T2D susceptibility loci: $\lambda=1.041$ and $\lambda=1.071$ for BMI unadjusted analysis, under recessive and dominant models respectively, and $\lambda=1.031$ and $\lambda=1.063$ for BMI adjusted analysis, under recessive and dominant models respectively.

DISCOVERY: Gene-based meta-analyses. For exome-array studies only, we aggregated association summary statistics under an additive model across studies, with and without adjustment for BMI, using RareMetal¹⁷. This approach uses score statistics and the variance-covariance matrix from the RareMetalWorker single-variant analysis to estimate the correlation in effect size estimates between variants due to LD. We performed gene-based analyses using a burden test (assuming all variants have same direction of effect on T2D susceptibility) and SKAT (allowing variants to have different directions of effect on T2D susceptibility). We used two previously defined filters for annotation and MAF¹⁸ to define group files: (i) strict filter, including 44,666 variants; and (ii) broad filter, including all variants from the strict filter, and 97,187 additional variants.

We assessed the contribution of each variant to gene-based signals by performing approximate conditional analyses. We repeated RareMetal analyses for the gene, excluding each variant in turn from the group file, and compared the strength of the association signal.

Fine-mapping of coding variant association signals with T2D susceptibility. We defined a locus as mapping 500kb up- and down-stream of each index coding variant (**Supplementary Table 5**), excluding the MHC. Our fine-mapping analyses aggregated association summary

statistics from 24 GWAS incorporating 50,160 T2D cases and 465,272 controls of European ancestry from the DIAGRAM Consortium (**Supplementary Table 9**). Each GWAS was imputed using miniMAC¹² or IMPUTEv2^{48,49} up to high-density reference panels: (i) 22 GWAS were imputed up to the Haplotype Reference Consortium²⁰; (ii) the UK Biobank GWAS was imputed to a merged reference panel from the 1000 Genomes Project (multi-ethnic, phase 3, October 2014 release)⁴⁴ and the UK10K Project⁹; and (iii) the deCODE GWAS was imputed up to the deCODE Icelandic population-specific reference panel based on whole-genome sequence data¹⁹. Association with T2D susceptibility was tested for each remaining variant using logistic regression, adjusting for age, sex, and study-specific covariates, under an additive genetic model. Analyses were performed with and without adjustment for BMI. For each study, variants with minor allele count<5 (in cases and controls combined) or those with imputation quality $r^2\text{-hat}<0.3$ (miniMAC) or proper-info<0.4 (IMPUTE2) were removed. Association summary statistics for each analysis were corrected for residual inflation by means of genomic control⁴¹, calculated after excluding variants mapping to established T2D susceptibility loci.

We aggregated association summary statistics across studies, with and without adjustment for BMI, in a fixed-effects inverse variance weighted meta-analysis, using METAL⁵¹. The BMI unadjusted meta-analysis was corrected for residual inflation by means of genomic control ($\lambda=1.012$)⁴¹, calculated after excluding variants mapping to established T2D susceptibility loci. No adjustment was required for BMI adjusted meta-analysis ($\lambda=0.994$). From the meta-analysis, variants were extracted that were present on the HRC panel and reported in at least 50% of total effective sample size.

We included 37 of the 40 identified coding variants in fine-mapping analyses, excluding three that were not amenable to fine-mapping in the GWAS data sets: (i) the locus in the major histocompatibility complex because of the extended and complex structure of LD across the region, which complicates fine-mapping efforts; (ii) the East Asian specific *PAX4* p.Arg190His (rs2233580) signal, since the variant was not present in European ancestry GWAS; and (iii) *ZHX3* p.Asn310Ser (rs4077129) because the variant was only weakly associated with T2D in the GWAS data sets used for fine-mapping.

To delineate distinct association signals in four regions, we undertook approximate conditional analyses, implemented in GCTA⁵⁴, to adjust for the index coding variants and non-coding lead GWAS SNPs: (i) *RREB1* p. Asp1171Asn (rs9379084), p.Ser1499Tyr

(rs35742417), and rs9505118; (ii) *HNF1A* p.Ile75Leu (rs1169288) and p.Ala146Val (rs1800574); (iii) *GIPR* p.Glu318Gln (rs1800437) and rs8108269; and (iv) *HNF4A* p.Thr139Ile (rs1800961) and rs4812831. We made use of summary statistics from the fixed-effects meta-analyses (BMI unadjusted for *RREB1*, *HNF1A*, and *HNF4A*, and BMI adjusted for *GIPR* as this signal was only seen in BMI adjusted analysis) and genotype data from 5,000 random individuals of European ancestry from the UK Biobank, as reference for LD between genetic variants across the region.

For each association signal, we first calculated an approximate Bayes' factor⁵⁵ in favour of association on the basis of allelic effect sizes and standard errors from the meta-analysis. 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 0.04⁵⁵.

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

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

In this expression, ρ_j denotes the prior probability that the j th variant drives the association signal, and the summation in the denominator is over all variants across the locus. We considered two prior models: (i) functionally unweighted, for which $\rho_j = 1$ for all variants; and (ii) annotation informed, for which ρ_j is determined by the functional severity of the variant. For the annotation informed prior, we considered five categories of variation¹⁶, such that: (i) $\rho_j = 165$ for PTVs; (ii) $\rho_j = 33$ for moderate-impact variants; (iii) $\rho_j = 3$ for low-impact variants; (iv) $\rho_j = 1.5$ for other variants mapping to DHS; and (v) $\rho_j = 0.5$ for all other variants.

For each locus, the 99% credible set²¹ under each prior was then constructed by: (i) ranking all variants according to their posterior probability of driving the association signal;

and (ii) including ranked variants until their cumulative posterior probability of driving the association attained or exceeded 0.99.

Functional impact of coding alleles. We used CADD³⁴ to obtain scaled Combined Annotation Dependent Depletion score (CADD-score) for each of the 40 significantly associated coding variants. The CADD method objectively integrates a range of different annotation metrics into a single measure (CADD-score), providing an estimate of deleteriousness for all known variants and an overall rank for this metric across the genome. We obtained the estimates of the intolerance of a gene to harbouring loss-of-function variants (pLI) from the ExAC data set³³. We used the Kolmogorov-Smirnov test to determine whether fine-mapping groups 1 and 2 have the same statistical distribution for each of these parameters.

T2D loci and physiological classification. To explore the different patterns of association between T2D and other anthropometric/metabolic/endocrine traits and diseases, we performed hierarchical clustering analysis. We obtained association summary statistics for a range of metabolic traits and other outcomes for 94 coding and non-coding variants that were significantly associated with T2D through collaboration or by querying publically available GWAS meta-analysis datasets. The z-score (allelic effect/SE) was aligned to the T2D-risk allele. We obtained the distance matrix amongst z-score of the loci/traits using the Euclidean measure and performed clustering using the complete agglomeration method. Clustering was visualised by constructing a dendrogram and heatmap.

DATA AVAILABILITY STATEMENT

Summary level data of the exome-array component of this project can be downloaded from the DIAGRAM consortium website <http://diagram-consortium.org/> and Accelerating Medicines Partnership T2D portal <http://www.type2diabetesgenetics.org/>.

MATERIALS & CORRESPONDENCE

Correspondence and requests for materials should be addressed to mark.mccarthy@drl.ox.ac.uk and anubha@well.ox.ac.uk. Reprints and permissions information is available at www.nature.com/reprints.

1260 38. Grove, M.L. *et al.* Best practices and joint calling of the HumanExome BeadChip: the
1261 CHARGE Consortium. *PLoS One* **8**, e68095 (2013).

1262 39. Price, A.L. *et al.* Long-range LD can confound genome scans in admixed populations.
1263 *Am J Hum Genet* **83**, 132-5; author reply 135-9 (2008).

1264 40. Weale, M.E. Quality control for genome-wide association studies. *Methods Mol Biol*
1265 **628**, 341-72 (2010).

1266 41. Devlin, B. & Roeder, K. Genomic control for association studies. *Biometrics* **55**, 997-
1267 1004 (1999).

1268 42. Kang, H.M. *et al.* Variance component model to account for sample structure in
1269 genome-wide association studies. *Nat Genet* **42**, 348-54 (2010).

1270 43. Eastwood, S.V. *et al.* Algorithms for the Capture and Adjudication of Prevalent and
1271 Incident Diabetes in UK Biobank. *PLoS One* **11**, e0162388 (2016).

1272 44. Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68-74
1273 (2015).

1274 45. Marchini, J. & Howie, B. Genotype imputation for genome-wide association studies.
1275 *Nat Rev Genet* **11**, 499-511 (2010).

1276 46. Hoffmann, T.J. *et al.* Next generation genome-wide association tool: design and
1277 coverage of a high-throughput European-optimized SNP array. *Genomics* **98**, 79-89
1278 (2011).

1279 47. Hoffmann, T.J. *et al.* Design and coverage of high throughput genotyping arrays
1280 optimized for individuals of East Asian, African American, and Latino race/ethnicity
1281 using imputation and a novel hybrid SNP selection algorithm. *Genomics* **98**, 422-30
1282 (2011).

1283 48. Howie, B.N., Donnelly, P. & Marchini, J. A flexible and accurate genotype imputation
1284 method for the next generation of genome-wide association studies. *PLoS Genet* **5**,
1285 e1000529 (2009).

1286 49. Howie, B., Fuchsberger, C., Stephens, M., Marchini, J. & Abecasis, G.R. Fast and
1287 accurate genotype imputation in genome-wide association studies through pre-
1288 phasing. *Nat Genet* **44**, 955-9 (2012).

1289 50. Winkler, T.W. *et al.* Quality control and conduct of genome-wide association meta-
1290 analyses. *Nat Protoc* **9**, 1192-212 (2014).

1291 51. Willer, C.J., Li, Y. & Abecasis, G.R. METAL: fast and efficient meta-analysis of
1292 genomewide association scans. *Bioinformatics* **26**, 2190-1 (2010).

1293 52. Cook, J.P., Mahajan, A. & Morris, A.P. Guidance for the utility of linear models in
1294 meta-analysis of genetic association studies of binary phenotypes. *Eur J Hum Genet*
1295 **25**, 240-245 (2017).

1296 53. Ioannidis, J.P., Patsopoulos, N.A. & Evangelou, E. Heterogeneity in meta-analyses of
1297 genome-wide association investigations. *PLoS One* **2**, e841 (2007).

1298 54. Yang, J. *et al.* Conditional and joint multiple-SNP analysis of GWAS summary statistics
1299 identifies additional variants influencing complex traits. *Nat Genet* **44**, 369-75, s1-3
1300 (2012).

1301 55. Wakefield, J. A Bayesian measure of the probability of false discovery in genetic
1302 epidemiology studies. *Am J Hum Genet* **81**, 208-27 (2007).

1303

1304 **URLs**

1305 Type 2 Diabetes Knowledge Portal: <http://www.type2diabetesgenetics.org/>