



et al., Min, J. L., Hemani, G., Hannon, E., deKleer, K., Castillo-Fernandez, J., Luijk, R., Carnero-Montoro, E., Lawson, D. J., Burrows, K., Suderman, M. J., Bretherick, A. D., Richardson, T. G., Klughammer, J., Lochkova, V., Sharp, G. C., Al Khleifat, A., Shatunov, A., Lacoangeli, A., ... Relton, C. L. (2021). Genomic and phenotypic insights from an atlas of genetic effects on DNA methylation. *Nature Genetics*, 53(9), 1311-1321. <https://doi.org/10.1038/s41588-021-00923-x>

Peer reviewed version

Link to published version (if available):
[10.1038/s41588-021-00923-x](https://doi.org/10.1038/s41588-021-00923-x)

[Link to publication record in Explore Bristol Research](#)
PDF-document

This is the author accepted manuscript (AAM). The final published version (version of record) is available online via Springer Nature at <https://doi.org/10.1038/s41588-021-00923-x>

University of Bristol - Explore Bristol Research

General rights

This document is made available in accordance with publisher policies. Please cite only the published version using the reference above. Full terms of use are available: <http://www.bristol.ac.uk/red/research-policy/pure/user-guides/ebr-terms/>

1 Genomic and phenomic insights from an atlas of 2 genetic effects on DNA methylation

3 Josine L Min^{1,2*}, Gibran Hemani^{1,2*}, Ellis Hannon³, Koen F Dekkers⁴, Juan Castillo-
4 Fernandez⁵, René Luijk⁴, Elena Carnero-Montoro^{5,6}, Daniel J Lawson^{1,2}, Kimberley
5 Burrows^{1,2}, Matthew Suderman^{1,2}, Andrew D Bretherick⁷, Tom G Richardson^{1,2},
6 Johanna Klughammer⁸, Valentina Iotchkova⁹, Gemma Sharp^{1,2}, Ahmad Al Khleifat¹⁰,
7 Aleksey Shatunov¹⁰, Alfredo Iacoangeli^{10,11}, Wendy L McArdle², Karen M Ho², Ashish
8 Kumar^{12,13,14}, Cilla Söderhäll¹⁵, Carolina Soriano-Tárraga¹⁶, Eva Giralt-Steinhauer¹⁶,
9 Nabila Kazmi^{1,2}, Dan Mason¹⁷, Allan F McRae¹⁸, David L Corcoran¹⁹, Karen
10 Sugden^{19,20}, Silva Kasela²¹, Alexia Cardona^{22,23}, Felix R Day²², Giovanni Cugliari^{24,25},
11 Clara Viberti^{24,25}, Simonetta Guarrera^{24,25}, Michael Lerro²⁶, Richa Gupta^{27,28}, Sailalitha
12 Bollepalli^{27,28}, Pooja Mandaviya²⁹, Yanni Zeng^{7,30,31}, Toni-Kim Clarke³², Rosie M
13 Walker^{33,34}, Vanessa Schmoll³⁵, Darina Czamara³⁵, Carlos Ruiz-Arenas^{36,37,38}, Faisal I
14 Rezwani³⁹, Riccardo E Marioni^{34,40}, Tian Lin¹⁸, Yvonne Awaloff³⁵, Marine Germain⁴¹,
15 Dylan Aïssi⁴², Ramona Zwamborn⁴³, Kristel van Eijk⁴³, Annelot Dekker⁴³, Jenny van
16 Dongen⁴⁴, Jouke-Jan Hottenga⁴⁴, Gonneke Willemsen⁴⁴, Cheng-Jian Xu⁴⁵, Guillermo
17 Barturen⁶, Francesc Català-Moll⁴⁶, Martin Kerick⁴⁷, Carol Wang⁴⁸, Phillip Melton⁴⁹,
18 Hannah R Elliott^{1,2}, Jean Shin^{50,51}, Manon Bernard⁵⁰, Idil Yet⁵, Melissa Smart⁵², Tyler
19 Gorrie-Stone⁵², BIOS Consortium, Chris Shaw^{10,53}, Ammar Al Chalabi^{10,53,54}, Susan M
20 Ring^{1,2}, Göran Pershagen¹², Erik Melén^{12,55}, Jordi Jiménez-Conde¹⁶, Jaume Roquer¹⁶,
21 Debbie A Lawlor^{1,2}, John Wright¹⁷, Nicholas G Martin⁵⁶, Grant W Montgomery¹⁸, Terrie
22 E Moffitt^{19,20,57,60}, Richie Poulton⁵⁸, Tõnu Esko^{21,59}, Lili Milani²¹, Andres Metspalu²¹, John
23 RB Perry²², Ken K Ong²², Nicholas J Wareham²², Giuseppe Matullo^{24,25}, Carlotta
24 Sacerdote^{25,61}, Salvatore Panico⁶², Avshalom Caspi^{19,20,57,60}, Louise Arseneault⁶⁰,
25 France Gagnon²⁶, Miina Ollikainen^{27,28}, Jaakko Kaprio^{27,28}, Janine F Felix^{63,64,65},
26 Fernando Rivadeneira²⁹, Henning Tiemeier^{66,67}, Marinus H van IJzendoorn^{68,69}, André G
27 Uitterlinden²⁹, Vincent WV Jaddoe^{63,64,67}, Chris Haley⁷, Andrew M McIntosh^{32,34}, Kathryn
28 L Evans^{33,34}, Alison Murray⁷⁰, Katri Räikkönen⁷¹, Jari Lahti⁷¹, Ellen A Nohr^{72,73}, Thorkild
29 IA Sørensen^{1,2,74,75}, Torben Hansen⁷⁴, Camilla Schmidt Morgen⁷⁶, Elisabeth B
30 Binder^{35,77}, Susanne Lucae³⁵, Juan Ramon Gonzalez^{36,37,38}, Mariona
31 Bustamante^{36,37,38,78}, Jordi Sunyer^{36,37,38,79}, John W Holloway^{39,80}, Wilfried Karmaus⁸¹,
32 Hongmei Zhang⁸¹, Ian J Deary³⁴, Naomi R Wray^{18,82}, John M Starr^{34,83}, Marian
33 Beekman⁴, Diana van Heemst⁸⁴, P Eline Slagboom⁴, Pierre-Emmanuel Morange⁸⁵,
34 David-Alexandre Trégouët⁴¹, Jan H Veldink⁴³, Gareth E Davies⁸⁶, Eco JC de Geus⁴⁴,
35 Dorret I Boomsma⁴⁴, Judith M Vonk⁸⁷, Bert Brunekreef^{88,89}, Gerard H Koppelman⁴⁵,
36 Marta E Alarcón-Riquelme^{6,12}, Rae-Chi Huang⁹⁰, Craig Pennell⁴⁸, Joyce van Meurs²⁹, M
37 Arfan Ikram⁶⁴, Alun D Hughes⁹¹, Therese Tillin⁹¹, Nish Chaturvedi⁹¹, Zdenka Pausova⁴⁹,
38 Tomas Paus⁹², Timothy D Spector⁵, Meena Kumari⁵², Leonard C Schalkwyk⁵², Peter M
39 Visscher^{18,82}, George Davey Smith^{1,2}, Christoph Bock⁸, Tom R Gaunt^{1,2}, Jordana T
40 Bell^{5‡}, Bastiaan T Heijmans^{4‡}, Jonathan Mill^{3‡}, Caroline L Relton^{1,2‡}

41

42 * These authors contributed equally to this research.

43 ‡These authors jointly supervised this work.

44

45 **Corresponding author:** Josine L Min, josine.min@bristol.ac.uk

46

47 Affiliations

48 ¹ MRC Integrative Epidemiology Unit, University of Bristol, Bristol, UK

49 ² Population Health Sciences, Bristol Medical School, University of Bristol, Bristol, UK

50 ³ University of Exeter Medical School, UK

51 ⁴ Molecular Epidemiology, Department of Biomedical Data Sciences, Leiden University
52 Medical Center, Leiden, The Netherlands

53 ⁵ Department of Twin Research and Genetic Epidemiology, King's College London,
54 London, UK

55 ⁶ Pfizer - University of Granada - Andalusian Government Center for Genomics and
56 Oncological Research (GENYO), Spain

57 ⁷ MRC Human Genetic Unit, Institute of Genetics and Molecular Medicine, University of
58 Edinburgh, Edinburgh, UK

59 ⁸ CeMM, Austrian Academy of Sciences, Vienna, Austria

60 ⁹ MRC Weatherall Institute of Molecular Medicine, Oxford, UK

61 ¹⁰ Department of Basic and Clinical Neuroscience, Maurice Wohl Clinical Neuroscience
62 Institute, London, UK

63 ¹¹ Department of Biostatistics and Health Informatics, King's College London, London,
64 UK

65 ¹² Institute of Environmental Medicine, Karolinska Institutet, Stockholm, Solna, Sweden

66 ¹³ Chronic Disease Epidemiology unit, Swiss Tropical and Public Health Institute, Basel,
67 Switzerland

68 ¹⁴ University of Basel, Basel, Switzerland

69 ¹⁵ Department of Women's and Children's Health, Karolinska Institutet, Stockholm,
70 Sweden

71 ¹⁶ Neurology Department, Hospital del Mar - IMIM (Institut Hospital del Mar
72 d'Investigacions Mèdiques), Barcelona, Spain

73 ¹⁷ Bradford Institute for Health Research, Bradford, UK

74 ¹⁸ Institute for Molecular Bioscience, University of Queensland, Australia

75 ¹⁹ Center for Genomic and Computational Biology, Duke University, Durham, NC, USA

76 ²⁰ Department of Psychology and Neuroscience, Duke University, Durham, NC, USA

77 ²¹ Estonian Genome Center, Institute of Genomics, University of Tartu, Estonia

78 ²² MRC Epidemiology Unit, University of Cambridge, School of Clinical Medicine,
79 Institute of Metabolic Science, Cambridge Biomedical Campus, Cambridge CB2 0QQ,
80 United Kingdom

81 ²³ Department of Genetics, University of Cambridge, Downing Street, Cambridge CB2
82 3EH, United Kingdom

83 ²⁴ Department of Medical Sciences, University of Turin, Turin, Italy

84 ²⁵ Italian Institute for Genomic Medicine (IIGM), Turin, Italy
85 ²⁶ University of Toronto, Toronto, Canada
86 ²⁷ Institute for Molecular Medicine, University of Helsinki, Helsinki, Finland
87 ²⁸ Department of Public Health, Faculty of Medicine, University of Helsinki, Helsinki,
88 Finland
89 ²⁹ Department of Internal Medicine, Erasmus University Medical Center, Rotterdam, The
90 Netherlands
91 ³⁰ Faculty of Forensic Medicine, Zhongshan School of Medicine, Sun Yat-Sen
92 University, Guangzhou, China
93 ³¹ Guangdong Province Key Laboratory of Brain Function and Disease, Zhongshan
94 School of Medicine, Sun Yat-Sen University, Guangzhou, China
95 ³² Division of Psychiatry, Royal Edinburgh Hospital, University of Edinburgh, Edinburgh
96 EH10 5HF, UK
97 ³³ Medical Genetics Section, Centre for Genomic and Experimental Medicine, Institute
98 of Genetics and Molecular Medicine, Western General Hospital, University of
99 Edinburgh, Crewe Road, Edinburgh EH4 2XU, UK
100 ³⁴ Centre for Cognitive Ageing and Cognitive Epidemiology, Department of Psychology,
101 University of Edinburgh, 7 George Square, Edinburgh EH8 9JZ, UK
102 ³⁵ Department of Translational Research in Psychiatry, Max-Planck-Institute of
103 Psychiatry, Munich, Germany
104 ³⁶ ISGlobal, Barcelona Global Health Institute, Barcelona, Spain
105 ³⁷ Universitat Pompeu Fabra (UPF), Barcelona, Spain
106 ³⁸ CIBER Epidemiología y Salud Pública (CIBERESP), Madrid, Spain
107 ³⁹ Human Development and Health, Faculty of Medicine, University of Southampton,
108 Southampton, UK
109 ⁴⁰ Institute of Genetics and Molecular Medicine, University of Edinburgh, Edinburgh, UK
110 ⁴¹ INSERM UMR_S 1219, Bordeaux Population Health Center, University of Bordeaux,
111 33076 Bordeaux Cedex, France
112 ⁴² Department of General and Interventional Cardiology, University Heart Center
113 Hamburg, Hamburg, Germany
114 ⁴³ Department of Neurology, Brain Center Rudolf Magnus, University Medical Center
115 Utrecht, Utrecht, 3584 CG, The Netherlands
116 ⁴⁴ Department of Biological Psychology, Amsterdam Public Health Research Institute,
117 Vrije Universiteit Amsterdam, Van Der Boechorststraat 7-9, 1081 BT, Amsterdam, The
118 Netherlands
119 ⁴⁵ University of Groningen, University Medical Center Groningen, Department of
120 Pediatric Pulmonology and Pediatric Allergology, Beatrix Children's Hospital, GRIAC
121 Research Institute Groningen, The Netherlands
122 ⁴⁶ Chromatin and Disease Group, Cancer Epigenetics and Biology Programme (PEBC),
123 Bellvitge Biomedical Research Institute (IDIBELL), 08908 L'Hospitalet de Llobregat,
124 Barcelona, Spain
125 ⁴⁷ Instituto de Parasitología y Biomedicina López Neyra, CSIC, Granada, Spain
126 ⁴⁸ School of Medicine and Public Health, Faculty of Medicine and Health, The University
127 of Newcastle, Newcastle, NSW, Australia

128 ⁴⁹ The Curtin/UWA Centre for Genetic Origins of Health and Disease, Faculty of Health
129 Sciences, School of Biomedical Sciences, Curtin University and School of Biomedical
130 Sciences, Faculty of Health and Medical Sciences, The University of Western Australia,
131 Perth Australia
132 ⁵⁰ The Hospital for Sick Children, University of Toronto, Toronto, Canada M5G 1X8
133 ⁵¹ Rotman Research Institute, University of Toronto, Toronto, Canada M6A 2E1
134 ⁵² University of Essex, Wivenhoe Park, Colchester, Essex, CO4 3SQ
135 ⁵³ Department of Neurology, King's College Hospital, London, UK
136 ⁵⁴ United Kingdom Dementia Research Institute, King's College London, London, UK
137 ⁵⁵ Department of Clinical Science and Education, Södersjukhuset, Karolinska Institutet,
138 Stockholm, Sweden
139 ⁵⁶ QIMR Berghofer Medical Research Institute, Brisbane, Australia
140 ⁵⁷ Department of Psychiatry and Behavioral Sciences, Duke University Medical School,
141 Durham, NC, USA
142 ⁵⁸ Dunedin Multidisciplinary Health and Development Research Unit, Department of
143 Psychology, University of Otago, Dunedin, New Zealand
144 ⁵⁹ Program in Medical and Population Genetics, Broad Institute, Broad Institute,
145 Cambridge, MA, USA
146 ⁶⁰ MRC Social, Genetic and Developmental Psychiatry Centre, Institute of Psychiatry,
147 Psychology and Neuroscience, King's College London, London, UK
148 ⁶¹ Piemonte Centre for Cancer Prevention, Turin, Italy
149 ⁶² Dipartimento Di Medicina Clinica E Chirurgia, Federico II University, Naples, Italy
150 ⁶³ The Generation R Study Group, Erasmus MC, University Medical Center Rotterdam,
151 Rotterdam, The Netherlands
152 ⁶⁴ Department of Epidemiology, Erasmus MC, University Medical Center Rotterdam,
153 Rotterdam, The Netherlands
154 ⁶⁵ Department of Pediatrics, Erasmus MC, University Medical Center Rotterdam,
155 Rotterdam, The Netherlands
156 ⁶⁶ Department of Child and Adolescent Psychiatry, Erasmus Medical Center, Rotterdam,
157 Netherlands
158 ⁶⁷ Department of Social and Behavioral Science, Harvard TH Chan School of Public
159 Health, Boston, USA
160 ⁶⁸ School of Clinical Medicine, University of Cambridge, UK;
161 ⁶⁹ Department of Psychology, Education and Child Studies, Erasmus University
162 Rotterdam, Rotterdam, The Netherlands
163 ⁷⁰ Institute of Medical Sciences, University of Aberdeen, Aberdeen, UK
164 ⁷¹ Department of Psychology and Logopedics, Faculty of Medicine, University of
165 Helsinki, Finland
166 ⁷² Research Unit for Gynaecology and Obstetrics, Institute of Clinical research,
167 University of Southern Denmark, Odense, Denmark
168 ⁷³ Centre of Women's, Family and Child Health, University of South-Eastern Norway,
169 Kongsberg, Norway
170 ⁷⁴ The Novo Nordisk Foundation Center for Basic Metabolic Research, Faculty of Health
171 and Medical Sciences, University of Copenhagen, Denmark

172 ⁷⁵ Department of Public Health (Section of Epidemiology), Faculty of Health and Medical
173 Sciences, University of Copenhagen, Copenhagen, Denmark.
174 ⁷⁶ The National Institute of Public Health, University of Southern Denmark, Copenhagen
175 ⁷⁷ Department of Psychiatry and Behavioral Sciences, Emory University School of
176 Medicine, Atlanta, GA, USA
177 ⁷⁸ Center for Genomic Regulation (CRG), Barcelona Institute of Science and
178 Technology, Barcelona, Spain
179 ⁷⁹ IMIM (Hospital del Mar Medical Research Institute), Barcelona, Spain
180 ⁸⁰ Clinical and Experimental Sciences, Faculty of Medicine, University of Southampton,
181 Southampton, UK
182 ⁸¹ Division of Epidemiology, Biostatistics, and Environmental Health Sciences, School of
183 Public Health, University of Memphis, Memphis, USA
184 ⁸² Queensland Brain Institute, University of Queensland, Australia
185 ⁸³ Alzheimer Scotland Dementia Research Centre, University of Edinburgh, University of
186 Edinburgh, UK
187 ⁸⁴ Department of Gerontology and Geriatrics, Leiden University Medical Center, Leiden,
188 The Netherlands
189 ⁸⁵ C2VN, Aix-Marseille University, INSERM, INRAE, Marseille, France
190 ⁸⁶ Avera Institute for Human Genetics, Sioux Falls, USA
191 ⁸⁷ University of Groningen, University Medical Center Groningen, Department of
192 Epidemiology, GRIAC Research Institute Groningen, Groningen, The Netherlands
193 ⁸⁸ Institute for Risk Assessment Sciences, Universiteit Utrecht, Utrecht, The Netherlands
194 ⁸⁹ Julius Center for Health Sciences and Primary Care, University Medical Center
195 Utrecht, Utrecht, The Netherlands
196 ⁹⁰ Telethon Kids Institute, University of Western Australia, Perth, WA, Australia
197 ⁹¹ UCL Institute of Cardiovascular Science, London, UK
198 ⁹² Bloorview Research Institute, Holland Bloorview Kids Rehabilitation Hospital and
199 Departments of Psychology and Psychiatry, University of Toronto, Toronto, Canada,
200 M4G 1R8
201

202 Contributions

203 **Project management:** G.H., G.S., J.L.M

204 **Designed individual studies and contributed data:**

205 A.A.C., A.Cas., A.D.H., A.G.U, A.Me., A.Mu., A.M.M., B.B., B.T.H.,
206 C.H., C.L.R., C.P., C.Sa., C.Sh., C.Sö., D.A.L., D.v.H., D.I.B., D.T., E.A.N., E.B.B.,
207 E.J.C.d.G, E.M., F.G., F.R., G.E.D, G.H.K., G.P., G.W.M., H.R.E., H.T., H.Z., I.J.D.,
208 J.F.F., J.H.V., J.J.C., J.Ka., J.L., J.M., J.M.S., J.M.V., J.v.M., J.R., J.R.B.P., J.R.G.,
209 J.Sh., J.T.B., J.W., J.W.H., K.K.O., K.L.E., K.R., L.A., L.C.S., L.M., M.A.I., M.Bee.,
210 M.Bu., M.E.A.R., M.H.v.IJ., M.Ke., M.O., N.C., N.G.M., N.J.W., N.R.W., P.E.S., P.Mo.,
211 P.M.V., R.H., R.P., S.L., S.P., T.D.S., T.E., T.E.M., T.I.A.S, T.P., T.T., V.W.V.J., W.K.,
212 Z.P.

213 **Generated and/or quality-controlled data:** A.A.K., A.I., A.S., C.S.M., H.R.E., J.L.M.,
214 K.B., K.M.H., N.K., S.M.R., T.H., R.M.W., W.L.M.
215 **Designed new statistical or bioinformatics tools:** G.H., J.L.M., M.Su., T.R.G., V.I.
216 **Analysed the data and/or provided critical interpretation of results:**
217 A.D.B, A.Car., A.D., A.F.M., A.K., B.T.H., C.B., C.H., C.L.R., C.R.A., C.Sor., C.V., C.X.,
218 C.W., D.A., D.C., D.J.L., D.L.C., D.M., E.C.M., E.G., E.H., E.M., F.C.M., F.I.R., F.R.D.,
219 G.B., G.C., G.D.S., G.H., G.H.K., G.M., G.W., I.Y., J.C.F., J.v.D., J.J.H., J.Ka., J.Kl.,
220 J.L.M., J.M., J.Su., J.T.B., K.B., K.v.E., K.F.D., K.S., L.C.S., M.Ber., M.Bu., M.H.v.IJ.,
221 M.G., M.Ku., M.L., M.Sm., M.Su., N.K., P.Me., P.Ma., P.M.V., R.E.M., R.G., R.L., R.Z.,
222 S.B., S.G., S.K., T.C., T.G., T.G.R., T.I.A.S., T.L., T.R.G., Y.A., Y.Z., V.I., V.S.
223 **Designed and/or managed the study:** B.T.H., C.B., C.L.R., J.M., J.T.B., T.R.G.
224 **Wrote the manuscript:** A.D.B., B.T.H., C.B., C.L.R., D.J.L., E.C.M, E.H., G.D.S., G.H.,
225 J.C.F., J.Kl., J.L.M., J.M., J.T.B., K.B., K.F.D., M.Su., P.M.V., R.L., T.G.R., T.R.G., V.I.

226 Competing interests

227 The authors declare no competing interests.

228 Financial disclosures

229 None of the authors have financial disclosures.

230

231 Abstract

232 Characterising genetic influences on DNA methylation (DNAm) provides an opportunity
233 to understand mechanisms underpinning gene regulation and disease. Here we
234 describe results of DNA methylation-quantitative trait loci (mQTL) analyses on 32,851
235 participants, identifying genetic variants associated with DNAm at 420,509 DNAm sites
236 in blood. We present a database of >270,000 independent mQTL of which 8.5%
237 comprise long-range (*trans*) associations. Identified mQTL associations explain 15-17%
238 of the additive genetic variance of DNAm. We reveal that the genetic architecture of
239 DNAm levels is highly polygenic and DNAm exhibits signatures of negative and positive
240 natural selection. Using shared genetic control between distal DNAm sites we construct
241 networks, identifying 405 discrete genomic communities enriched for genomic
242 annotations and complex traits. Shared genetic factors are associated with both DNAm
243 levels and complex diseases but only in a minority of cases these associations reflect
244 direct causal relationships from DNAm to trait or vice versa indicating a more complex
245 genotype-phenotype map than previously anticipated.

246 Main

247 The role of common inter-individual variation in DNA methylation (DNAm) on disease
248 mechanisms is not yet well characterised. It has, however, been hypothesised to serve
249 as a viable biomarker for risk stratification, early disease detection and the prediction of
250 disease prognosis and progression.¹ Because genetic influences on DNAm in blood
251 have been shown to be widespread²⁻⁴, a powerful avenue into researching the
252 functional consequences of changes in DNAm levels is to map genetic differences
253 associated with population-level variation, identifying DNA methylation quantitative trait
254 loci, (mQTL) that include both local (*cis* mQTL) and distal (*trans* mQTL) effects. We can
255 harness mQTL as natural experiments, allowing us to observe randomly perturbed
256 DNAm levels in a manner that is not confounded with environmental factors^{5,6}. In this
257 regard, mapping even very small genetic effects on DNAm is valuable for gaining power
258 to evaluate whether its variation has a substantial causal role in disease and other
259 biological processes.

260
261 To date, only a small fraction of the total genetic variation estimated to influence DNAm
262 across the genome has been identified⁷, and the proportion of *trans* heritability
263 explained by *trans* mQTL (defined as more than 1Mb from the DNAm site) is much
264 smaller than the proportion of *cis* heritability explained by *cis* mQTL. Therefore, the
265 majority of genetic effects are likely to act in *trans*, have small effect sizes^{5,7-9}, while
266 being potentially more informative in the biological insights they provide.^{8,10} Much larger
267 sample sizes are required to map associations involving small genetic effects in order to
268 permit greater understanding of the genetic architecture and the biological processes
269 underlying DNAm⁷. To this end, we established the Genetics of DNA Methylation
270 Consortium (GoDMC), an international collaboration of human epidemiological studies
271 that comprises >30,000 study participants with genetic, phenotypic and DNAm data.

272
273 Importantly, the unrivalled sample size and coverage of our study enables us to identify
274 a large number of *cis* and *trans* mQTL to gain biological insights that were previously
275 impossible. First, we use this extensive resource to uncover the genetic architecture of
276 DNAm and to study natural selection pressures. Second, we learn about how *cis*- and
277 *trans*-acting variants and DNAm sites interact through the development of new network
278 approaches. Third, we interrogate the potential role of DNAm in disease mechanisms by
279 exhaustively mapping the causal relationships of DNAm with 116 complex traits and
280 diseases in a bi-directional manner. A database of our results is available as a resource
281 to the community at <http://mqtlldb.godmc.org.uk/>.

282 Genetic variants influence 45% of tested DNAm sites

283 In order to map genetic influences on DNAm, we established an analysis workflow that
284 enabled standardized meta-analysis and data integration across 36 population-based
285 and disease datasets with genotype and DNAm data. Using a two-phase discovery
286 study design, we analyzed ~10 million genotypes imputed to the 1000 Genomes

287 reference panel¹¹ and 420,509 DNAm sites measured by Infinium HumanMethylation
288 BeadChips in whole blood derived from 27,750 European participants (**Figures 1A** and
289 **S1-S5, Table S1-S2, Supplementary Note 1, Supplementary Information**).

290
291 Using linkage disequilibrium (LD) clumping, we identified 248,607 independent *cis*-
292 mQTL associations ($p < 1e-8$, $< 1\text{Mb}$ from the DNAm site, **Figure S4**) with a median
293 distance between single nucleotide polymorphisms (SNP) and DNAm sites of 36kb
294 (IQR=118 kb, **Figure S3A**). We found 23,117 independent *trans* mQTL associations
295 (using a conservative threshold of $p < 1e-14^7$, **Figure S4, Supplementary**
296 **Information**). These mQTL involved 190,102 DNAm sites, representing 45.2% of all
297 those tested (**Figure 1B**) which is a 1.9x increase of sites with a *cis* association ($p < 1e-$
298 8) and 10x increase of sites with a *trans* association ($p < 1e-14$) over a previous study
299 whose sample size was 7x smaller⁸. As expected, mQTL effect sizes for each DNAm
300 site (the maximum absolute additive change in DNAm level measured in standard
301 deviation (SD) per allele) were lower for sites with a *trans* association (as compared to
302 sites with a *cis* association (per allele SD change = -0.02 (s.e.=0.002, $p = 2.1e-14$,
303 **Figure S6**). The differential improvement in yield between *cis* and *trans* associations is
304 revealing in terms of the genetic architecture – relatively small sample sizes are
305 sufficient to uncover the majority of large *cis* effects, whereas much larger sample sizes
306 are required to identify the polygenic *trans* component.

307
308 The majority of *trans* associations (80%) were inter-chromosomal. Of the intra-
309 chromosomal *trans* associations, 34% were $>5\text{ Mb}$ from the DNAm site, **Figure S7**). We
310 then compared the rate of inter-chromosomal *trans* associations to the rate of intra-
311 chromosomal *trans* associations (excluding chromosome 6) and found a substantially
312 lower number of inter-chromosomal *trans* associations per 5 Mb region (1.59) than intra-
313 chromosomal associations ($>1\text{ Mb}$: 7.95; $>6\text{ Mb}$ 4.81).

314
315 Next, using conditional analysis¹² we explored the potential for multiple independent
316 SNPs operating within the locus of each mQTL, identifying 758,130 putative
317 independent variants. Each DNAm site, for which a mQTL in *cis* had been detected, had
318 a median of 2 independent variants (IQR=4 variants, **Figure S8**). For all subsequent
319 analyses, we used index SNPs from clumping procedures to be conservative and
320 unbiased due to the non-independence of genetic variants.

321
322 The microarray technology used in the majority of cohorts limited us to analyse $<2\%$ of
323 sites across the genome¹³, which are biased to promoters and strongly
324 underrepresented regulatory elements. To explore the impact of expanding the
325 coverage of arrays, we calculated the linear relationship between the median number of
326 probes by gene on the 450k array and the median number of *cis* and *trans* mQTL. For
327 each probe, we found an increase of 0.76 *cis* mQTL ($p < 9.03e-16$) and 0.05 *trans* mQTL
328 ($p < 1.47e-05$) (**Figure S9**). A similar increase was seen in non-genic regions. This
329 indicates that expanding coverage will increase mQTL yield although this will depend on
330 the genetic contribution of the DNAm site and cell type specificity.

331
332 We sought to replicate the mQTL using the Generation Scotland (GS) cohort (n = 5,101)
333 for which mQTL results were previously generated using an independent analysis
334 pipeline (**Supplementary Information, Supplementary Note 1**). Data were available to
335 allow us to test for replication of 188,017 of our discovery mQTL (137,709 sites) and we
336 found a very strong correlation of effect sizes for both *cis* and *trans* effects ($r=0.97$,
337 $n=155,191$ and 0.96 , $n=14,465$ at $p<1e-3$, respectively; **Figure 1C**); 99.6% of the
338 associations had a consistent sign (further discussion in **Supplementary Information**).
339 At an approximate Bonferroni corrected threshold of $0.05/188,017$, 142,727 of the
340 discovery mQTL replicated in the GS cohort (76%); the replication rate for *cis* and *trans*
341 mQTL were 76% and 79%, respectively. To evaluate whether our replication rate was in
342 line with expectations given the smaller replication sample size, we estimated that under
343 the assumption that the discovery mQTL are true positives 171,824 mQTL would be
344 expected to replicate at a nominal threshold of $1e-3$. In very close agreement we found
345 that the actual number of mQTL replicating at this level was 169,656, indicating that the
346 majority of our discovery mQTL are likely to be true positives (**Table S3**,
347 **Supplementary Information**). Our findings support that there is little between-study
348 heterogeneity in our analysis and that genetic effects on DNAm are highly stable across
349 cohorts (**Figure S2, Table S2**).

350
351 Overall, the variance explained by replicated genetic effects was small. For 99% of the
352 associations in *cis* and *trans*, mQTL explained less than 21% and 16% of the DNAm
353 variation respectively (**Figure S10**). Aggregating across all 420,509 tested DNAm sites,
354 our replicated mQTL associations explain 1.3% of the total assayed DNAm variation,
355 8% of this being due to *trans*-associations. Restricting to sites that have at least one *cis*-
356 effect or *trans*-effect, however, we explain 4.2% and 2.5% of the DNAm variance,
357 respectively.

358
359 We then investigated how much of the heritability of variable DNAm can be explained
360 by our mQTL associations on the 450k array using family-based heritability studies of
361 DNAm^{2,14}. We found a strong positive relationship between variance explained by
362 replication mQTL estimates (127,680 sites in GS) and heritability for both studies
363 (family: $r=0.41$ across, 121,582 available sites; twin: $r=0.37$ across 118,955 available
364 sites) (**Figure 1D, Table S4**). The mQTL that we identified explain 15%-17% of the
365 additive genetic variance of DNAm (**Figure S11**). Finally, there were strong positive
366 relationships between the heritability of DNAm levels at a DNAm site and the number of
367 independent mQTL (**Figure S12**), heritability and effect size (**Figure S13**), variance
368 explained and the number of independent mQTL (**Figure S14**) and variance explained
369 and distribution of DNAm levels (**Figure S15**). Overall, our results support a mixed
370 genetic architecture of polygenic genome-wide effects and larger *cis* effects.

371
372 The coverage of the mQTL search in this study was limited by the computational
373 necessity of a multiple stage study design (**Figure S16**). Those mQTL that we
374 discovered with r^2 less than 1% are likely a small fraction of all the mQTL in this

375 category expected to exist (**Figure S17**). Across these DNAm sites, and within the
376 range of mQTL detected in our study ($r^2 > 0.22\%$) we estimate that there are twice as
377 many *cis* mQTL and 22.5 times more *trans* mQTL yet to discover (**Figure S17**). This
378 would likely not explain all estimated heritability, indicating that a substantial set of the
379 heritability is due to causal variants with smaller effects than those detectable given our
380 study size or due to rare variants not represented in our imputed genotype data.
381

382 *Cis* and *trans* mQTL operate through distinct mechanisms

383 We analysed how inter-individual DNAm changes are associated to genetic variation in
384 a context way which has so far mainly focused on *cis* mQTL^{7,8,15-17}. The statistical power
385 of the mQTL analysis allowed us to identify SNPs only associated with DNAm in *cis*
386 ($n=157,095$, 69.9%), only associated with DNAm in *trans* ($n=794$, 0.35%), or associated
387 with DNAm in both *cis* and *trans* ($n=66,759$, 29.7%). Similarly, of the 190,102 DNAm
388 sites influenced by a SNP, 170,986 DNAm sites (89.9%) were *cis-only*, 11,902 DNAm
389 sites (6.3%) were *cis+trans*, and 7,214 DNAm sites (3.8%) were *trans-only*. This
390 categorisation allowed us to infer biological properties of *trans*-features that were not
391 due to their *cis*-effects.
392

393 Here, we first compared the distribution of DNAm levels (weighted mean DNAm level
394 across 36 studies (defined as low (<20%), intermediate (20%-80%) or high (>80%)
395 between the *cis* and *trans* DNAm sites (**Figure 1B**). We then performed enrichment
396 analyses on the mQTL SNPs and DNAm sites using 25 combinatorial chromatin states
397 from 127 cell types (including 27 blood cell types)¹⁸ and gene annotations (**Figure 2A**,
398 **S18-S21, Tables S5-S8**). Consistent with previous studies^{7,8,17}, we found that *cis only*
399 sites are represented in high (32%), low (28%) and intermediate (40%) DNAm levels
400 and these sites are mainly enriched for enhancer chromatin states (mean OR=1.37),
401 CpG islands (OR=1.25) and shores (OR=1.26).

402 For *cis+trans* sites, we found that the majority of these sites (66%) have intermediate
403 DNAm levels. By replicating this finding in two isolated white-blood-cell subsets (**Figure**
404 **S22**), we showed that this is due to cell-to-cell variability^{18,19} or sub cell type differences
405 which may indicate that these loci contribute to the divergence into further sub cell
406 types. In line with the observation that intermediate levels of DNAm are found at distal
407 regulatory sequences^{20,21}, these sites were enriched for enhancer (mean OR=1.65) and
408 promoter states (mean OR=1.41). However for *trans only* sites, we found a pattern of
409 low DNAm (for 55% of sites) and enrichments for promoter states (mean OR=1.39)
410 especially TssA promoter state (mean OR=2.03). We demonstrated that these
411 inferences about *cis* and *trans* enrichments were not sensitive to the definition of *trans*
412 associations, by showing that the patterns were consistent if we restricted to only inter-
413 chromosomal associations (**Supplemental Information, Figure S23**).
414

415 We continued by analysing the differences in properties between SNPs that have local
416 versus long-range DNAm influences. We found that *cis only* and *cis+trans* SNPs were

417 enriched for active chromatin states and genic regions whereas *trans only* SNPs were
418 enriched for intergenic regions and the heterochromatin state (**Figure 2A, S20-S21,**
419 **Tables S7-S8**). Our analysis shows that *trans-only* sites and SNPs have different
420 properties as *cis+trans* SNPs and sites, indicating that enrichments of general *trans*
421 categories are dominated by their *cis* functionality. Overall, these results highlight that a
422 complex relationship between molecular features is underlying the mQTL categories
423 and the biological contexts are substantially different between *cis* and *trans* features.
424

425 We found that these inferences were often shared across other tissues. For example,
426 DNAm sites with low or intermediate DNAm levels have similar DNAm distributions in
427 12 tissues (**Figure S24-26**). However, while SNP and DNAm site enrichments were
428 typically present in multiple tissues, enrichments were stronger in blood datasets for the
429 enhancer states (SNP: difference in mean OR=0.055, $p=0.038$; sites: difference in
430 mean OR=0.21, $p < 2e-16$) and DNase state (SNP: difference in mean OR=0.13,
431 $p=0.004$; sites: difference in mean OR=0.41 $p=9.65e-16$) indicating some level of tissue
432 specificity for mQTL in these regions (**Figure S18, S20, S27**).
433

434 To investigate the question of tissue specificity further, we compared the correlation of
435 effect estimates of *cis* and *trans* mQTL in blood against adipose tissue ($n=603$)²² and
436 brain ($n=170$)⁹ (**Supplementary Information, Table S9**). We found a larger extent of
437 QTL sharing of blood and adipose tissue as compared to blood and brain which might
438 be explained by shared cell types in line with *cis* eQTL findings²³. Generally, the
439 between tissue effect correlations were high, in line with a recent comparison of *cis*-
440 mQTL effects between brain and blood²⁴. However, we found that the highest
441 correlations were for associations involving *trans-only* sites (Adipose $r_b=0.92$ (se
442 $=0.004$); Brain $r_b=0.88$ (se=0.009)) despite having on average smaller effect sizes than
443 *cis only* associations, implying that they are *less* tissue specific than *cis* effects (Adipose
444 $r_b=0.73$ (se $=0.002$); Brain $r_b=0.59$ (se=0.004)) which is line with the notion that
445 promoters are less tissue-specific. Stratifying the mQTL categories to low, intermediate
446 and high DNAm, showed that the brain-blood correlations are the lowest for
447 intermediate DNAm categories and adipose-blood correlations are lowest for high
448 DNAm categories, which may suggest cellular heterogeneity for high DNAm levels
449 (**Table S9**). These results show the value of large sample sizes in blood to detect *trans*
450 mQTL regardless of the tissue.

451 *Trans* mQTL SNPs and DNAm exhibit patterned TF binding

452 Recent studies have uncovered multiple types of transcription factor (TFs)/DNA
453 interactions with DNAm including the binding of DNAm-sensitive TFs²⁵⁻²⁷. Epigenetic
454 editing studies have revealed that local methylation and demethylation activities are
455 affected by TF binding and cooperativity between TFs^{26,28}. To gain insights into how
456 SNPs induce long-range DNAm changes, we mapped enrichments for DNAm sites and
457 SNPs across binding sites for 171 TFs in 27 cell types^{29,30}. We found strong
458 enrichments for the majority of TFs amongst DNAm sites with a *trans* association

459 (*cis+trans*: 55%; *trans only*: 80%; *cis only*: 18%) which is in line with the observation that
460 loss of DNAm at promoters is usually associated with gene activation³¹, and amongst
461 *cis-acting* SNPs (*cis only*: 96%, *cis+trans*: 91%, *trans only*: 1%) (**Figures 2B, S28, S29**).
462 Consistent with the observation that *trans only* DNAm sites are enriched for CpG
463 islands (**Figure S19**), sites that overlap TFBS were relatively hypomethylated
464 independent of tissue (weighted mean DNAm levels = 21% vs 52%, $p < 2.2e-16$) (**Figure**
465 **S30**) and we found that generally the TFBS enrichments were not tissue specific (**Table**
466 **S10-11, Figure S28-29**).

467 Next we investigated a possible mechanism that may be responsible for these *trans*
468 mQTL. We hypothesized that if a *trans* mQTL is driven by TF activity^{8,10} then particular
469 TF-TF pairs may exhibit preferential enrichment³². A mQTL has a pair of TFBS
470 annotations³⁰, one for the SNP and one for the DNAm site. Using a novel approach
471 (two-dimensional functional enrichment, **Figure S31**), we evaluated if the annotation
472 pairs amongst 18,584 inter-chromosomal *trans*-mQTL were associated to TF binding in
473 a non-random pattern (**Supplementary Information**). We found that 6.1% (22,962 of
474 378,225) of possible pairwise combinations of SNP-DNAm site annotations were more
475 over- or under-represented than expected by chance after strict multiple testing
476 correction (**Supplementary Information, Table S12, Figure 2C-D**).

477
478 After accounting for abundance and other characteristics, the strongest pairwise
479 enrichments involved sites close to TFBS for proteins in the cohesin complex, for
480 example CTCF, SMC3 and RAD21, as well as TFs such as GATA2 related to cohesin³³.
481 Bipartite analysis showed that these clustered due to being related to similar sets of
482 SNP annotations (**Figure 2C**). Other clusters were also found, for example, sites close
483 to TFBS for interferon regulatory factor 1 (*IRF1*), a gene for which *trans*-acting
484 regulatory networks³⁴, and enrichment amongst causally interacting caQTL³⁵ have been
485 previously reported were more likely to be influenced by SNPs near TFBS for EZH2,
486 SMC3, ATF3, BCL3, TR4 and MAX. The relationship between IRF1 and these other
487 proteins has been documented previously³⁶⁻³⁸. For example EZH2 mediates the
488 silencing of IRF1³⁹; BCL3 and IRF1 are co-down-regulated during inflammation³⁶; and
489 ATF3 is a negative regulator of cytokines which themselves induce IRF1^{37,38}.

490
491 Previous studies have indicated chromosomal interactions (genomic regions that have
492 been shown to spatially colocalise within the cell⁴⁰) as alternative mechanism for *trans*
493 coordination^{8,41}. We compared the locations of inter-chromosomal *trans* mQTL
494 ($n=18,584$) to known regions of chromatin interactions. We found 1175 overlaps for 637
495 SNP-DNAm site pairs (3.4%) where the LD region of the mQTL SNP and the
496 corresponding site overlapped with any interacting regions (525 SNPs, 602 sites) as
497 compared to a mean of 473 SNP-DNAm site pairs in 1000 permuted datasets
498 ($OR=1.36$, $p_{Fisher}=6.5e-7$, $p_{empirical} < 1e-3$) (**Figure S32**). To summarise, our results show
499 that *trans* mQTL are in part driven by long-range cooperative TF interactions and, that
500 for a small proportion of interchromosomal *trans* mQTL the spatial distance *in vivo* is
501 likely to be small.

502 Communities of DNAm sites are identified by shared *trans*-genetic 503 effects

504 Genetic variation can perturb chromatin activity^{32,35,41}, DNAm⁸ or gene expression⁴²
505 across multiple sites in *cis* and *trans* revealing coordinated activity between regulatory
506 elements and genes. *Trans*-mQTL provide an opportunity to infer how distal genomic
507 regions are functionally related, but the polygenic nature of DNAm variation could lead
508 to apparent shared genetic effects that arise from distinct causal variants rather than
509 shared genetic factors. We observed that there were 1,728,873 instances where a SNP
510 acting in *trans* also influenced a *cis* DNAm site (before LD pruning). Genetic
511 colocalization analysis indicated that 278,051 of these instances were due to the *cis* and
512 *trans* sites sharing a genetic factor, representing 3,573 independent *cis-trans* genomic
513 region pairs, of which 3,270 were inter-chromosomal (**Table S13**, see **Supplementary**
514 **Information** for sensitivity analysis for the colocalization method used in the context of
515 the two-stage mQTL discovery design). These pairs consisted of 1,755 independent
516 SNPs and 5,109 independent DNAm sites across the genome, indicating that some
517 sites with *cis* associations shared genetic factors with multiple sites with *trans*
518 associations revealing distal coordination between mQTL. From the *cis-trans* pairs we
519 constructed a network linking these genomic regions which elucidated 405
520 “communities” of genomic regions that were substantially connected (**Supplementary**
521 **Information**). Fifty-six of these communities comprised 10 or more sites, and the
522 largest community comprised 253 sites (**Figure 3A**).

523
524 We hypothesised that *cis* sites were causally influencing multiple *trans* sites within their
525 communities (i.e. a causal chain of mQTL to DNAm at a *cis* site to DNAm at a *trans*
526 site). We evaluated whether the estimated causal effect (obtained from the *trans*-mQTL
527 effect divided by the *cis*-mQTL effect i.e. the Wald ratio) of the *cis* site on the *trans* site
528 was consistent with the observational correlation between the *cis*- and *trans*-site. While
529 there was an association, the relationship was weak ($r=0.096$, $p=1.73e-6$, **Figure S33**),
530 indicating that changes in *cis* sites causing changes in *trans* sites is likely not the
531 predominant mechanism. We did observe that the *cis-trans* DNAm levels were more
532 strongly correlated than we would expect by chance (**Figure S34**), which supports the
533 notion that they are jointly regulated without generally being causally related.

534
535 To gain functional insights into these communities, we evaluated if DNAm sites within
536 each community were enriched for regulatory annotations and/or gene ontologies
537 (**Table S14-S17**, **Figure S35-36**). Multiple communities showed enrichments (FDR P
538 <0.001); for example community 9 DNAm sites were strongly enriched for TFBS
539 annotations relating to the cohesin complex in multiple cell types, community 22 DNAm
540 sites were enriched for NFKB and EBF1 in B lymphocytes and community 76 DNAm
541 sites were enriched for EZH2 and SUZ12 and bivalent promotor and repressed
542 polycomb states (**Figure 3B**). Community 2 (comprising 253 sites) was enriched for
543 active enhancer state in 3 cell types and for lymphocyte activation (GO:0046649 FDR p

544 = 0.016) and multiple KEGG pathways including the JAK-STAT signalling pathway
545 (I04630: FDR $p=8.53e-7$) (**Table S16, Table S17**).

546
547 Regulatory features within a network may share a set of biological features that are
548 related to complex traits. We performed enrichment analysis to evaluate if the loci
549 tagged by DNAm sites in a community were related to each of 133 complex traits
550 (**Table S18**), accounting for non-random genomic properties of the selected loci.
551 Restricting the analysis to only the 56 communities with ten or more sites, we found
552 eleven communities that tagged genomic loci that were enriched for small p-values with
553 22 complex traits (FDR < 0.05) (**Figure 3C, Table S19**). Blood related phenotypes were
554 overrepresented (11 out of 23 enrichments being related to metal levels or
555 haematological measures, binomial test p-value = $4.2e-5$). Amongst the communities
556 enriched for GWAS signals, community 16 was highly associated with iron and
557 haemoglobin traits. Community 9 was associated to plasma cortisol ($p = 8.27e-5$).
558 Finally, we performed enrichment analysis on 36 blood cell count traits⁴³ and found
559 enrichments for two communities. Community 16 was enriched for hematocrit ($p=4.34e-$
560 10) and hemoglobin concentration ($p=1.99e-8$) and community 5 was enriched for
561 reticulocyte traits ($p=1.67e-6$) (**Figure S37**). The enrichments found for these DNAm
562 communities indicate that a potentially valuable utility of mapping *trans*-mQTL is to
563 indicate how distal regions of the genome are functionally related.

564 mQTL can be used to identify shared genetic influences with 565 disease

566 The majority of GWA loci map to non-coding regions⁴⁴ and *cis* mQTL are enriched
567 amongst GWA^{16,45,46}. Here we investigated the value of the large number of mQTL
568 especially *trans* mQTL to annotate functional consequences of GWA loci. We first
569 tested genome-wide enrichment of GWAS associations (SNPs at $p < 5e-8$ for a given
570 complex trait) amongst mQTL SNPs, performing separate analysis for mQTL acting in
571 *cis*, *cis* and *trans* and *trans*. We utilized genome-wide summary statistics for 37
572 phenotypes related to 11 disease/trait categories with 41 publicly available GWAS
573 datasets (**Table S20**). After accounting for non-random genomic distribution of mQTL⁴⁷
574 and multiple testing, we identified enrichments for 35% of the complex traits (**Figure**
575 **S38, Table S20, Supplementary Information**) mainly for studies with a larger number
576 of GWA signals. The *cis+trans* mQTL were most strongly enriched for low p-values
577 across multiple traits. Six phenotypes across 4 disease categories were associated with
578 *cis* mQTL, nine phenotypes across 5 disease categories were associated with *cis+trans*
579 mQTL. Inflammatory bowel disease and Crohn's disease were associated with both
580 sets. Height was associated across all three categories of mQTL but interestingly was
581 depleted for mQTL in the *trans only* group (OR=0.354, $p=7.31e-8$). The distribution of
582 enrichment effect estimates (ORs) of *trans* mQTL was substantially closer to the null or
583 in depletion when compared to mQTL that included *cis* effects (**Figure 2E**). These
584 enrichments correspond to the results reported earlier, in which *trans*-SNPs were

585 typically depleted for enhancer and promoter regions, whereas complex trait loci are
586 enriched for coding and regulatory regions⁴⁸.

587
588 Though the mQTL discovery pipeline adjusted for predicted cell types^{49,50} and non-
589 genetic DNAm PCs, there is a possibility that residual cell-type heterogeneity remains.
590 We performed another set of GWAS enrichment analysis, this time using 36 blood cell
591 traits⁴³, and found enrichments. These were strongest amongst *cis+trans* mQTL, as
592 seen in the previous enrichments (**Figure S39**). Interrogating this further, we found that
593 for 98.9-100% of the mQTL, mQTL SNPs explained more variation in DNAm than they
594 explain variation in blood cell counts suggesting a causal chain of mQTL to blood trait⁵¹.
595 Alternatively, a systematic measurement error difference could explain these
596 observations, where DNAm captures blood cell counts more accurately than
597 conventional measures.

598
599 The enrichments suggest that overlaps are not due to chance which motivated us to a
600 much more in-depth analysis on a much larger number of traits/diseases. We searched
601 for instances of DNAm sites sharing the same genetic factors against each of 116
602 complex traits and diseases, and initially found 23,139 instances of an mQTL strongly
603 associating with a complex trait (**Figure 4**). To evaluate the extent to which these were
604 due to shared genetic factors (and not, for example, LD between independent causal
605 variants), we performed genetic colocalization analysis⁵² (**Table S18, Table S21**).
606 Excluding genetic variants in the *MHC* region, we found 1,373 putative examples in
607 which at least one DNAm site putatively shared a genetic factor with at least one of 71
608 traits (including 19 diseases). Those DNAm sites that had a shared genetic factor with a
609 trait were 6.9 times more likely to be present in a community compared to any other
610 DNAm site with a known mQTL (Fisher's exact test 95% CI 4.8-9.7, $p = 9.2e-19$). Next,
611 we evaluated how often the DNAm site that colocalised with a known GWAS hit was the
612 closest DNAm site to the lead GWAS variant by physical distance. Notably, in only
613 18.1% of the cases where a GWAS signal and an assayed 450k DNAm site colocalised,
614 was that DNAm site the closest DNAm site to the signal. This finding is similar to results
615 found for gene expression⁵³, but the converse has been found for protein levels⁵⁴.

616
617 It has previously been difficult to conclude whether genetic colocalisation between
618 DNAm and complex traits indicates a) a causal relationship where the DNAm level is on
619 the pathway from genetic variant to trait (vertical pleiotropy) or b) a non-causal
620 relationship where the variant influences the trait and DNAm independently through
621 different pathways (horizontal pleiotropy)⁵⁵. In Mendelian randomisation (MR) it is
622 reasoned that under a causal model, multiple independent genetic variants influencing
623 DNAm should exhibit consistent causal effects on the complex trait⁵⁶. Amongst the
624 putative colocalising signals, 440 (32%) involved a DNAm site that had at least one
625 other independent mQTL. We cannot determine with certainty the causal relationship of
626 any specific site with a trait. To test if there was a general trend of DNAm sites causally
627 influencing a trait we evaluated if the MR effect estimate based on the colocalising
628 signals were consistent with those obtained based on the secondary signals. There

629 were substantially more large genetic effects of the secondary mQTL on respective
630 traits than expected by chance (70 with $p < 0.05$, binomial test $p = 2.4e-16$). However
631 only 41 (59%) of these had effect estimates in the same direction as the primary
632 colocalising variant, which is not substantially better than chance (binomial test $p =$
633 0.19). Twelve of the 41 mQTL were located in the *HLA* region. Of the remaining mQTL,
634 27 were associated with anthropometric (*ESR1* and birth weight), immune response
635 (*IRF5* and systemic lupus erythematosus) and lipid traits (*TBL2* and triglycerides). We
636 then performed systematic colocalization analysis of all mQTL against 36 blood cell
637 traits⁴³. Here we discovered 94,738 instances of a DNAm site and a blood cell trait
638 sharing a causal variant. In 28,138 instances the colocalising DNAm site had an
639 independent secondary mQTL, and with these associations we again tested for a
640 general trend of DNAm sites causally influencing the blood trait. The association
641 between independent signals was very weak ($R^2 = 0.008$), suggesting that the general
642 causal model is not supported. Together, across the sites that were analysable in this
643 manner, these results indicate that those blood measured DNAm sites that have shared
644 genetic factors with traits cannot be typically thought of as mediating the genetic
645 association to the trait (**Figure S40-S41, Table S22**). Instead, if DNAm is a coregulatory
646 phenomenon then the colocalising signals between DNAm sites and complex traits may
647 be due to a common cause, for example genetic variants primarily acting on TF
648 binding.^{8,10}

649 The influence of traits on DNAm variation

650 Previous studies have not been adequately powered to estimate the causal influences
651 of complex traits on DNAm variation through MR, as the sample size of the outcome
652 variable (DNAm) is a predominant factor in statistical power^{52,57}. We systematically
653 analysed 109 traits for causal effects on DNAm using two-sample MR^{58,59}, where each
654 trait was instrumented using SNPs obtained from their respective previously published
655 GWAS (**Supplemental Note 2, Table S18**). Included amongst the traits were 35
656 disease traits, which when used as exposure variables in MR must be interpreted in
657 terms of the influence of liability rather than presence/absence of disease. The sample
658 size used to estimate SNP effects in DNAm was up to 27,750 (**Figure 4**).

659
660 We initially identified 4785 associations where risk factors or genetic liability to disease
661 influences DNAm levels (multiple testing threshold $p < 1.4e-7$). However, MR analysis
662 on omic variables can lead to false positives due to violations in assumptions. We
663 developed a filtering process involving a novel causal inference method to help protect
664 against these invalid associations (**Supplementary Information, Supplementary Note**
665 **2, Figure S42**). This left 85 associations (involving 84 DNAm sites) in which DNAm
666 sites were putatively influenced by 13 traits (nine risk factors or four diseases) (**Table**
667 **S23**). Further filtering that would exclude traits that were predominantly instrumented by
668 variants in the *HLA* region or driven by one SNP would reduce the total number of
669 associations substantially from 84 to 19. We replicated five associations for triglycerides

670 influencing DNAm sites near *CPTA1* and *ABCG1*⁶⁰ and found associations for
671 transferrin saturation/iron influencing DNAm sites near *HFE*.

672
673 We next evaluated if there was evidence for small, widespread changes in DNAm levels
674 in response to complex trait variation, by calculating the genomic control inflation factor
675 (GC_{in}) for the p-values obtained from the MR analyses of each trait against all DNAm
676 sites. Five traits (fasting glucose, age at menarche, cigarettes smoked per day,
677 immunoglobulin G index levels, serum creatinine), showed GC_{in} values above 1.05
678 (**Figure S43**). A high GC_{in} value can be the result of the trait that has an influence on a
679 few sites or has a widespread effect on DNAm. GC_{in} calculations were performed at
680 each chromosome singly for each trait (**Figure S44**) and in a leave-one-chromosome-
681 out analysis (**Figure S45**). The GC_{in} remained consistent (except for immunoglobulin G
682 index levels), indicating that the traits have small but widespread influences on DNAm
683 levels across the genome.

684
685 While most of the traits (n=105, 96%) tested did not appear to induce genome-wide
686 enrichment this does not rule out the possibility of them having many localised small
687 effects. For example, the smallest MR p-value for the analysis of body mass index on
688 DNAm levels was $2.27e-6$, which did not withstand genome-wide multiple testing
689 correction, and GC_{in} was 0.95. However, restricting GC_{in} to 187 sites known to
690 associate with body mass index from previous epigenome-wide association studies
691 (EWAS)¹⁹ indicated a strong enrichment of low p-values (median $GC_{in} = 3.95$). A similar
692 pattern was found for triglycerides, in which genome-wide median $GC_{in} = 0.94$ but the
693 10 sites known to associate with triglycerides from previous EWAS⁶¹ had an MR p-value
694 of $8.3e-70$ (Fisher's combined probability test). These results indicate that traits causally
695 influencing DNAm levels in blood is the most likely mechanism that gives rise to these
696 EWAS hits. It also indicates that the general finding that there were very few filtered
697 putative causal effects of risk factors or genetic liability to disease on DNAm could be
698 due to true positives being generally very small, even to the extent that our sample size
699 of up to 27,750 individuals was insufficient to find them.

700 DNAm sites influenced by genetic variation are under selection

701 Natural selection has modified the allele frequency of complex trait associated variants
702 through their beneficial or deleterious effects on fitness⁶²⁻⁶⁵. Here we investigate
703 whether mQTL SNPs are frequent targets of natural selection utilizing selection scores
704 acting through different timescales and mechanisms to each SNP in 1000G: a
705 population differentiation method (global F_{st}), several haplotype-based methods
706 (integrated haplotype score (iHS), Cross Population Extended Haplotype Homozygosity
707 (XPEHH) and the singleton density score (SDS) (**Table S24, Supplementary**
708 **Information**).

709
710 We then tested whether there is enrichment of mQTL associations (Bonferroni adjusted
711 $p < 0.01$) among SNPs that show evidence of positive selection for each metric while

712 controlling for non-random genomic distribution⁴⁷ (excluding two regions (*HLA* and *LCT*)
713 known to be under high selective pressure). We found enrichments of positive selection
714 signatures among SNPs with *cis only* (F_{st} : $p=7.87e-23$, $OR=1.31$, SDS : $p=4.43e-10$,
715 $OR=1.42$) and *cis+trans* (F_{st} : $p=7.1e-21$, $OR=1.35$, SDS : $p=4.35e-11$, $OR=1.53$, XPEHH
716 (CEU vs CHB): $p=7.7e-7$, $OR=1.53$) associations (**Figure 2F, Table S25**). The strong
717 enrichments for *cis+trans* ($n=107-1585$) and *cis only* ($n=1186-4980$) indicating that
718 positive selection is most likely to operate on *cis* acting variants. However, there is less
719 power to detect these enrichments for *trans only* SNPs ($n=14-102$).

720
721 We next examined whether there was a relationship between the mQTL effect sizes
722 (allele frequency adjusted) and the selection scores as a proxy for the estimated
723 strength of selection. Using a linear model for each of the selection metrics (accounting
724 for the number of proxies, distance to TSS, CpG and GC frequency), we found that the
725 strongest mQTL effect size was positively associated with F_{st} ($p<1.1e-05$) but not with
726 recent changes in allele frequency (measured by SDS) with consistent directions across
727 the mQTL categories (*cis only*, *cis+trans* and *trans only*) (**Figure S46**). These results
728 may indicate that DNA sites might either the primary target of selection or the mQTL
729 SNP have pleiotropic effects on fitness⁶⁶.

730
731 Enrichment of F_{st} amongst mQTL could also be due to negative selection. Evidence for
732 negative selection can be inferred from the strong negative relationship between mQTL
733 SNP effect size and MAF (difference in mQTL SNP effect size= -0.56 , $p=2.2e-308$,
734 **Figure S46**). To confirm that this relationship is not an artefact of having defined the
735 SNP effect via the maximum effect each SNP has on any DNAm site, we developed a
736 novel method (**Supplementary Information, Figure S47**) to quantify the relationship
737 for the strongest acting SNPs at a given frequency, allowing for a majority of unselected
738 SNPs. SNPs with a higher frequency have a smaller average effect ($S=0.4$, CI 0.325-
739 0.475), where $S=0$ corresponds to no selection and $S=1$ corresponds to strong negative
740 selection. We found similar relationships across the mQTL categories (*cis only*,
741 *cis+trans* and *trans only*) (**Figure S48**) though there was insufficient power to quantify
742 selection for *trans only* SNPs. These results can be interpreted that predominantly
743 genetic regions that regulate DNAm are under negative or balancing selection^{66,67} and
744 thus, retain the ancestral DNAm structure. However, a minority of regions containing
745 DNAm sites have experienced positive selection.

746
747 Alleles showing evidence of selection are likely to be biologically meaningful⁶⁸. To
748 investigate whether genetic variants underlying DNAm implicated in selection are linked
749 to diseases/traits, we examined whether GWAS-associated variants from 42 datasets
750 across 11 disease categories were enriched for *cis* mQTL SNPs overlapping extreme
751 SDS scores. After accounting for non-random genomic distribution⁴⁷, we found that
752 GWAS-associated variants from 19/42 traits were overlapping with at least one *cis*
753 mQTL SNP with extreme SDS . We found an enrichment of mQTL SNPs overlapping
754 extreme SDS scores ($p<2.6e-3$) among variants associated with five traits including
755 extreme height ($OR=17.2$, $p=1.08e-7$), Crohn's disease ($OR=11.3$, $p=4.42e-5$), height

756 (OR=1.99, p=6.76e-5), schizophrenia (OR=5.28, p=1.21e-3) and cardiovascular disease
757 (OR=9.85, p=1.67e-3) (**Table S26**). A comparison showed that the genetic variance for
758 cardiovascular disease associated mQTL or height associated mQTL with extreme SDS
759 was higher when compared to all trait associated SNPs (**Figure S49**). To summarize,
760 our results provide the first evidence that selection may have shaped the landscape of
761 DNAm values of the 450k sites although the mechanism for the selection signals that
762 exist at these loci remains unknown.

763 Implications

764 A map of hundreds of thousands of genetic associations has enabled novel biological
765 insights related to DNAm variation. Using a rigorous analytical framework enabled us to
766 minimise heterogeneity and expand sample sizes for large omic data. This revealed a
767 genetic architecture of DNAm that is polygenic. Given the diverse ranges of age, gender
768 proportions and geographical origins between the cohorts in this analysis, the minimal
769 extent of heterogeneity across datasets indicates that genetic effects on DNAm are
770 relatively stable across contexts. We show that *cis* and *trans* mQTL operate through
771 distinct mechanisms, as their genomic properties are distinct. A driver of long-range
772 associations may be co-regulated through TF binding and nuclear organisation.

773
774 Though we found substantial sharing of genetic signals between DNAm sites and
775 complex traits, we were able to demonstrate that this was not predominantly due to
776 DNAm variation being on the causal path from genotype to phenotype. While our results
777 include <2% of the DNAm sites in the genome and are limited by the two-phase design,
778 these findings have several implications especially in the context of EWAS studies that
779 are often based on the same tissue and DNAm array. First, we anticipate that some
780 previously reported EWAS associations are likely due to reverse causation e.g. the risk
781 factor or genetic liability to disease state itself alters DNAm and not vice versa, or
782 confounding. Second, having found there are strong negative and positive selection
783 pressures acting on mQTL, this may be explained through selection acting on complex
784 traits first. Third, the genetic effects on DNAm that overlap with complex traits likely
785 primarily influence other regulatory factors which in turn influence complex traits and
786 DNAm through diverging pathways. Fourth, DNAm might be on the causal pathway in a
787 disease-relevant cell type or context. Fifth, if the path from genotype to complex traits is
788 non-linear, for example involving the statistical interactions between different regulatory
789 features¹⁵, then our results indicate that large individual-level multi-omic datasets will be
790 required to dissect such mechanisms.

791
792 Future studies may be more fruitful in finding causal relationships with complex disease.
793 Either EPIC arrays¹³ or low-cost sequencing technologies⁶⁹ will expedite detailed
794 interrogations of enhancer and other regulatory regions. Especially, single molecule
795 long-read sequencing promises to expand the genetic and epigenetic spectrum by
796 allowing the detection of complex genetic variation such as allele specific DNAm and

797 structural variation and different types of DNA modifications. However, as disease
798 relevant signals and regulatory regions may be cell type specific, new analytical tools
799 are required to infer cell type specific mQTL from bulk tissue. Given our projection of
800 mQTL yields expected for future studies, pleiotropy involving mQTL is likely to be
801 increasingly important to model when interpreting genotype-trait pathways.

802

803 Overall our data and results have resulted in the most comprehensive atlas of genetic
804 effects to date. We expect that this atlas will be of use to the scientific community for
805 studies of genome regulation, contribute to the control of confounding in EWAS and to
806 perform causality analysis.

Figure 1: Discovery and replication of mQTL

a) Study Design. In the first phase, 22 cohorts performed a complete mQTL analysis of up to 480,000 sites against up to 12 million variants; retaining their results for $p < 1e-5$. In the second phase, 120 million SNP-DNA_m site pairs selected from the first phase, and GWA catalog SNPs against 345k DNA_m sites, were tested in 36 studies (including 20 phase 1 studies) and meta-analysed. **b) Distributions of the weighted mean of DNA_m across 36 cohorts for *cis only*, *cis+trans* and *trans only* sites.** Plots are coloured with respect to the genomic annotation. *Cis only* sites showed a bimodal distribution of DNA_m. *Cis+trans* sites showed intermediate levels of DNA_m. *Trans only* sites showed low levels of DNA_m. **c) Discovery and replication effect size estimates** between GoDMC (n=27,750) and Generation Scotland (n=5,101) for 169,656 mQTL associations. The regression coefficient is 1.13 (se=0.0007). **d) Relationship between DNA_m site heritability estimates and DNA_m variance explained in Generation Scotland.** The regression coefficient for the twin family study was 3.16 (se=0.008) and for the twin study 2.91 (se=0.008) across 403,353 DNA_m sites. The variance explained for DNA_m sites with missing r^2 (n=277,428) and/or $h^2=0$ (Twin family: n=80,726 Twin: 34,537) were set to 0.

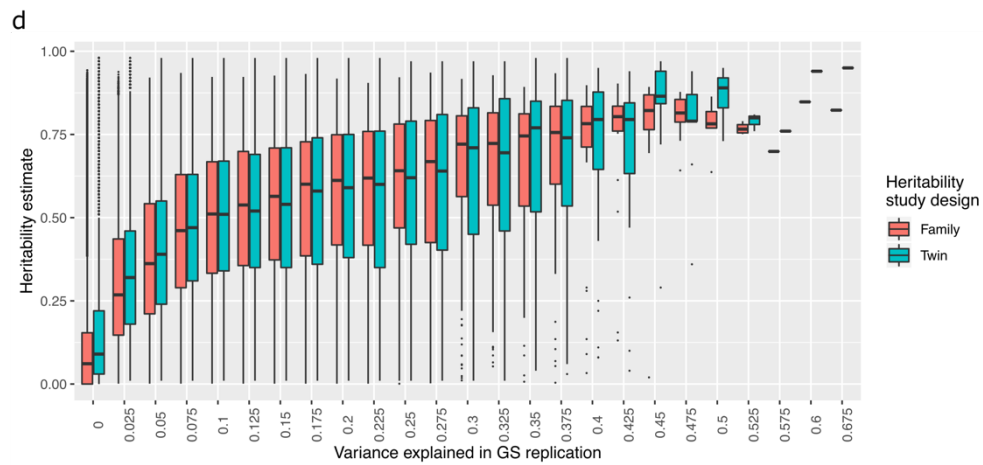
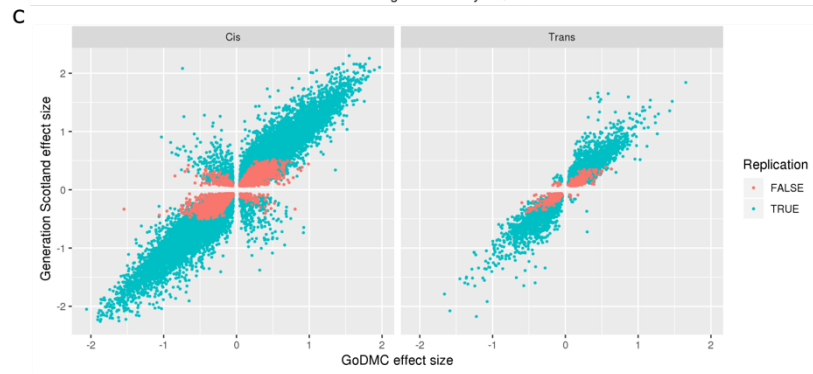
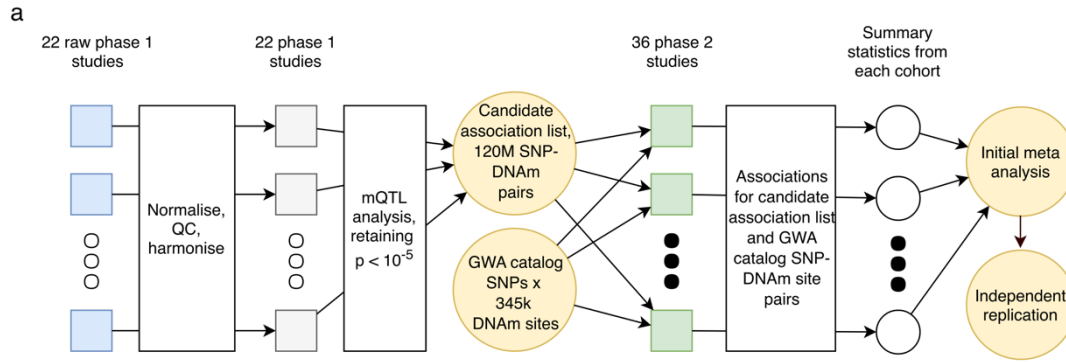


Figure 2: *Cis* and *trans* mQTL operate through distinct mechanisms

a) Distributions of enrichments for chromatin states and gene annotations among mQTL sites and SNPs. The heatmap represents the distribution of odds ratios for *cis only*, *trans only*, or *cis+trans* sites and SNPs. Significance has been categorised as: *=FDR<0.001; **=FDR<1e-10; ***=FDR<1e-50

b) Distributions of enrichment for occupancy of TFBS among mQTL sites and SNPs. Each density curve represents the distribution of odds ratios for *cis only*, *trans only*, or *cis+trans* sites (left) and SNPs (right).

c) A bipartite graph of the two-dimensional enrichment for *trans*-mQTL. SNPs annotations (blue) with $p_{\text{emp}} < 0.01$ after multiple testing correction co-occur with particular site annotations (red).

d) Distribution of two-dimensional enrichment values of *trans*-mQTL. There was substantial departure from the null in the real dataset for all tissues indicating that the TFBS of a site depended on the TFBS of the SNP that influenced it.

e) Distributions of enrichment of mQTL among 41 complex traits and diseases. Each density curve represents the distribution of odds ratios for *cis only*, *trans only*, or *cis+trans* SNPs.

f) Enrichment of selection signals among mQTL SNPs. Radial lines show odds ratios for the different selection metrics (F_{st} , SDS, iHS, XPEHH (CEU vs CHB) and XPEHH (CEU vs YRI) by site annotation (*cis any*, *cis only*, *cis+trans*, *trans only*, *trans any*). Dots in the inner ring of the outer circle denote enrichment (if present) at thresholds $p < 1e-11$ (outermost) to $p < 1e-14$ (innermost).

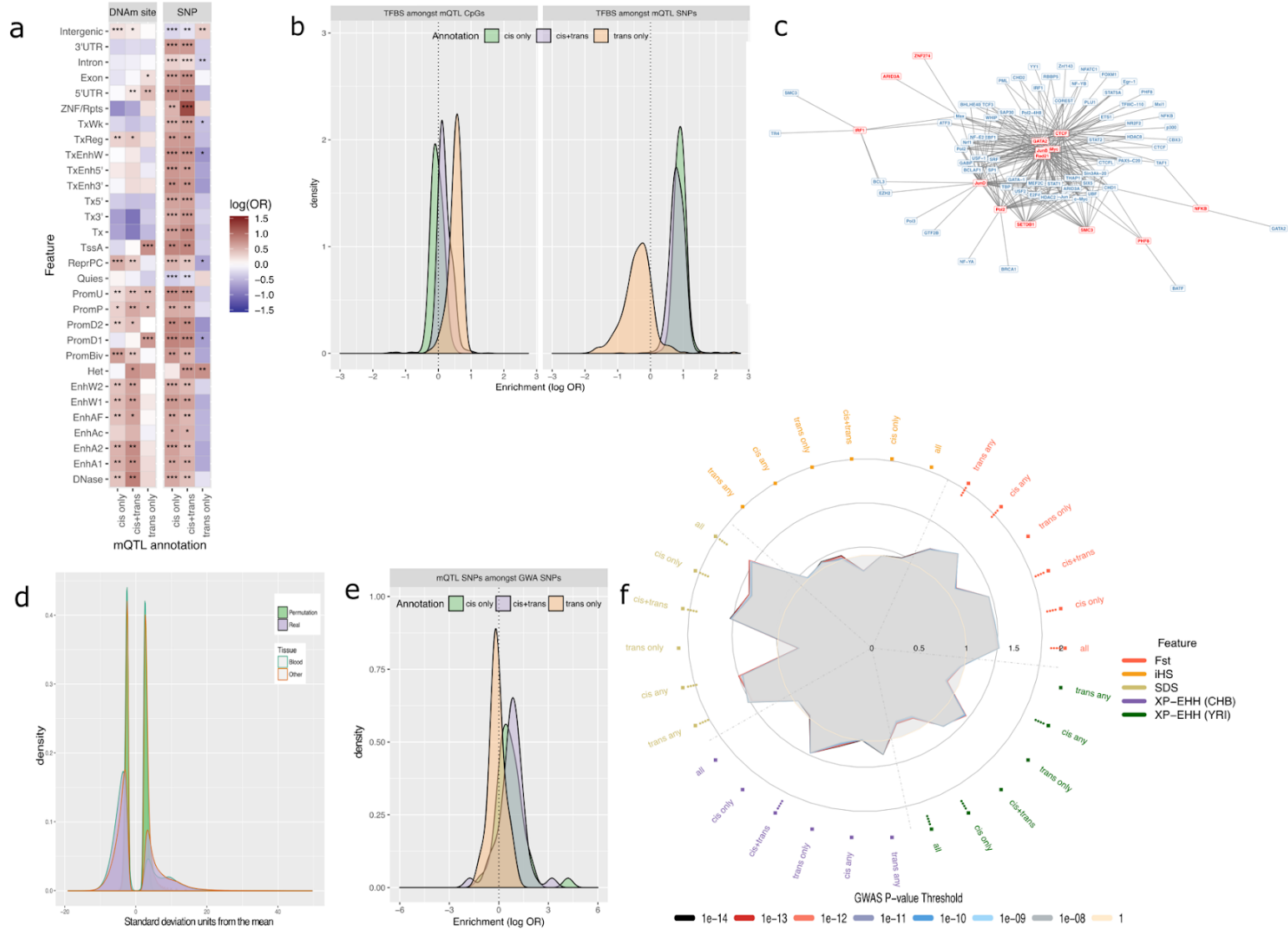


Figure 3: Communities constructed from *trans*-mQTL. **a) A network depicting all communities in which there were twenty or more sites.** Random walks were used to generate communities (colours), so occasionally a DNA site connects different communities. **b) The relationship between genomic annotations, mQTL and communities.** Communities 9 and 22 are comprised of DNAm sites that are related through shared genetic factors. The sankey plots show the genomic annotations for the genetic variants (left) and for the DNAm sites (right). The DNAm sites comprising these communities are enriched for TFBS related to the cohesin complex and NFkB, respectively. **c) Enrichment of GWA traits among community SNPs.** The genomic loci for each of the 56 largest communities were tested for enrichment of low p-values in 133 complex trait GWAS (y-axis). The x-axis depicts the $-\log_{10}$ p-value for enrichment, with the 5% FDR shown by the vertical dotted line. Enrichments were particularly strong for blood related phenotypes (including circulating metal levels).

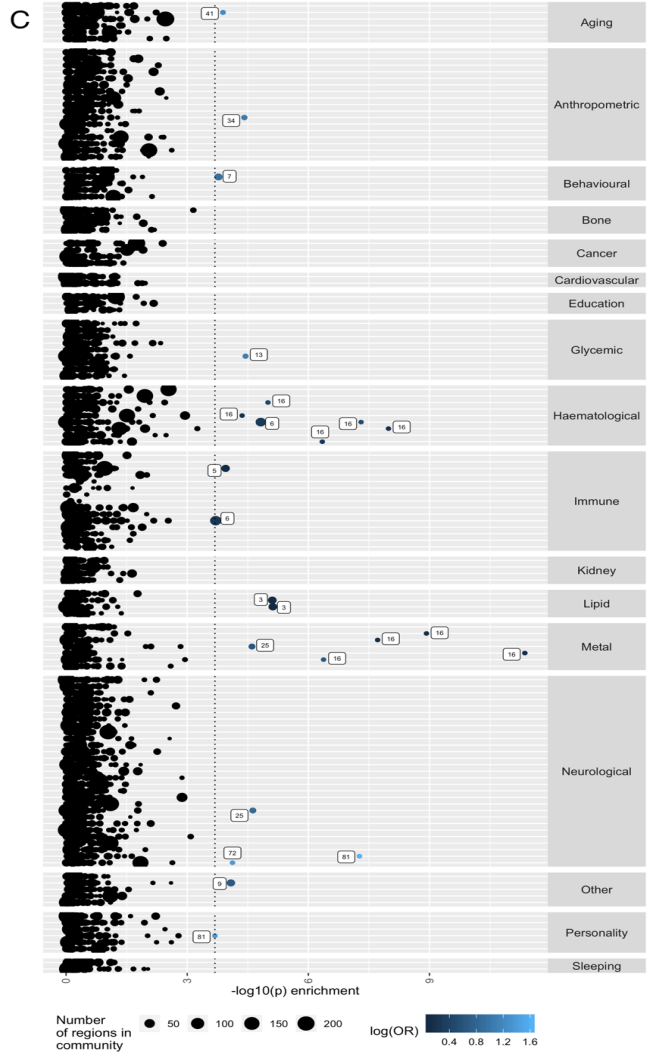
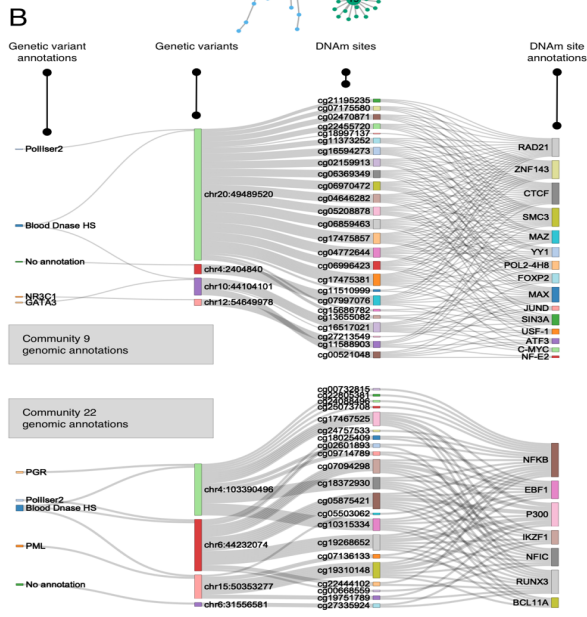
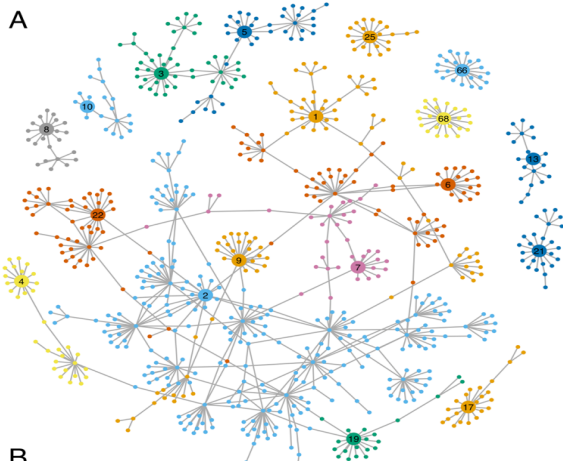
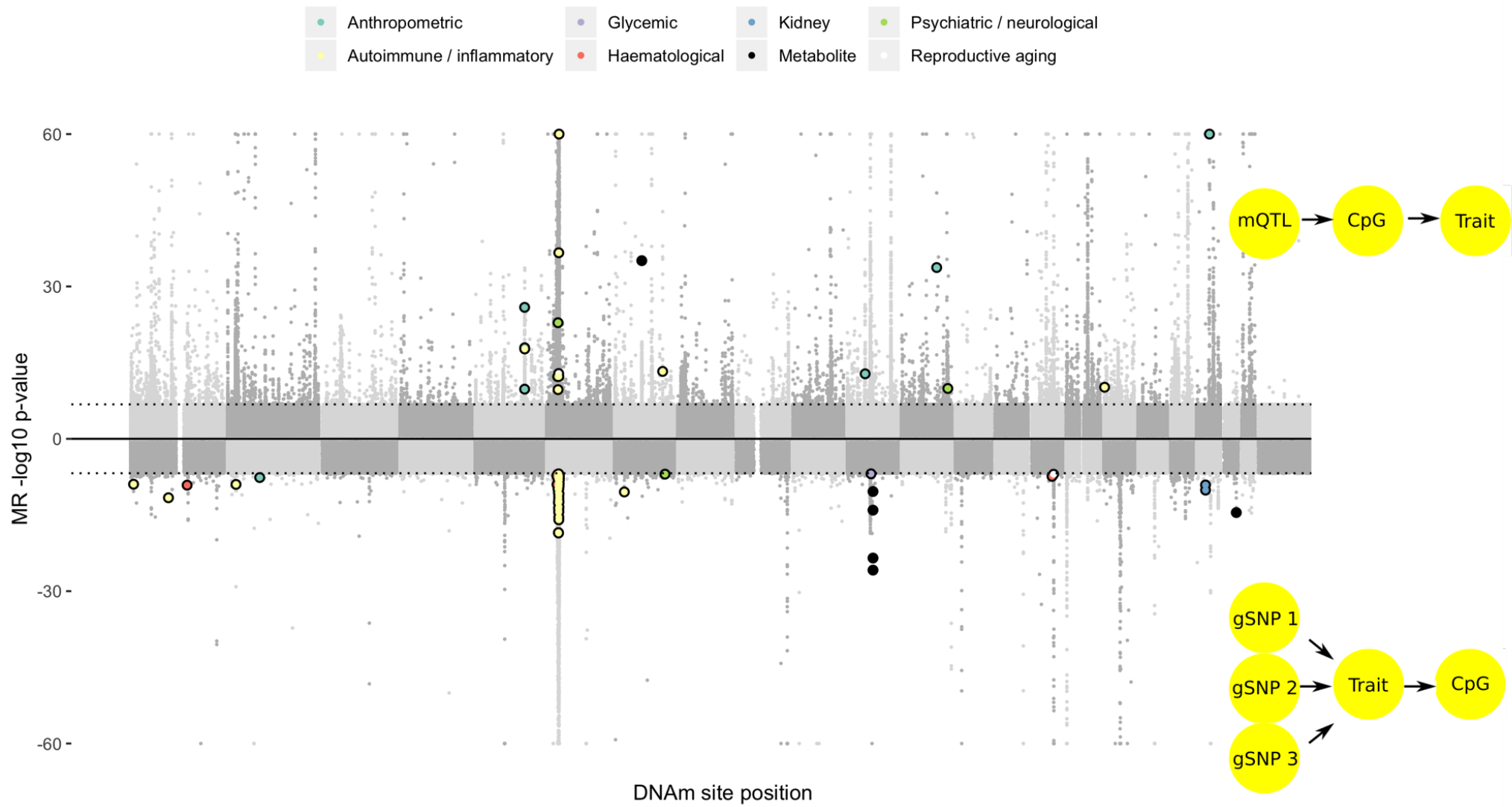


Figure 4: Identifying putative causal relationships between sites and traits using bi-directional MR. Aggregated results from a systematic bi-directional MR analysis between DNAm sites and 116 complex traits. The top plot depicts results from tests of DNAm sites colocalising with complex traits. The light grey points represent MR estimates that either did not surpass multiple testing, or shared small p-values at both the DNAm site and complex trait but had weak evidence of colocalisation. Bold, coloured points are those that showed strong evidence for colocalisation ($H_4 > 0.8$). The bottom plot shows the $-\log_{10}$ p-values from MR analysis of risk factor or genetic liability of disease on DNAm levels. Extensive follow up was performed on DNAm site-trait pairs with putative associations, and those that pass filters are plotted in bold and colored according to the trait category. A substantial number of MR results in both directions exhibited very strong effects but failed to withstand sensitivity analyses.



References

1. Petronis, A. Epigenetics as a unifying principle in the aetiology of complex traits and diseases. *Nature* **465**, 721-7 (2010).
2. van Dongen, J. *et al.* Genetic and environmental influences interact with age and sex in shaping the human methylome. *Nat Commun* **7**, 11115 (2016).
3. Hannon, E. *et al.* Characterizing genetic and environmental influences on variable DNA methylation using monozygotic and dizygotic twins. *PLoS Genet* **14**, e1007544 (2018).
4. Kerkel, K. *et al.* Genomic surveys by methylation-sensitive SNP analysis identify sequence-dependent allele-specific DNA methylation. *Nat Genet* **40**, 904-8 (2008).
5. Schadt, E.E. *et al.* Genetics of gene expression surveyed in maize, mouse and man. *Nature* **422**, 297-302 (2003).
6. Davey Smith, G. & Hemani, G. Mendelian randomization: genetic anchors for causal inference in epidemiological studies. *Hum Mol Genet* **23**, R89-98 (2014).
7. Gaunt, T.R. *et al.* Systematic identification of genetic influences on methylation across the human life course. *Genome Biol* **17**, 61 (2016).
8. Bonder, M.J. *et al.* Disease variants alter transcription factor levels and methylation of their binding sites. *Nat Genet* **49**, 131-138 (2017).
9. Hannon, E. *et al.* Methylation QTLs in the developing brain and their enrichment in schizophrenia risk loci. *Nat Neurosci* **19**, 48-54 (2016).
10. Hop, P.J. *et al.* Genome-wide identification of genes regulating DNA methylation using genetic anchors for causal inference. *Genome Biol* **21**, 220 (2020).
11. Abecasis, G.R. *et al.* An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**, 56-65 (2012).
12. Yang, J. *et al.* Conditional and joint multiple-SNP analysis of GWAS summary statistics identifies additional variants influencing complex traits. *Nat Genet* **44**, 369-75, S1-3 (2012).
13. Pidsley, R. *et al.* Critical evaluation of the Illumina MethylationEPIC BeadChip microarray for whole-genome DNA methylation profiling. *Genome Biol* **17**, 208 (2016).
14. Shah, S. *et al.* Genetic and environmental exposures constrain epigenetic drift over the human life course. *Genome Res* **24**, 1725-33 (2014).
15. Gutierrez-Arcelus, M. *et al.* Passive and active DNA methylation and the interplay with genetic variation in gene regulation. *Elife* **2**, e00523 (2013).
16. Chen, L. *et al.* Genetic Drivers of Epigenetic and Transcriptional Variation in Human Immune Cells. *Cell* **167**, 1398-1414.e24 (2016).
17. McRae, A.F. *et al.* Identification of 55,000 Replicated DNA Methylation QTL. *Sci Rep* **8**, 17605 (2018).
18. Kundaje, A. *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317-30 (2015).
19. Wahl, S. *et al.* Epigenome-wide association study of body mass index, and the adverse outcomes of adiposity. *Nature* **541**, 81-86 (2017).
20. Elliott, G. *et al.* Intermediate DNA methylation is a conserved signature of genome regulation. *Nat Commun* **6**, 6363 (2015).
21. Feldmann, A. *et al.* Transcription factor occupancy can mediate active turnover of DNA methylation at regulatory regions. *PLoS Genet* **9**, e1003994 (2013).
22. Grundberg, E. *et al.* Global analysis of DNA methylation variation in adipose tissue from twins reveals links to disease-associated variants in distal regulatory elements. *Am J Hum Genet* **93**, 876-90 (2013).
23. Kim-Hellmuth, S. *et al.* Cell type-specific genetic regulation of gene expression across human tissues. *Science* **369**(2020).
24. Qi, T. *et al.* Identifying gene targets for brain-related traits using transcriptomic and methylomic data from blood. *Nat Commun* **9**, 2282 (2018).

25. Yin, Y. *et al.* Impact of cytosine methylation on DNA binding specificities of human transcription factors. *Science* **356**(2017).
26. Domcke, S. *et al.* Competition between DNA methylation and transcription factors determines binding of NRF1. *Nature* **528**, 575-9 (2015).
27. Baubec, T. *et al.* Genomic profiling of DNA methyltransferases reveals a role for DNMT3B in genic methylation. *Nature* **520**, 243-7 (2015).
28. Ginno, P.A. *et al.* A genome-scale map of DNA methylation turnover identifies site-specific dependencies of DNMT and TET activity. *Nat Commun* **11**, 2680 (2020).
29. Sánchez-Castillo, M. *et al.* CODEX: a next-generation sequencing experiment database for the haematopoietic and embryonic stem cell communities. *Nucleic Acids Res* **43**, D1117-23 (2015).
30. Consortium, E.P. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57-74 (2012).
31. Schübeler, D. Function and information content of DNA methylation. *Nature* **517**, 321-6 (2015).
32. Waszak, S.M. *et al.* Population Variation and Genetic Control of Modular Chromatin Architecture in Humans. *Cell* **162**, 1039-50 (2015).
33. Viny, A.D. *et al.* Dose-dependent role of the cohesin complex in normal and malignant hematopoiesis. *J Exp Med* **212**, 1819-32 (2015).
34. Battle, A. *et al.* Genetic effects on gene expression across human tissues. *Nature* **550**, 204-213 (2017).
35. Kumasaka, N., Knights, A.J. & Gaffney, D.J. High-resolution genetic mapping of putative causal interactions between regions of open chromatin. *Nat Genet* **51**, 128-137 (2019).
36. Chinenov, Y., Coppo, M., Gupte, R., Sacta, M.A. & Rogatsky, I. Glucocorticoid receptor coordinates transcription factor-dominated regulatory network in macrophages. *BMC Genomics* **15**, 656 (2014).
37. Geller, D.A. *et al.* Cytokine induction of interferon regulatory factor-1 in hepatocytes. *Surgery* **114**, 235-42 (1993).
38. Gilchrist, M. *et al.* Systems biology approaches identify ATF3 as a negative regulator of Toll-like receptor 4. *Nature* **441**, 173-8 (2006).
39. Yang, J., Tian, B., Sun, H., Garofalo, R.P. & Brasier, A.R. Epigenetic silencing of IRF1 dysregulates type III interferon responses to respiratory virus infection in epithelial to mesenchymal transition. *Nat Microbiol* **2**, 17086 (2017).
40. Rao, S.S. *et al.* A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell* **159**, 1665-80 (2014).
41. Delaneau, O. *et al.* Chromatin three-dimensional interactions mediate genetic effects on gene expression. *Science* **364**(2019).
42. Vosa, U. Unraveling the polygenic architecture of complex traits using blood eQTL metaanalysis. *bioRxiv* (2018).
43. Astle, W.J. *et al.* The Allelic Landscape of Human Blood Cell Trait Variation and Links to Common Complex Disease. *Cell* **167**, 1415-1429.e19 (2016).
44. Maurano, M.T. *et al.* Systematic localization of common disease-associated variation in regulatory DNA. *Science* **337**, 1190-5 (2012).
45. Tachmazidou, I. *et al.* Whole-Genome Sequencing Coupled to Imputation Discovers Genetic Signals for Anthropometric Traits. *Am J Hum Genet* **100**, 865-884 (2017).
46. Kato, N. *et al.* Trans-ancestry genome-wide association study identifies 12 genetic loci influencing blood pressure and implicates a role for DNA methylation. *Nat Genet* **47**, 1282-1293 (2015).
47. Iotchkova, V. *et al.* GARFIELD classifies disease-relevant genomic features through integration of functional annotations with association signals. *Nat Genet* **51**, 343-353 (2019).
48. Finucane, H.K. *et al.* Partitioning heritability by functional annotation using genome-wide association summary statistics. *Nat Genet* **47**, 1228-35 (2015).

49. Reinius, L.E. *et al.* Differential DNA methylation in purified human blood cells: implications for cell lineage and studies on disease susceptibility. *PLoS One* **7**, e41361 (2012).
50. Houseman, E.A. *et al.* Model-based clustering of DNA methylation array data: a recursive-partitioning algorithm for high-dimensional data arising as a mixture of beta distributions. *BMC Bioinformatics* **9**, 365 (2008).
51. Hemani, G., Tilling, K. & Davey Smith, G. Orienting the causal relationship between imprecisely measured traits using GWAS summary data. *PLoS Genet* **13**, e1007081 (2017).
52. Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of genetic association studies using summary statistics. *PLoS Genet* **10**, e1004383 (2014).
53. Zhu, Z. *et al.* Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. *Nat Genet* **48**, 481-7 (2016).
54. Zheng, J. *et al.* Phenome-wide Mendelian randomization mapping the influence of the plasma proteome on complex diseases. *Nat Genet* (2020).
55. Richardson, T.G. *et al.* Systematic Mendelian randomization framework elucidates hundreds of CpG sites which may mediate the influence of genetic variants on disease. *Hum Mol Genet* **27**, 3293-3304 (2018).
56. Hemani, G., Bowden, J. & Davey Smith, G. Evaluating the potential role of pleiotropy in Mendelian randomization studies. *Hum Mol Genet* **27**, R195-R208 (2018).
57. Brion, M.J., Shakhbazov, K. & Visscher, P.M. Calculating statistical power in Mendelian randomization studies. *Int J Epidemiol* **42**, 1497-501 (2013).
58. Pierce, B.L. & Burgess, S. Efficient design for Mendelian randomization studies: subsample and 2-sample instrumental variable estimators. *Am J Epidemiol* **178**, 1177-84 (2013).
59. Hemani, G. *et al.* The MR-Base platform supports systematic causal inference across the human phenome. *Elife* **7**(2018).
60. Dekkers, K.F. *et al.* Blood lipids influence DNA methylation in circulating cells. *Genome Biol* **17**, 138 (2016).
61. Braun, K.V.E. *et al.* Epigenome-wide association study (EWAS) on lipids: the Rotterdam Study. *Clin Epigenetics* **9**, 15 (2017).
62. Byars, S.G. *et al.* Genetic loci associated with coronary artery disease harbor evidence of selection and antagonistic pleiotropy. *PLoS Genet* **13**, e1006328 (2017).
63. Turchin, M.C. *et al.* Evidence of widespread selection on standing variation in Europe at height-associated SNPs. *Nat Genet* **44**, 1015-9 (2012).
64. Field, Y. *et al.* Detection of human adaptation during the past 2000 years. *Science* **354**, 760-764 (2016).
65. Zeng, J. *et al.* Signatures of negative selection in the genetic architecture of human complex traits. *Nat Genet* **50**, 746-753 (2018).
66. Johnson, T. & Barton, N. Theoretical models of selection and mutation on quantitative traits. *Philos Trans R Soc Lond B Biol Sci* **360**, 1411-25 (2005).
67. Eyre-Walker, A. Evolution in health and medicine Sackler colloquium: Genetic architecture of a complex trait and its implications for fitness and genome-wide association studies. *Proc Natl Acad Sci U S A* **107 Suppl 1**, 1752-6 (2010).
68. Vitti, J.J., Grossman, S.R. & Sabeti, P.C. Detecting natural selection in genomic data. *Annu Rev Genet* **47**, 97-120 (2013).
69. Simpson, J.T. *et al.* Detecting DNA cytosine methylation using nanopore sequencing. *Nat Methods* **14**, 407-410 (2017).