

Meta-analysis of up to 622,409 individuals identifies 40 novel smoking behaviour associated genetic loci

A Mesut Erzurumluoglu ^{1*}, Mengzhen Liu ^{2*}, Victoria E Jackson ^{3,4,1*}, Daniel R Barnes ⁵, Gargi Datta ^{2,6}, Carl A Melbourne ¹, Robin Young ⁵, Chiara Batini ¹, Praveen Surendran ⁵, Tao Jiang ⁵, Sheikh Daud Adnan ⁷, Saima Afaq ⁸, Arpana Agrawal ⁹, Elisabeth Altmaier ¹⁰, Antonis C Antoniou ¹¹, Folkert W Asselbergs ^{12,13,14,15}, Clemens Baumbach ¹⁰, Laura Beirut ¹⁶, Sarah Bertelsen ¹⁷, Michael Boehnke ¹⁸, Michiel L Bots ^{19,20}, David M Brazel ^{6,21}, John C Chambers ^{22,8,23,24}, Jenny Chang-Claude ^{25,26}, Chu Chen ^{27,28}, Yi-Ling Chou ⁹, Janie Corley ^{29,30}, Sean P David ³¹, Rudolf A de Boer ³², Christiaan A de Leeuw ³³, Joe G Dennis ¹¹, Anna F Dominiczak ³⁴, Alison M Dunning ³⁵, Douglas F Easton ^{11,35}, Charles Eaton ²⁸, Paul Elliott ^{36,37,38,39}, Evangelos Evangelou ^{8,40}, Tatiana Foroud ⁴¹, Alison Goate ⁴², Jian Gong ⁴³, Hans J Grabe ⁴⁴, Jeff Haessler ⁴³, Christopher Haiman ⁴⁵, Göran Hallmans ⁴⁶, Anke R Hammerschlag ³³, Sarah E Harris ^{29,47}, Andrew Hattersley ⁴⁸, Andrew Heath ⁹, Chris Hsu ⁴⁹, William G Iacono ², Stavroula Kanoni ^{50,51}, Manav Kapoor ¹⁷, Jaakko Kaprio ^{52,53}, Sharon L Kardinaal ⁵⁴, Fredrik Karpe ^{55,56}, Jukka Kontto ⁵⁷, Jaspal S Kooner ^{23,24,37,58}, Charles Kooperberg ^{43,59}, Kari Kuulasmaa ⁵⁷, Markku Laakso ⁶⁰, Dongbing Lai ⁴¹, Claudia Langenberg ⁶¹, Nhung Le ⁶², Guillaume Lettre ^{63,64}, Anu Loukola ^{52,53}, Jian'an Luan ⁶¹, Pamela A F Madden ⁹, Massimo Mangino ⁶⁵, Riccardo E Marioni ^{29,47}, Eirini Marouli ^{50,51}, Jonathan Marten ⁶⁶, Nicholas G Martin ⁶⁷, Matt McGue ², Kyriaki Michailidou ^{68,11}, Evelin Mihailov ⁶⁹, Alireza Moayyeri ⁷⁰, Marie Moitry ⁷¹, Martina Müller-Nurasyid ^{72,73,74}, Aliya Naheed ⁷⁵, Matthias Nauck ^{76,77}, Matthew J Neville ^{55,56}, Sune Fallgaard Nielsen ⁷⁸, Kari North ⁷⁹, Jessica D Faul ⁸⁰, Markus Perola ^{52,57}, Paul D P Pharoah ^{11,35}, Giorgio Pistis ⁸¹, Tinca J Polderman ³³, Danielle Posthuma ^{33,82}, Neil Poulter ^{8,83}, Beenish Qaiser ^{52,53}, Asif Rasheed ⁸⁴, Alex Reiner ^{43,28}, Frida Renström ^{85,86}, John Rice ⁸⁷, Rebecca Rohde ⁸⁸, Olov Rolandsson ⁸⁹, Nilesh J Samani ⁹⁰, Maria Samuel ⁸⁴, David Schlessinger ⁹¹, Steven H Scholte ⁹², Robert A Scott ⁶¹, Peter Sever ^{58,83}, Yaming Shao ⁸⁸, Nick Shrine ¹, Jennifer A Smith ⁵⁴, John M Starr ^{29,93}, Kathleen Stirrups ^{94,50}, Danielle Stram ⁹⁵, Heather M Stringham ¹⁸, Ioanna Tachmazidou ⁹⁶, Jean-Claude Tardif ^{63,64}, Deborah J Thompson ¹¹, Hilary A Tindle ⁹⁷, Vinicius Tragante ⁹⁸, Stella Trompet ^{99,100}, Valerie Turcot ^{63,64}, Jessica Tyrrell ⁴⁸, Ilonca Vaartjes ^{19,20}, Andries R van der Leij ⁹², Peter van der Meer ³², Tibor V Varga ⁸⁵, Niek Verweij ^{32,101}, Henry Völzke ^{102,77}, Nicholas J Wareham ⁶¹, Helen R Warren ^{103,104}, David R Weir ⁸⁰, Stefan Weiss ^{105,77}, Leah Wetherill ⁴¹, Hanieh Yaghoobkar ⁴⁸, Ersin Yavas ^{106,107}, Yu Jiang ¹⁰⁸, Fang Chen ¹⁰⁸, Xiaowei Zhan ¹⁰⁹, Weihua Zhang ^{8,110}, Wei Zhao ¹¹¹, Wei Zhao ⁵⁴, Kaixin Zhou ¹¹², Philippe Amouyel ¹¹³, Stefan Blankenberg ^{114,115}, Mark J Caulfield ^{103,104}, Rajiv Chowdhury ⁵, Francesco Cucca ⁸¹, Ian J Deary ^{29,30}, Panos Deloukas ^{116,96,117}, Emanuele Di Angelantonio ^{118,5}, Marco Ferrario ¹¹⁹, Jean Ferrières ¹²⁰, Paul W Franks ^{85,121}, Tim M Frayling ⁴⁸, Philippe Frossard ⁸⁴, Ian P Hall ¹²², Caroline Hayward ⁶⁶, Jan-Håkan Jansson ¹²³, J Wouter Jukema ^{124,125}, Frank Kee ¹²⁶, Satu Männistö ⁵⁷, Andres Metspalu ⁶⁹, Patricia B Munroe ^{103,104}, Børge Grønne Nordestgaard ⁷⁸, Colin N A Palmer ¹²⁷, Veikko Salomaa ⁵⁷, Naveed Sattar ¹²⁸, Timothy Spector ¹²⁹, David Peter Strachan ¹³⁰, Understanding Society Scientific Group, EPIC-CVD, GSCAN, Consortium for Genetics of Smoking Behaviour, CHD Exome+ consortium, Pim van der Harst ^{32,131}, Eleftheria Zeggini ⁹⁶, Danish

36 Saleheen ^{132,133,134,5}, Adam S Butterworth ⁵, Louise V Wain ^{1,135}, Goncalo R Abecasis ¹⁸, John Danesh ^{5,96},
37 Martin D Tobin ^{1,135†}, Scott Vrieze ^{2†}, Dajiang J Liu ^{108†#}, Joanna M M Howson ^{5†#}

38

39 1. Department of Health Sciences, University of Leicester, Leicester, UK

40 2. Department of Psychology, University of Minnesota

41 3. Population Health and Immunity Division, The Walter and Eliza Hall Institute of Medical Research, 1G
42 Royal Pde, 3052 Parkville, Australia

43 4. Department of Medical Biology, University of Melbourne, 3010 Parkville, Australia

44 5. Cardiovascular Epidemiology Unit, Department of Public Health and Primary Care, University of
45 Cambridge, Cambridge, UK

46 6. Institute for Behavioral Genetics, University of Colorado Boulder

47 7. National Institute of Cardiovascular Diseases, Sher-e-Bangla Nagar, Dhaka, Bangladesh

48 8. Department of Epidemiology and Biostatistics, Imperial College London, London W2 1PG, UK

49 9. Department of Psychiatry, Washington University

50 10. Research Unit of Molecular Epidemiology, Helmholtz Zentrum München-German Research Center for
51 Environmental Health, Neuherberg, Germany

52 11. Centre for Cancer Genetic Epidemiology, Department of Public Health and Primary Care, University of
53 Cambridge, Cambridge, UK, CB1 8RN

54 12. Department of Cardiology, Division Heart & Lungs, University Medical Center Utrecht, University of
55 Utrecht, the Netherlands

56 13. Durrer Center for Cardiovascular Research, Netherlands Heart Institute, Utrecht, the Netherlands

57 14. Institute of Cardiovascular Science, Faculty of Population Health Sciences, University College London,
58 London, United Kingdom

59 15. Farr Institute of Health Informatics Research and Institute of Health Informatics, University College
60 London, London, United Kingdom

61 16. Department of Psychiatry, Washington University School of Medicine

62 17. Department of Neuroscience, Icahn School of Medicine at Mount Sinai

63 18. Department of Biostatistics and Center for Statistical Genetics, University of Michigan, Ann Arbor,
64 Michigan

65 19. Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, 3508GA Utrecht,
66 the Netherlands

67 20. Center for Circulatory Health, University Medical Center Utrecht, 3508GA Utrecht, the Netherlands

68 21. Department of Molecular, Cellular, and Developmental Biology, University of Colorado Boulder

69 22. Lee Kong Chian School of Medicine, Nanyang Technological University, Singapore 308232, Singapore

70 23. Department of Cardiology, Ealing Hospital, Middlesex UB1 3HW, UK

71 24. Imperial College Healthcare NHS Trust, London W12 0HS, UK

72 25. Division of Cancer Epidemiology, German Cancer Research Centre (DKFZ), Heidelberg, Germany

- 73 26. Cancer Epidemiology Group, University Medical Centre Hamburg-Eppendorf, University Cancer Centre
74 Hamburg (UCCH), Hamburg, Germany
- 75 27. Public Health Sciences Division, Fred Hutchinson Cancer Research Center
- 76 28. Department of Epidemiology, University of Washington, Seattle, WA
- 77 29. Centre for Cognitive Ageing and Cognitive Epidemiology, University of Edinburgh, Edinburgh, UK, EH8
78 9JZ
- 79 30. Psychology, University of Edinburgh, Edinburgh, UK, EH8 9JZ
- 80 31. Department of Medicine, Stanford University, Stanford, CA
- 81 32. University Medical Center Groningen, University of Groningen, Department of Cardiology, the
82 Netherlands
- 83 33. Department of Complex Trait Genetics, Center for Neurogenomics and Cognitive Research, Amsterdam
84 Neuroscience, VU University Amsterdam
- 85 34. Institute of Cardiovascular and Medical Sciences, College of Medical, Veterinary and Life Sciences,
86 University of Glasgow, Glasgow, UK
- 87 35. Centre for Cancer Genetic Epidemiology, Department of Oncology, Cambridge Centre, University of
88 Cambridge, Cambridge, UK, CB1 8RN
- 89 36. Department of Epidemiology and Biostatistics, Imperial College London, London, UK
- 90 37. MRC-PHE Centre for Environment and Health, Imperial College London, London, W2 1PG, UK
- 91 38. National Institute for Health Research Imperial Biomedical Research Centre, Imperial College Healthcare
92 NHS Trust and Imperial College London, London, UK
- 93 39. UK Dementia Research Institute (UK DRI) at Imperial College London, London, UK
- 94 40. Department of Hygiene and Epidemiology, University of Ioannina Medical School, Ioannina, Greece.
- 95 41. Department of Medical and Molecular Genetics, Indiana University School of Medicine
- 96 42. Department of Neuroscience, Icahn School of Medicine at Mount Sinai, New York, USA
- 97 43. Public Health Sciences Division, Fred Hutchinson Cancer Research Center
- 98 44. Department of Psychiatry and Psychotherapy, University Medicine Greifswald, 17475 Greifswald,
99 Germany
- 100 45. Department of Preventative Medicine, Keck School of Medicine, University of Southern California
- 101 46. Department of Public Health and Clinical Medicine, Nutritional research, Umeå University, Sweden
- 102 47. Centre for Genomic and Experimental Medicine, University of Edinburgh, Edinburgh, UK, EH4 2XU
- 103 48. Genetics of Complex Traits, University of Exeter Medical School, Exeter, United Kingdom
- 104 49. University of Southern California
- 105 50. William Harvey Research Institute, Barts and The London School of Medicine and Dentistry, Queen Mary
106 University of London, London, UK, EC1M 6BQ
- 107 51. Centre for Genomic Health, Queen Mary University of London, London EC1M 6BQ, UK
- 108 52. Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Helsinki, Finland
- 109 53. Department of Public Health, University of Helsinki, Helsinki, Finland

- 110 54. Department of Epidemiology, School of Public Health, University of Michigan
111 55. Oxford Centre for Diabetes, Endocrinology and Metabolism, University of Oxford
112 56. Oxford National Institute for Health Research, Biomedical Research Centre, Churchill Hospital, Oxford,
113 UK
114 57. Department of Public Health Solutions, National Institute for Health and Welfare, FI-00271, Helsinki,
115 Finland
116 58. National Heart and Lung Institute, Imperial College London, London W12 0NN, UK
117 59. Department of Biostatistics, University of Washington School of Medicine, Seattle, WA
118 60. University of Eastern Finland, Finland
119 61. MRC Epidemiology Unit, Institute of Metabolic Science, University of Cambridge School of Clinical
120 Medicine, Cambridge, CB2 0QQ, UK
121 62. Department of Medical Microbiology, Immunology and Cell Biology, Southern Illinois University School
122 of Medicine
123 63. Montreal Heart Institute, Montreal, Quebec, H1T 1C8, Canada
124 64. Department of Medicine, Faculty of Medicine, Universite de Montreal, Montreal, Quebec, H3T 1J4,
125 Canada
126 65. Twin Research & Genetic Epidemiology Unit, Kings College, London
127 66. MRC Human Genetics Unit, MRC Institute of Genetics and Molecular Medicine, University of
128 Edinburgh, Edinburgh, UK
129 67. Queensland Institute for Medical Research
130 68. Department of Electron Microscopy/Molecular Pathology, The Cyprus Institute of Neurology and
131 Genetics, 1683 Nicosia, Cyprus
132 69. Estonian Genome Center, University of Tartu, Tartu, Estonia
133 70. Institute of Health Informatics, University College London, London, UK
134 71. Department of Epidemiology and Public health, Strasbourg University hospital, University of Strasbourg,
135 France
136 72. Institute of Genetic Epidemiology, Helmholtz Zentrum München - German Research Center for
137 Environmental Health, Neuherberg, Germany
138 73. Department of Medicine I, Ludwig-Maximilians-University Munich, Munich, Germany
139 74. DZHK (German Centre for Cardiovascular Research), Partner Site Munich Heart Alliance, Munich,
140 Germany
141 75. Initiative for Noncommunicable Diseases, Health Systems and Population Studies Division, International
142 Centre for Diarrhoeal Disease Research , Bangladesh (icddr,b) International Centre for Diarrhoeal Disease
143 Research , Bangladesh (icddr,b)
144 76. Institute of Clinical Chemistry and Laboratory Medicine, University Medicine Greifswald, 17475
145 Greifswald, Germany
146 77. DZHK (German Centre for Cardiovascular Research), Partner Site Greifswald, University Medicine,

147 Greifswald, Germany
148 78. Department of Clinical Biochemistry Herlev Hospital, Copenhagen University Hospital, Herlev Ringvej
149 74, DK-2730 Herlev, Denmark
150 79. Department of Epidemiology, University of North Carolina, Chapel Hill
151 80. Survey Research Center, Institute for Social Research, University of Michigan
152 81. Istituto di Ricerca Genetica e Biomedica, Consiglio Nazionale delle Ricerche (CNR), Monserrato,
153 Cagliari, Italy
154 82. Department of Clinical Genetics, VU University Medical Centre Amsterdam, Amsterdam Neuroscience
155 83. International Centre for Circulatory Health, Imperial College London, UK
156 84. Centre for Non-Communicable Diseases, Karachi, Pakistan
157 85. Genetic and Molecular Epidemiology Unit, Lund University Diabetes Centre, Department of Clinical
158 Sciences, Skåne University Hospital, Lund University, SE-214 28, Malmö, Sweden
159 86. Department of Biobank Research, Umeå University, SE-901 87, Umeå, Sweden
160 87. Departments of Psychiatry and Mathematics, Washington University St. Louis
161 88. University of North Carolina, Chapel Hill
162 89. Department of Public Health & Clinical Medicine, Section for Family Medicine, Umeå universitet, SE-
163 90185 Umeå, Sweden
164 90. Department of Cardiovascular Sciences, University of Leicester, Cardiovascular Research Centre,
165 Glenfield Hospital, Leicester, LE3 9QP, UK
166 91. National Institute on Aging, National Institutes of Health
167 92. Department of Psychology, University of Amsterdam & Amsterdam Brain and Cognition, University of
168 Amsterdam
169 93. Alzheimer Scotland Research Centre, University of Edinburgh, Edinburgh, UK, EH8 9JZ
170 94. Department of Haematology, University of Cambridge, Cambridge, UK, CB2 0PT
171 95. Department of Preventative Medicine, Keck School of Medicine, University of Southern California
172 96. Wellcome Trust Sanger Institute, Hinxton, Cambridge, CB10 1SA, UK
173 97. Department of Medicine, Vanderbilt University, Nashville, TN
174 98. Department of Cardiology, Division Heart and Lungs, University Medical Center Utrecht, Utrecht
175 University, 3508GA Utrecht, the Netherlands
176 99. Department of gerontology and geriatrics, Leiden University Medical Center, Leiden, The Netherlands
177 100. Department of cardiology, Leiden University Medical Center, Leiden, The Netherlands
178 101. Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, 301 Binney Street,
179 Cambridge, MA 02142, USA
180 102. Institute for Community Medicine, University Medicine Greifswald, 17475 Greifswald
181 103. Clinical Pharmacology, William Harvey Research Institute, Queen Mary University of London, London,
182 EC1M 6BQ, UK
183 104. NIHR Barts Cardiovascular Biomedical Research Unit, Queen Mary University of London, London,

184 EC1M 6BQ, UK
185 105. Interfaculty Institute for Genetics and Functional Genomics; University Medicine and Ernst-Moritz-
186 Arndt-University Greifswald, 17475 Greifswald, Germany
187 106. Department of Neuroscience, Psychology and Behaviour, University of Leicester, Leicester, UK
188 107. Department of Biomedical Engineering, The Pennsylvania State University, University Park 16802, USA
189 108. Institute of Personalized Medicine, Penn State College of Medicine
190 109. Department of Clinical Science, Center for Genetics of Host Defense, University of Texas Southwestern
191 110. Department of Cardiology, Ealing Hospital, London North West Healthcare NHS Trust, Middlesex UB1
192 3HW, UK
193 111. Department of Biostatistics and Epidemiology, University of Pennsylvania, USA
194 112. School of Medicine, University of Dundee, Dundee, UK
195 113. Department of Epidemiology and Public Health, Institut Pasteur de Lille, Lille, France
196 114. Department of General and Interventional Cardiology, University Heart Center Hamburg, Germany
197 115. University Medical Center Hamburg Eppendorf
198 116. William Harvey Research Institute, Barts and The London School of Medicine and Dentistry, Queen
199 Mary University of London, EC1M 6BQ UK
200 117. Princess Al-Jawhara Al-Brahim Centre of Excellence in Research of Hereditary Disorders (PACER-HD),
201 King Abdulaziz University, Jeddah 21589, Saudi Arabia
202 118. NIHR Blood and Transplant Research Unit in Donor Health and Genomics, Department of Public Health
203 and Primary Care, University of Cambridge, Cambridge CB1 8RN, UK
204 119. EPIMED Research Centre, Department of Medicine and Surgery, University of Insubria at Varese, Italy
205 120. Department of Epidemiology, UMR 1027- INSERM, Toulouse University-CHU Toulouse, Toulouse,
206 France
207 121. Department of Nutrition, Harvard T. H. Chan School of Public Health, Boston, MA 02115, USA
208 122. Division of Respiratory Medicine, University of Nottingham, Nottingham, UK
209 123. Department of Public Health and Clinical Medicine, Skellefteå Research Unit, Umeå University, Sweden
210 124. Department of Cardiology, Leiden University Medical Center, Leiden, The Netherlands
211 125. The Interuniversity Cardiology Institute of the Netherlands, Utrecht, The Netherlands
212 126. UKCRC Centre of Excellence for Public Health, Queens, University, Belfast
213 127. Medical Research Institute, University of Dundee, Ninewells Hospital and Medical School, Dundee, UK
214 128. University of Glasgow, Glasgow, UK
215 129. Department of Genetic Epidemiology, Kings College, London
216 130. Population Health Research Institute, St George's, University of London, London SW17 0RE, UK
217 131. University of Groningen, University Medical Center Groningen, Department of Genetics, Groningen,
218 The Netherlands
219 132. Department of Biostatistics and Epidemiology, Perelman School of Medicine, University of
220 Pennsylvania, USA

221 133. Center for Non-Communicable Diseases, Karachi, Pakistan
222 134. Department of Public Health and Primary Care, University of Cambridge, UK
223 135. National Institute for Health Research Leicester Respiratory Biomedical Research Centre, Glenfield
224 Hospital, Leicester, UK
225
226
227 *Indicates that these authors contributed equally to this work and share the first author position
228 † Indicates that these authors contributed equally to this work and share the last author position
229 # Correspondence: Joanna M M Howson, jmmh2@medschl.cam.ac.uk and Dajiang J Liu dxl46@psu.edu

230 **Abstract**

231 Smoking is a major heritable and modifiable risk factor for many diseases, including cancer, common
232 respiratory disorders and cardiovascular diseases. Fourteen genetic loci have previously been associated with
233 smoking behaviour-related traits. We tested up to 235,116 single nucleotide variants (SNVs) on the Exome-
234 array for association with smoking initiation, cigarettes per day, pack-years, and smoking cessation in a fixed
235 effects meta-analysis of up to 61 studies (346,813 participants). In a subset of 112,811 participants, a further
236 one million SNVs were also genotyped and tested for association with the four smoking behaviour traits.
237 SNV-trait associations with $P < 5 \times 10^{-8}$ in either analysis were taken forward for replication in up to 275,596
238 independent participants from UK Biobank. Lastly, a meta-analysis of the discovery and replication studies
239 was performed.

240 Sixteen SNVs were associated with at least one of the smoking behaviour traits ($P < 5 \times 10^{-8}$) in the discovery
241 samples. Ten novel SNVs, including rs12616219 near *TMEM182*, were followed-up and five of them
242 (rs462779 in *REV3L*, rs12780116 in *CNNM2*, rs1190736 in *GPR101*, rs11539157 in *PJAI*, and rs12616219
243 near *TMEM182*) replicated at a Bonferroni significance threshold ($P < 4.5 \times 10^{-3}$) with consistent direction of
244 effect. A further 35 SNVs were associated with smoking behaviour traits in the discovery plus replication
245 meta-analysis (up to 622,409 participants) including a rare SNV, rs150493199, in *CCDC141* and two low-
246 frequency SNVs in *CEP350* and *HDGFRP2*. Functional follow-up implied that decreased expression of
247 *REV3L* may lower the probability of smoking initiation. The novel loci will facilitate understanding the
248 genetic aetiology of smoking behaviour and may lead to identification of potential drug targets for smoking
249 prevention and/or cessation.

250

251 **Introduction**

252 Smoking is a major risk factor for many diseases, including common respiratory disorders such as chronic
253 obstructive pulmonary disease (COPD)^{1,2}, cancer³ and cardiovascular diseases⁴, and is reported to cause 1 in
254 10 premature deaths worldwide⁵. A greater understanding of the genetic aetiology of smoking behaviour has
255 the potential to lead to new therapeutic interventions to aid smoking prevention and cessation, and thereby
256 reduce the global burden of such diseases.

257 Previous genome-wide association studies (GWASs) identified 14 common SNVs^{1,6-12} (with minor allele
258 frequency, MAF>0.01) robustly associated with smoking behaviour related traits ($P<5\times 10^{-8}$). The 15q25
259 (*CHRNA3/5-CHRNA4*) region has the largest effect, explaining ~1% and 4-5% of the phenotypic variance of
260 smoking quantity¹³ and cotinine, a biomarker of nicotine intake¹⁴, respectively. Overall, genetic loci identified
261 to date explain ~2% of the estimated genetic heritability of smoking behaviour⁶, which is reported to be
262 between 40-60%¹⁵⁻¹⁷. A recent study suggested that an important proportion (~3.3%) of the phenotypic
263 variance of smoking behaviour related traits was explained by rare nonsynonymous variants (MAF<0.01)¹⁸.
264 Hence, well-powered studies of rare variants are needed.

265 To investigate the effect of rare coding variants on smoking behaviour, we studied 346,813 participants (of
266 which 324,851 were of European ancestry) from 62 cohorts (**Supp. Tables 1 and 2**) at up to 235,116 SNVs
267 from the exome array. As we had access to UK Biobank, we also interrogated SNVs present on the UK
268 Biobank and UK BiLEVE Axiom arrays to identify additional associations across the genome beyond the
269 exome array. To our knowledge, these datasets are an order of magnitude larger than the previous studies⁶,
270 and constitute the most powerful exome-array study of smoking behaviour to date.

271 **Materials and Methods**

272 **Participants**

273 Our study combined study-level summary association data from up to 59 studies of European ancestry and
274 two studies of South Asian ancestry from three consortia (CGSB (Consortium for Genetics of Smoking
275 Behaviour), GWAS & Sequencing Consortium of Alcohol and Nicotine use (GSCAN) and the Coronary
276 Heart Disease (CHD) Exome+ consortium) INTERVAL and UK Biobank. In total, up to 324,851 individuals
277 of European ancestry and 21,962 South Asian individuals were analysed in the discovery stage (**Figure 1**).
278 Further information about the participating cohorts and consortia is given in **Supp. Table 1** and the Supp.
279 Material. All participants provided written informed consent and studies were approved by local Research
280 Ethics Committees and/or Institutional Review boards.

281 **Phenotypes**

282 We chose to analyse the following four smoking behaviour related traits because of their broad availability in
283 existing epidemiological and medical studies, as well as their biological relevance for addiction behaviours:

- 284 i) Smoking initiation (binary trait: ever vs never smokers). Ever smokers were defined as
285 individuals who have smoked >99 cigarettes in their lifetime, which is consistent with the
286 definition by the Centre for Disease Control¹⁹;
- 287 ii) Cigarettes per day (CPD; quantitative trait: average number of cigarettes smoked per day by ever
288 smokers);
- 289 iii) Pack-years (quantitative trait; Packs per day x Years smoked, with a pack defined as 20
290 cigarettes); years smoked is typically formed from age at smoking commencement to current age
291 for current smokers or age at cessation for former smokers.
- 292 iv) Smoking cessation (binary trait: former vs current smokers).

293 In UK Biobank, phenotypes were defined using phenotype codes 1239, 1249, and 2644 for smoking initiation
294 and smoking cessation, and 1239, 3436, 3456 for CPD and pack-years. CPD was inverse normal transformed
295 in the CHD Exome+, INTERVAL and CGSB studies and categorised (1-10, 11-20, 21-30, and 31+ CPD) by
296 the GSCAN studies and UK Biobank (**Supp. Table 2**). All studies performed an inverse normal
297 transformation of pack-years. Summary statistics of study level phenotype distributions are provided in **Supp.**
298 **Table 1.**

299 **Genotyping and quality control**

300 Fifty-nine cohorts were genotyped using exome arrays (up to 235,116 SNVs) and two (UK Biobank and
301 INTERVAL) were genotyped using Axiom Biobank Arrays (up to 820,000 SNVs; **Supp. Table 2**). In total,
302 ~1.06M SNVs were analysed including ~64,000 SNVs on both the Axiom and Exome Arrays. Furthermore,
303 two studies (NAGOZALC and GFG) genotyped their participants using arrays with custom content,
304 increasing the total number of variants analysed to 1,207,583 SNVs. Individual studies performed quality
305 control (QC; **Supp. Material, Supp. Table 2**) and additional QC was conducted centrally (i) to ensure alleles
306 were consistently aligned, (ii) that there were no major sample overlaps between contributing studies, and (iii)
307 variants conformed to Hardy-Weinberg equilibrium and call rate thresholds. We also examined the
308 distribution of the effect sizes and test statistics across cohorts to ensure the test statistics were well-calibrated.

309 **Study level analyses**

310 Each study (including the case-cohort studies²⁰) undertook analyses of up to four smoking traits using
311 RAREMETALWORKER²¹ or RVTESTS²² (**Supp. Table 2**), which generated single variant score statistics
312 and their covariance matrices within sliding windows of 1Mb. CPD and pack-years were analysed using linear
313 models or linear mixed models. Smoking initiation and smoking cessation were analysed using logistic

314 models or linear mixed models. All studies adjusted each trait for age, sex, at least three genetic principal
315 components and any study-specific covariates (**Supp. Table 2**). Chromosome X variants were analysed using
316 the above described approach, but coding males as 0/2. This coding scheme ensures that on average females
317 and males have equal dosages and so is optimal for genes that are inactivated (due to X chromosome
318 inactivation) and is valid for genes that do not undergo X chromosome activation. Males and females were
319 analysed together adjusting for sex as a covariate.

320 **Single variant meta-analyses**

321 Fixed effects meta-analyses across the individual contributing studies of single variant associations were
322 undertaken using the Cochran-Mantel-Haenszel method in RAREMETAL. Z-score statistics were used in the
323 meta-analysis to ensure that the association results are robust against potentially different units of
324 measurement in the phenotype definitions across studies²³. We performed genomic control correction on the
325 meta-analysis results. Variants with $P < 1 \times 10^{-6}$ in tests of heterogeneity were excluded. Variants with $P \leq 5 \times 10^{-8}$
326 were taken forward for replication. In addition, rs12616219 was also taken forward for replication as its P -
327 value was very close to this threshold (smoking initiation, $P = 5.49 \times 10^{-8}$). None of the rare SNVs were genome-
328 wide significant, therefore we also took forward the rare variant with the smallest association P -value,
329 rs141611945 ($P = 2.95 \times 10^{-7}$; MAF < 0.0001).

330 **Replication and combined meta-analysis of discovery and replication data**

331 As UK biobank genetic data were released in two phases, we took the opportunity to replicate findings from
332 the discovery stage in a further 275,596 individuals made available in the phase two release of UK Biobank
333 genetic data. To avoid potential relatedness between discovery and replication samples, the replication
334 samples were screened and individuals with relatedness closer than second degree with the discovery sample
335 in the UK Biobank were removed²⁴. Phenotypes were defined in the same way as the discovery samples
336 (described above). Since the exome array and the UK Biobank Axiom arrays do not fully overlap, we used
337 both genotyped exome variants (approx. 64,000) as well as the additional ~90,000 well imputed exome array
338 variants from UK Biobank (imputation quality score > 0.3) for replication of single variant and gene-based
339 tests. The rare *ATF6* variant was absent from the UK Biobank array and is more prevalent in Africans
340 (MAF = 0.01) than Europeans (MAF = 0.0007). Therefore, replication was sought in 1,437 individuals of
341 African American-ancestry from the HRS and COGA studies. Analysis methods for replication cohorts were
342 the same as for discovery cohorts, including methods to analyse chromosome X (**Supp. Table 2**). The criteria
343 set for the replication were (i) the same direction of effect as the discovery analysis and (ii) $P \leq 0.0045$ in the
344 replication studies (i.e. Bonferroni-adjusted for eleven SNVs at $\alpha = 0.05$).

345 Finally, in order to fully utilise all available data, we carried out a combined meta-analysis of the discovery
346 and replication samples across the exome array content using the same protocols mentioned above.

347 **Conditional analyses**

348 To identify conditionally independent variants associations within previously reported and novel loci a
349 sequential forward stepwise selection was performed²⁵. A 1MB region was defined around the reported or
350 novel sentinel variant (500kb either side) and conditional analyses performed with all variants within the
351 region. If a conditionally independent variant was identified, ($P < 5 \times 10^{-6}$; Bonferroni adjusted for ~10,000
352 independent variants in the test region) the analysis was repeated conditioning on both the most significant
353 conditionally independent variant and the sentinel variant. This stepwise approach was repeated (conditioning
354 on the variants identified in current and earlier iterations) until there were no variants remaining in the region
355 that were conditionally independent. The same protocol was followed for the novel SNVs identified in this
356 study.

357 **Gene-based analyses**

358 For discovery gene-based meta-analyses, we utilised three statistical methods as part of the RAREMETAL
359 package: the Weighted Sum Test (WST)²⁶, the burden test²⁷ and the Sequence Kernel Association test
360 (SKAT)²⁸. EFACTS (v.3.3.0)²⁹ was used to annotate variants (for use in gene-based meta-analyses), as
361 recommended by RAREMETAL. Two MAF cut-offs were used, one used low frequency (MAF < 0.05) and
362 rare variants, the second only used rare variants (MAF < 0.01). Nonsynonymous, stop gain, splice site, start
363 gain, start loss, stop loss, and synonymous variants were selected for inclusion. A sensitivity analysis to
364 exclusion of synonymous variants was also performed. Gene-level associations with $P < 8 \times 10^{-7}$ were deemed
365 statistically significant (Bonferroni-adjusted for ~20,000 genes and three tests at $\alpha = 0.05$). To examine if the
366 gene associations were driven by a single variant, the gene tests were conducted conditional on the SNV with
367 the smallest P -value in the gene, using the shared single variant association statistic and covariance matrices²¹.
368 ²⁵.

369 **Mendelian Randomization analyses**

370 To evaluate the causal effect of SI and CPD on BMI, schizophrenia and educational attainment (EA), we
371 conducted Mendelian randomization (MR) analyses using three complementary approaches available in MR-
372 Base³⁰: inverse variance weighted regression³¹, MR-Egger^{32, 33}, and weighted median³⁴. We used both the
373 previously reported smoking associated SNVs and the SNVs from the current report (as provided in **Tables 1-**
374 **3 and Supp. Table 3**) as instrumental variables. The BMI³⁵, schizophrenia³⁶ and educational attainment³⁷ data
375 came from previously published publicly available data. To assess possible reverse causation, we also used

376 outcome associated SNVs as instrumental variables and conducted MR analyses using SI and CPD as
377 outcome. We considered $P < 0.05/3 = 0.017$ as statistically significant (Bonferroni adjusted for three traits).

378 ***In silico* functional follow up of associated SNVs**

379 To identify whether the (replicated) SNVs identified here affected other traits, we queried the GWAS
380 Catalog³⁸ (version: e91/28/02/2018, downloaded on 01/03/18) for genome-wide significant ($P < 5 \times 10^{-8}$)
381 associations using all proxy SNVs ($r^2 \geq 0.8$) within 2Mb of the top variant in our study.

382 eQTL lookups were carried out in the 13 brain tissues available in GTEx V7³⁹, Brain xQTL (dorsolateral
383 prefrontal cortex)⁴⁰ and BRAINEAC⁴¹ databases, all of which had undergone QC by the individual studies.
384 We did not perform additional QC on these data. In brief, GTEx used Storey's q-value method to correct the
385 FDR for testing multiple transcripts based upon the empirical P -values for the most significant SNV for each
386 transcript⁴³. BRAINEAC calculated the number of tests per transcript and used Benjamini-Hochberg
387 procedure to calculate FDR per transcript using a FDR < 1% as significant. BRAINxQTL used $P < 8 \times 10^{-8}$ as a
388 cut-off for significance for any given transcript. SNVs that met the study specific significance and FDR
389 thresholds, which were in LD ($r^2 > 0.8$ in 1000 Genomes Europeans) with the top eQTL or the sentinel eQTL
390 for a given tissue/transcript combination were considered significant. The genes implicated by these eQTL
391 databases and/or coding changes (e.g. missense and nonsense SNVs) were put into ConsensusPathDB⁴⁴ to
392 identify whether these genes were over-represented in any known biological pathways. Replicated missense
393 SNVs were also put into PolyPhen-2⁴⁵ and FATHMM (unweighted)⁴⁶ to obtain variant effect prediction.

394 **Results**

395 **Single variant associations**

396 In the discovery meta-analyses, we identified 15 common SNVs that were genome-wide significant ($P < 5 \times 10^{-8}$)
397 for one or more of the smoking behaviour traits, of which 9 were novel (**Table 1, Supp. Table 3**). Seven
398 novel loci were identified for smoking initiation, one for both CPD and pack-years and one for smoking
399 cessation (**Figures 1, 2, Table 1 and Supp. Figure 1**). Results for the significant loci were consistent across
400 participating cohorts and there was at least nominal evidence of association ($P < 0.05$) at the novel loci within
401 each of the contributing consortia (**Supp. Table 4**). Full association results for all novel SNVs across the four
402 traits are provided in **Supp. Table 5**. No rare variants were genome-wide significant; the rare variant with the
403 smallest P -value was a missense variant in *ATF6*, rs141611945 (MAF < 0.0001, CPD $P = 2.95 \times 10^{-7}$).

404 Eleven SNVs (including rs12616219 near *TMEM182* with $P = 5.49 \times 10^{-8}$, and the rare variant, rs141611945)
405 were taken forward for replication in independent samples (**Table 1**). The latest release of European UK
406 Biobank individuals not included in the discovery stage (smoking initiation, $n = 275,596$; smoking cessation
407 $n = 123,851$; CPD $n = 80,015$; pack-years $n = 78,897$), was used for replication of the common variants (**Figure**

408 1). Five of the common variants replicated (four for smoking initiation and one with CPD and pack-years) at
409 $P < 0.0045$. Two coding variants (rs11539157, rs1190736) were predicted to be ‘probably damaging’ by
410 PolyPhen-2 and FATHMM. The remaining five SNVs were at least nominally associated ($P < 0.01$) in the
411 replication samples and had consistent direction of effect across discovery and replication. Replication for the
412 rare variant rs141611945 could not be carried out in UK Biobank as the SNV nor its proxies ($r^2 > 0.3$) were
413 available. Thus we initiated replication in African American samples of the COGA ($n=476$) and HRS ($n=961$)
414 cohorts (overall MAF ≈ 0.01). The direction of effect was consistent in the two replication cohorts and
415 consistent with the discovery meta-analysis but a meta-analysis of the two replication cohorts yielded a
416 $P=0.28$. Further data are required to replicate this association.

417 We also performed a meta-analysis combining the discovery and replication samples (up to 622,409
418 individuals). LD score regression showed that the λ (intercept) for all traits was ~ 1.00 , which indicated that
419 confounding factors inflating the results was not an issue^{47, 48}. The combined analysis identified 35 additional
420 novel SNV-smoking trait associations, 33 with smoking initiation, one with CPD and one with smoking
421 cessation at $P < 5 \times 10^{-8}$ (**Table 2**). We note that among our four SNVs that did not replicate, rs216195 (in
422 *SMG6*) was genome-wide significant in the combined meta-analysis of discovery and replication studies
423 ($P=2.41 \times 10^{-9}$; **Table 2**).

424 We also calculated the phenotypic variance explained for novel and known variants. Results can be found in
425 the ‘Calculation of Phenotypic Variance Explained’ section in the **Supplementary Material**.

426 **Associations at known smoking behaviour loci**

427 We assessed evidence for associations at the 14 SNVs previously reported for smoking behaviour-related
428 traits. Seven were genotyped on the exome array and proxies ($r^2 > 0.3$; $\pm 2\text{Mb}$) were identified for the remaining
429 seven (**Supp. Table 3**). All showed nominal evidence of association at $P < 0.05$ and six of these were genome-
430 wide significant in the meta-analysis of the trait for which it was previously reported (**Supp. Table 3 and 5**).

431 Conditional analyses identified five independent associations within three previously reported loci and all five
432 replicated (**Table 3**). At the 19q13 (*RAB4B*) locus, there were three variants in or near *CYP2A6* associated
433 with CPD independently of the established variant (rs7937) and each other: rs8102683 (conditional
434 $P=4.53 \times 10^{-16}$), rs28399442 (conditional $P=2.63 \times 10^{-12}$) and rs3865453 (conditional $P=4.96 \times 10^{-10}$) and
435 rs28399442 was a low frequency variant. The same SNVs also showed evidence of independent effects with
436 pack-years, albeit with larger P -values ($P < 5 \times 10^{-6}$; **Supp. Table 5**). At the *TEX41/PABPC1P2* locus,
437 rs11694518 (conditional $P=3.43 \times 10^{-7}$) was associated with smoking initiation independently of the established
438 variant (rs10427255). At 15q25, rs938682 ($P=7.78 \times 10^{-21}$) was associated with CPD independently of the

439 established variant (rs1051730) and (in agreement with a previous report⁴⁹) is an eQTL for *CHRNA5* in brain
440 putamen basal ganglia tissues in GTEx.

441 **Gene-based association studies**

442 Gene-based collapsing tests using MAF<0.01 variants, did not identify any associated genes at the pre-
443 specified $P<8\times 10^{-7}$ threshold. Of the top four gene associations, three were novel (*CHRNA2*, *MMP17*, and
444 *CRCP*) and one was known (*CHRNA5*), and had $P<7\times 10^{-4}$, with CPD and/or pack-years (**Supp. Table 6**).
445 Analyses conditional on the variant with the smallest P -value in the gene, revealed the associations at
446 *CHRNA2*, *MMP17* and *CRCP* were due to more than one rare variant (conditional $P<0.05$; **Supp. Table 6**).
447 In contrast, the *CHRNA5* gene association was attributable to a single variant (rs2229961).

448 **Mendelian Randomization analyses**

449 We conducted MR analyses to elucidate the potential causal impact of SI and CPD on BMI, schizophrenia and
450 EA using the MR-Egger, median weighted and inverse variance weighted methods. We found a causal
451 association between SI and EA using both the median weighted and inverse variance weighted methods
452 ($P<0.0001$; **Supp. Table 7**) but not with MR-Egger ($P=0.2$). There was an association of SI with BMI using
453 MR-Egger only ($P=0.01$; **Supp. Table 7**), but there was evidence of horizontal pleiotropy ($P=0.001$) and no
454 support from the other methods. Similarly, increased CPD was only associated with reduced BMI using the
455 weighted median approach ($P=0.009$) and not the other methods ($P>0.017$). We also tested if schizophrenia,
456 EA or BMI causally influence CPD or SI using SNVs associated with schizophrenia, EA and BMI,
457 respectively, as instrumental variables. No evidence of such reverse causation was found (**Supp. Table 7**).
458 These results were consistent with previous analyses⁵⁰. There was no evidence of a causal effect of SI on
459 schizophrenia, or CPD on educational attainment (**Supp. Table 7**).

460 **Functional characterization of novel loci**

461 Using proxies with $r^2\geq 0.8$ in 1000 Genomes Europeans, we queried the GWAS catalogue³⁸ ($P\leq 5\times 10^{-8}$) for
462 pleiotropic effects of our novel sentinel SNVs. Two, rs11539157 and rs3001723 were previously associated
463 with schizophrenia³⁶, suggesting shared biological pathways between schizophrenia and smoking behaviours
464 (**Table 2**). This fits with the known association of smoking with schizophrenia⁵¹. Two, rs1514175 and
465 rs2947411 have previously been associated with BMI⁵², and extreme obesity⁵³.

466 eQTL lookups in GTEx V7 (13 Brain tissues with ≥ 80 samples)³⁹, Brain xQTL⁴⁰ and BRAINEAC⁴¹ databases
467 revealed that the A allele at rs462779, which decreases risk of smoking initiation, also decreased expression of
468 *REV3L* in cerebellum in GTEx (A allele $P=4.8\times 10^{-8}$; $\beta=-0.40$) and was in strong LD with the top eQTL for
469 *REV3L* in cerebellum ($r^2=0.86$ with rs9487668 in 1000 Genomes Europeans). The smoking initiation-

470 associated SNV, rs12780116, was an eQTL for *BORCS7* in four brain tissues, and *NT5C2* in the cerebellar
471 hemisphere (A allele $P=4.5 \times 10^{-7}$; $\beta=-0.32$) and the cerebellum ($P=5.6 \times 10^{-6}$; $\beta=-0.415$; in strong LD with the
472 top eQTL, $r^2=0.97$ with rs11191546). The G allele of a second variant in the region, rs7096169 (intronic to
473 *BORCS7* and only in weak LD with rs12780116, $r^2=0.18$ in 1000G Europeans) increases smoking initiation
474 and reduces expression of *BORCS7* and *AS3MT* in eight brain tissues (including dorsolateral prefrontal cortex
475 in the Brain xQTL and was the top *BORCS7* eSNP in GTEx in the Cerebellar Hemisphere, Cerebellum, and
476 Spinal cord cervical-C1). The same variant also reduced expression of *ARL3* in cerebellum in GTEx (**Table**
477 **2**).

478 Biological pathway enrichment analyses carried out in ConsensusPathDB⁴⁴ using the genes implicated by the
479 eQTL databases (**Table 2**) and/or a coding SNVs (i.e. *PJAI*, *GPR101*) showed that the (i) pyrimidine
480 metabolism and (ii) activation of nicotinic acetylcholine receptors pathways are enriched for these smoking
481 behaviour associated genes (false discovery rate <0.01 ; $P<0.0001$).

482 **Discussion**

483 Smoking is the most important preventable lifestyle risk factor for many diseases, including cancers^{3,54}, heart
484 disease^{4,55} and many respiratory diseases such as COPD^{1,2}. Not initiating is the best way to prevent smoking-
485 related diseases and genetics can play a considerable part in smoking behaviours including initiation. We have
486 performed the largest exome-wide genetic association study of smoking behaviour-related traits to date
487 involving up to 622,409 individuals, and identified and replicated five associations, including two on the X-
488 chromosome (**Table 1**). We identified a further 35 novel associations in a meta-analysis of discovery and
489 replication cohorts (**Table 2**). We validated 14 previously reported SNV-smoking trait associations (**Supp.**
490 **Table 3**) and identified secondary independent associations at three loci, including three in the 19q13 region
491 (rs8102683, rs28399442, and rs3865453; **Table 3**).

492 Gene-based tests improve power by aggregating effects of rare variants. While no genes reached our
493 Bonferroni-adjusted P -value threshold, we identified three candidate genes with multiple rare variant
494 associations for future replication: calcitonin gene-related peptide-receptor component (*CRCP*) with CPD and
495 *CHRNA2* and *MMP17* with pack-years (**Supp. Table 6**; also see ‘Genes of Interest’ section in **Supp.**
496 **Material**). *CRCP*’s protein product is expressed in brain tissues amongst others and functions as part of a
497 receptor complex for a neuropeptide that increases intracellular cyclic adenosine monophosphate levels⁵⁶.
498 *MMP17* encodes a matrix metalloproteinase that is also expressed in the brain and is a member of the
499 peptidase M10 family, and proteins in this family are involved in the breakdown of extracellular matrix in
500 normal physiological processes⁵⁷. Given, we were not able conclusively to identify rare variant associations,
501 even larger studies, are required to identify rare variants associated with smoking behaviours. In addition,
502 phenotypes such as cotinine levels⁵⁸ and nicotine metabolism speed⁵⁹ could be interrogated using methods
503 such as MTAG⁶⁰ to improve power.

504 As recommended by UK Biobank, we analysed UK Biobank samples by adjusting for genotyping array
505 because a subset of (extreme smokers in) UK Biobank were genotyped on a different array (UK BiLEVE).
506 However, this adjustment could potentially introduce collider bias in analyses of smoking traits. Given that the
507 UK BiLEVE study is relatively small compared to the full study, and the genetic effect sizes for smoking
508 associated variants are small, we expect the influence of collider bias to be small⁷⁰. Nevertheless, we
509 performed sensitivity analyses to assess the impact of collider bias. Firstly, we performed a meta-analysis
510 excluding the UK BiLEVE samples, and secondly, we re-analysed UK Biobank without adjusting for
511 genotype array. As expected, the estimated genetic effects from these additional analyses were very similar to
512 our reported results suggesting collider bias is not a concern (**Suppl. Table 8**).

513 Follow-up of the replicated SNVs in the literature and eQTL databases implicated some potentially interesting
514 genes: *NT5C2* is known to hydrolyse purine nucleotides and be involved in maintaining cellular nucleotide
515 balance, and was previously associated with schizophrenia⁶¹. *REV3L*, encodes the catalytic subunit of DNA
516 polymerase ζ (zeta) which is involved in translesion DNA synthesis. Previously, polymorphisms in a
517 microRNA target site of *REV3L* were shown to be associated with lung cancer susceptibility⁶². We showed
518 that decreased expression of *REV3L* may also lower the probability of smoking initiation. The SNV,
519 rs11776293, intronic in *EPHX2*, associated with reduced SI in the combined meta-analysis, and is in LD with
520 rs56372821 ($r^2=0.83$), which is associated with reduced cannabis use disorder⁶³. rs216195 (in *SMG6*) was
521 genome-wide significant in the discovery and the combined meta-analysis. *SMG6* is a plausible candidate
522 gene as it was previously shown to be less methylated in current smokers compared to never smokers⁶⁴. The
523 combined meta-analysis also identified a rare missense variant in *CCDC141*, rs150493199 (MAF<0.01; **Table**
524 **2**). Coding variants in *CCDC141* were previously associated with heart rate⁶⁵ and blood pressure^{66, 67}.

525 Smoking behaviours represent a complex phenotype that are linked to an array of socio-cultural and familial,
526 as well as genetic determinants. Kong *et al.*, recently reported that ‘genetic-nurture’ i.e. effects of non-
527 transmitted parental alleles, affect educational attainment⁶⁸. They also show that there is an effect of
528 educational attainment and genetic nurture on smoking behaviour. Four of our sentinel SNVs (or a strong
529 proxy; $r^2>0.8$) were associated with years of educational attainment³⁷ (rs2292239, rs3001723 ($P<5\times 10^{-8}$),
530 rs9320995 ($P=8.90\times 10^{-7}$), and rs13022438 ($P=3.79\times 10^{-6}$), in agreement with this paradigm and our MR
531 analyses indicated that initiating smoking reduced years in education. Future family studies will be required to
532 disentangle how much of the variance explained in the current analysis is due to direct versus genetic
533 nurturing effects.

534 Our study primarily focused on European ancestry, but we also included two non-European studies but these
535 non-European studies lacked statistical power on their own to identify ancestry specific effects. Therefore, we
536 did not perform ancestry specific meta-analyses. Nevertheless, our results offered cross ancestry replication.
537 One of the associations identified in the conditional analyses, rs8102683 (near *CYP2A6*), confirmed an

538 association with CPD that was previously identified by Kumasaka *et al.* in a Japanese population⁶⁹ but this is
539 the first time it was associated in Europeans (rs8102683 is also correlated with rs56113850 ($r^2=0.43$), a SNV
540 identified previously by Loukola *et al.*⁵⁹ in a genetic association study of nicotine metabolite ratio in
541 Europeans). As more non-European studies become available, it would be of great interest to perform non-
542 European ancestry studies, in order to fine-map causal variants for smoking related traits.

543 CPD and pack-years are two correlated measures of smoking. In the ~40,000 individuals from UK Biobank
544 with CPD and pack-years calculated, correlation between CPD and pack-years was 0.640. Interestingly, while
545 pack-years was inversely correlated with smoking cessation (-0.18) i.e. the more years a smoker has been
546 smoking the less likely they were to cease, CPD was positively correlated with smoking cessation (0.13) i.e.
547 heavier smokers were more likely to stop smoking. In contrast, the *DBH* SNV, rs3025343, (first identified via
548 its association with increased smoking cessation⁶) was associated with increased pack-years ($P=1.29 \times 10^{-14}$)
549 and increased CPD ($P=2.93 \times 10^{-9}$) in our study. The association at *DBH* also represents the first time that a
550 SNV has a smaller *P*-value for pack-years (n=131,892) compared to CPD (n=128,746). These findings may
551 help elucidate the genetic basis of these correlated addiction phenotypes.

552 We performed the largest exome-wide genetic association study of smoking behaviour-related traits to date
553 and nearly doubled the number of replicated associations to 24 (including conditional analyses) including
554 associations on the X-chromosome for the first time, which merit further study. We also identified a further 35
555 novel smoking trait associated SNVs in the combined meta-analysis. The novel loci identified in this study
556 will substantially expand our knowledge of the smoking addiction related traits, facilitate understanding the
557 genetic aetiology of smoking behaviour and may lead to identification of drug targets of potential relevance to
558 prevent individuals from initiating smoking and/or aid smokers to stop smoking.

559

560 **Conflict of Interest Statement**

561 Paul W. Franks has been a paid consultant for Eli Lilly and Sanofi Aventis and has received research support
562 from several pharmaceutical companies as part of European Union Innovative Medicines Initiative (IMI)
563 projects. Neil Poulter has received financial support from several pharmaceutical companies that manufacture
564 either blood pressure lowering or lipid lowering agents or both and consultancy fees. Peter Sever has received
565 research awards from Pfizer. Mark J. Caulfield is Chief Scientist for Genomics England, a UK government
566 company.

Figure legends

Figure 1: Study design including discovery and replication stages. NB: Gene-based studies, conditional analyses, and replication in African American ancestry samples not shown here for clarity. *GFG and NAGOZALC studies contributed additional custom content.

Figure 2: A concentric Circos plot of the association results for Smoking Initiation (SI; outer ring), Cigarettes per day (CPD) and Smoking Cessation (SC; inner ring) for chromosomes 1 to 22 (Pack-years results, which can be found in **Supp. Figure 1**, are omitted for clarity). Each dot represents a SNV, with the X and Y axes corresponding to genomic location in Mb and $-\log_{10} P$ -values, respectively. Labels show the nearest gene to the novel sentinel variants identified in the discovery stage and taken forward to replication. The top signals were truncated at 10^{-10} for clarity. Novel and previously reported signals are highlighted in red and dark blue, respectively. Grey rings on the y-axis increase by increments of 2 (initial ring corresponding to $P=0.001$, then 0.00001 etc.); and the outer and inner red rings correspond to the genome-wide significance level ($P=5 \times 10^{-8}$) and $P=5 \times 10^{-7}$, respectively. Image was created using Circos (v0.65).

Tables

Table 1: Association results for SNVs identified in single variant association meta-analyses and taken forward to replication are provided. Novel smoking trait associated SNVs that replicated with $P < 0.005$ and had consistent direction of effect in discovery and replication are highlighted in **bold**. The replication sample size for smoking initiation (SI), CPD, pack-years (PY), and smoking cessation (SC) were 275,596, 80,015, 78,897, and 123,851 respectively. Chromosome (Chr) and position (Pos) for hg19 build 37. EA: Effect allele; OA: other allele; Gene: closest gene; N: number of individuals; EAF: Effect allele frequency in the pooled samples; MAC: Minor allele count; DoE: Direction of effect; SE: Standard error. All SNVs had heterogeneity $P > 0.02$ in the discovery stage. *Replication was sought in 1,437 individuals of African American-ancestry from the HRS and COGA studies; ** The replication-stage beta(se) for the association of rs1190736 with PY in the replication stage was -0.026 (0.0039).

dbSNP ID (Exome ID)	Chr:Pos	EA/OA	Gene	Consequence	Trait	Discovery stage				Replication stage	
						N	EAF	DoE	P-value	Beta (SE)	P-value
rs141611945 (exm118559)	1:161771868	G/A	ATF6	Missense	CPD	128,746	0.0065% MAC=9	+	2.95x10 ⁻⁷	0.184 (0.169)	*P=0.276 in African American samples
rs1190736 ** (exm1659559)	X:136113464	A/C	GPR101	Missense	CPD (PY)	99,037 (96,824)	46.6% (47.0%)	-	1.40x10⁻¹¹ (4.98E-09)	-0.028 (0.0041) -0.027 (0.0049) -0.028 (0.0073)	All samples: 8.20E-12 (2.7E-11) Males only: 1.90E-08 (6.0E-08) Female only: 1.10E-04 (7.1E-04)
rs462779 (exm572256)	6:111695887	A/G	REV3L	Missense	SI	346,682	80.1%	-	4.52x10⁻⁸	-0.023 (0.0034)	9.7E-12
rs216195 (exm1276230)	17:2203167	G/T	SMG6	Missense	SI	335,406	27.3%	-	2.80x10⁻⁸	-0.008 (0.0029)	8.5E-03
rs11539157 (exm1643833)	X:68381264	A/C	PJA1	Missense	SI	289,917	16.5%	+	1.39x10⁻¹¹	0.022 (0.0026) 0.0158 (0.0033) 0.0185 (0.0039)	All samples: 5.40E-17 Males only: 1.30E-06 Females only: 2.20E-06
Non-Exome-chip SNVs											
rs12616219	2:104352495	A/C	TMEM182	Intergenic	SI	112,811	46.4%	-	5.49x10 ⁻⁸	-0.015 (0.0027)	5.5E-08
rs1150691	6:28168033	G/A	ZSCAN9	Missense	SI	112,811	34.8%	-	4.95x10⁻⁸	-0.007 (0.0028)	8.0E-03
rs2841334	9:128122320	A/G	GAPVD1	Intronic	SI	112,811	20.9%	-	2.28x10⁻⁸	-0.009 (0.0033)	7.5E-03
rs202664	22:41813886	C/T	TOB2	Intergenic	SC	51,043	19.9%	-	1.02x10⁻⁸	-0.011 (0.0050)	2.1E-02
rs11895381	2:60053727	A/G	BCL11A	Intergenic	SI	112,811	34.2%	-	5.61x10⁻⁹	-0.007 (0.0028)	1.2E-02
rs12780116	10:104821946	A/G	CNNM2	Intronic	SI	112,811	13.9%	+	9.19x10⁻¹⁰	0.017 (0.0039)	1.1E-05

Table 2: Association results for novel SNVs identified in the combined meta-analysis of the discovery and replication cohorts. Chromosome (Chr) and position (Pos) for each SNV is given for hg19 build 37. Only SNVs reaching genome-wide significance ($P < 5 \times 10^{-8}$) in the combined meta-analysis are shown. Magnitude of the effect size estimates are not presented as traits were transformed in differently by the three consortia analysed. SNVs identified in the discovery stage of this study (see Table 1) are denoted #. The discovery sample size for smoking initiation (SI), CPD, pack-years (PY), and smoking cessation (SC) were 346,813, 128,746, 131,892, and 121,543, respectively; and the replication sample size for SI, CPD, PY, and SC were 275,596, 80,015, 78,897, and 123,851, respectively. NB: rs6673752 (intronic to *UBAP2L*) was not available in the discovery cohorts. EA: Effect allele; OA: other allele. Beta(se): beta and standard error for association in the replication stage. All SNVs had heterogeneity $P > 0.0001$.

dbSNP ID (Exome-chip ID)	Chr:Pos	EA/OA	Gene	Consequence	Trait	EAF	Beta (se) in replication stage	P-value in combined meta- analysis (P-value in Discovery/Replication stage)	Notes
Combining only genotyped Exome-chip content on the Axiom array									
rs1514175	1:74991644	G/A	<i>TNNI3K</i>	Intronic	SI	0.57	-0.011 (0.003)	5.42x10⁻⁹ (9.03x10 ⁻⁵ /1.0x10 ⁻⁵)	Previously associated with BMI
rs7096169	10:104618695	G/A	<i>BORCS7</i> (<i>CNNM2</i> # in Table 1)	Intronic	SI	0.31	0.016 (0.003)	2.17x10⁻¹³ (3.38x10 ⁻⁷ /7.3x10 ⁻⁹)	r ² =0.28 between rs7096169 and rs12780116 (Table 1) in 1000 Genomes EUR. Previously associated with Schizophrenia. rs7096169 an eQTL for <i>ARL3</i> , <i>BORCS7</i> , and <i>AS3MT</i> in ≥1 of the brain tissues in GTEx
rs2292239	12:56482180	G/T	<i>ERBB3</i>	Intronic	SI	0.66	0.0121 (0.003)	2.78x10⁻⁸ (7.56x10 ⁻⁵ /1.5x10 ⁻⁵)	Previously associated with type-1 diabetes and years of educational attainment. rs2292239 is an eQTL for <i>RPS26</i> and <i>SUOX</i> in ≥4 of the brain tissues in GTEx
rs216195	17:2203167	G/T	<i>SMG6</i> #	Missense	SI	0.29	-0.0076 (0.003)	2.41x10⁻⁹ (2.80x10 ⁻⁸ /8.5x10 ⁻³)	Same SNV as in Table 1
Combining well-imputed Exome-chip content on the Axiom array									
rs2960306 (exm383568)	4:2990499	T/G	<i>GRK4</i>	Missense	CPD	0.34	-0.024 (0.005)	1.06x10⁻⁹ (3.99x10 ⁻⁵ /3.8x10 ⁻⁶)	rs2960306 is an eQTL for <i>GRK4</i> in four of the brain tissues in GTEx
rs4908760	1:8526142	A/G	<i>RERE</i>	Intronic	SI	0.35	0.0078 (0.003)	1.76x10⁻⁸ (3.36x10 ⁻⁶ /4.7x10 ⁻³)	Previously associated with Vitiligo
rs6692219 (exm127721)	1:179989584	C/G	<i>CEP350</i>	Missense	SI	0.028	-0.0257 (0.008)	4.69x10⁻⁹ (1.08x10 ⁻⁶ /1.3x10 ⁻³)	
rs11971186	7:126437897	G/A	<i>GRM8</i>	Intronic	SI	0.20	-0.0080 (0.003)	1.45x10⁻⁸ (1.38x10 ⁻⁶ /3.9x10 ⁻³)	
rs150493199 (exm249655)	2:179721072	A/T	<i>CCDC141</i>	Missense	SC	0.0098	0.048 (0.134)	1.28x10⁻⁸ (6.45x10 ⁻⁸ /0.72)	
Non-Exome-chip SNVs									

rs3001723	1:44037685	A/G	<i>PTPRF</i>	Intronic	SI	0.21	0.0159 (0.003)	6.64x10⁻¹¹ (0.00015/4.1x10 ⁻⁸)	Previously associated with Schizophrenia and Years of educational attainment
rs1937455	1:66416939	G/A	<i>PDE4B</i>	Intronic	SI	0.30	-0.0146 (0.0027)	1.23x10⁻⁹ (0.00073/5.6x10 ⁻⁸)	
rs72720396	1:91191582	G/A	<i>BARHL2</i>	Intergenic	SI	0.16	-0.0150 (0.003)	9.86x10⁻⁹ (5.63x10 ⁻⁵ /1.9x10 ⁻⁶)	
rs6673752	1:154219177	C/G	<i>UBAP2L</i>	Intronic	SI	0.055	-0.027 (0.004)	1.1x10⁻¹¹ (NA/1.1x10 ⁻¹¹)	
rs2947411	2:614168	G/A	<i>TMEM18</i>	Intergenic	SI	0.83	0.0189 (0.004)	4.97x10⁻¹⁰ (0.00017/7.1x10 ⁻⁸)	Previously associated with BMI
rs528301	2:45154908	A/G	<i>SIX3</i>	Intergenic	SI	0.38	0.0136 (0.002)	4.12x10⁻¹¹ (1.77x10 ⁻⁶ /3.8x10 ⁻⁷)	
rs6738833	2:104150891	T/C	<i>TMEM182[#]</i>	Intergenic	SI	0.33	-0.018 (0.003)	8.66x10⁻¹⁴ (1.63x10 ⁻⁶ /4.4x10 ⁻¹¹)	r ² =0.69 between rs6738833 and rs12616219 (Table 1) in European samples of the 1000 Genomes Project
rs13026471	2:137564022	T/C	<i>THSD7B</i>	Intronic	SI	0.18	0.0127 (0.003)	2.45x10⁻⁸ (0.00028/3.0x10 ⁻⁵)	
rs6724928	2:156005991	C/T	<i>KCNJ3</i>	Intergenic	SI	0.32	-0.011 (0.003)	4.47x10⁻⁸ (0.0019/4.8x10 ⁻⁵)	
rs13022438	2:162800372	G/A	<i>SLC4A10</i>	Intronic	SI	0.27	0.0146 (0.003)	1.41x10⁻¹¹ (0.0005/8.1x10 ⁻⁸)	
rs1869244	3:5724531	A/G	<i>LOC105376939</i>	Intergenic	SI	0.32	0.0123 (0.003)	2.76x10⁻⁹ (0.00040/4.1x10 ⁻⁶)	
rs35438712	3:85588205	T/C	<i>CADM2</i>	Intronic	SI	0.25	0.017 (0.003)	1.99x10⁻¹³ (1.15x10 ⁻⁵ /3.2x10 ⁻¹⁰)	
rs6883351	5:22193967	T/C	<i>CDH12</i>	Intronic	SI	0.34	0.0129 (0.003)	4.69x10⁻⁸ (0.0010/1.4x10 ⁻⁶)	
rs6414946	5:87729711	C/A	<i>TMEM161B</i>	Intronic	SI	0.32	-0.0137 (0.003)	5.27x10⁻¹⁰ (3.63x10 ⁻⁵ /2.8x10 ⁻⁷)	
rs11747772	5:166992708	C/T	<i>TENM2</i>	Intronic	SI	0.25	0.0144 (0.003)	6.20x10⁻⁹ (0.011/2.2x10 ⁻⁷)	
rs9320995	6:98726381	G/A	<i>POU3F2</i>	Intergenic	SI	0.18	0.0150 (0.003)	1.70x10⁻⁸ (0.00079/6.1x10 ⁻⁷)	
rs10255516	7:1675621	G/A	<i>ELFN1</i>	Intergenic	SI	0.33	-0.0139 (0.003)	2.86x10⁻¹⁰ (0.0021/1.8x10 ⁻⁷)	
rs10807839	7:3344629	G/A	<i>SDK1</i>	Intronic	SI	0.19	0.0162 (0.003)	8.93x10⁻¹¹ (0.0026/4.4x10 ⁻⁸)	
rs6965740	7:117514840	T/G	<i>CTTNBP2</i>	Intergenic	SI	0.31	-0.0126 (0.003)	9.66x10⁻⁹ (5.56x10 ⁻⁶ /2.8x10 ⁻⁶)	
rs11776293	8:27418429	T/C	<i>EPHX2</i>	Intronic	SI	0.12	-0.0200 (0.003)	2.23x10⁻¹² (0.00011/8.9x10 ⁻⁹)	rs11776293 is an eQTL for <i>CHRNA2</i> in cerebellum in GTEx
rs1562612	8:59817068	G/A	<i>TOX</i>	Intronic	SI	0.35	-0.0112 (0.003)	1.15x10⁻⁹ (1.42x10 ⁻⁵ /2.9x10 ⁻⁵)	
rs3857914	8:93184065	C/T	<i>RUNX1T1</i>	Intergenic	SI	0.19	0.0157 (0.003)	1.54x10⁻⁹ (0.065/7.1x10 ⁻⁸)	
rs2799849	9:86752641	C/T	<i>RMI1</i>	Intergenic	SI	0.22	-0.0156 (0.003)	1.94x10⁻⁸ (0.026/4.8x10 ⁻⁸)	
rs6482190	10:22037809	A/G	<i>LOC107984214</i>	Intronic	SI	0.17	0.0146 (0.003)	8.85x10⁻⁹ (0.0021/9.5x10 ⁻⁷)	
rs4523689	11:7950797	G/A	<i>OR10A6</i>	Intergenic	SI	0.27	-0.012 (0.003)	7.77x10⁻⁹ (0.00030/2.2x10 ⁻⁵)	

rs933006	13:38350193	A/G	<i>TRPC4</i>	Intronic	SI	0.32	-0.0143 (0.003)	3.50x10⁻⁸ (0.022/9.6x10 ⁻⁸)
rs557899	15:47643795	A/C	<i>SEMA6D</i>	Intronic	SI	0.26	0.0157 (0.003)	2.99x10⁻¹³ (4.46x10 ⁻⁵ /1.0x10 ⁻⁸)
rs76608582	19:4474725	A/C	<i>HDGFRP2</i>	Intronic	SI	0.029	-0.0360 (0.007)	8.50x10⁻⁹ (0.012/4.3x10 ⁻⁸)

Table 3: Results from conditional analyses at previously reported smoking behaviour loci. SNVs with $P < 5 \times 10^{-8}$ are highlighted in **bold**. The discovery sample size for smoking initiation (SI) and CPD was 346,813 and 128,746, respectively. The replication sample size for SI and CPD were 275,596 and 80,015, respectively. Chr: Chromosome; Pos: position for hg19 build 37; EA: Effect allele; OA: other allele; EAF: Effect allele frequency in the pooled samples; DoE: Direction of effect.

Gene region	dbSNP ID	Chr:Pos	EA/OA	Consequence	Trait	EAF	<i>P</i> (unconditional)	SNV(s) conditioned on	Discovery Conditional <i>P</i> [DoE]	Conditional <i>P</i> in replication [DoE]
19q13 (<i>RAB4B</i>)	rs8102683	19:41363765	C/T	Intergenic	CPD	74.8%	4.53×10^{-16}	rs7937	1.44×10^{-13} [+]	3.5×10^{-4} [+]
	rs28399442	19:41354458	A/C	Intronic (<i>CYP2A6</i>)	CPD	1.3%	2.27×10^{-12}	rs7937, rs8102683	2.63×10^{-12} [+]	8.1×10^{-14} [+]
	rs3865453	19:41338556	T/C	Intergenic	CPD	6.54%	2.96×10^{-12}	rs7937, rs8102683, rs28399442	4.96×10^{-10} [-]	2.3×10^{-13} [-]
<i>TEX41-PABPC1P2</i>	rs11694518	2:146125523	T/C	Intergenic	SI	29.5%	2.90×10^{-9}	rs10193706	3.43×10^{-7} [-]	4.0×10^{-31} [-]
15q25 (<i>CHRNA3</i>)	rs938682	15:78882925	A/G	Intronic (<i>CHRNA3</i>)	CPD	76.4%	1.83×10^{-69}	rs1051730	7.77×10^{-21} [+]	1.0×10^{-13} [+]

References

1. Wain LV, Shrine N, Miller S, Jackson VE, Ntalla I, Soler Artigas M *et al.* Novel insights into the genetics of smoking behaviour, lung function, and chronic obstructive pulmonary disease (UK BiLEVE): a genetic association study in UK Biobank. *Lancet Respir Med* 2015; **3**(10): 769-781.
2. Wain LV, Shrine N, Artigas MS, Erzurumluoglu AM, Noyvert B, Bossini-Castillo L *et al.* Genome-wide association analyses for lung function and chronic obstructive pulmonary disease identify new loci and potential druggable targets. *Nature genetics* 2017; **49**(3): 416-425.
3. McKay JD, Hung RJ, Han Y, Zong X, Carreras-Torres R, Christiani DC *et al.* Large-scale association analysis identifies new lung cancer susceptibility loci and heterogeneity in genetic susceptibility across histological subtypes. *Nature genetics* 2017; **49**(7): 1126-1132.
4. O'Donnell CJ, Nabel EG. Genomics of Cardiovascular Disease. *New England Journal of Medicine* 2011; **365**(22): 2098-2109.
5. Reitsma MB, Fullman N, Ng M, Salama JS, Abajobir A, Abate KH *et al.* Smoking prevalence and attributable disease burden in 195 countries and territories, 1990-2015: a systematic analysis from the Global Burden of Disease Study 2015. *The Lancet* 2017.
6. Tobacco and Genetics Consortium. Genome-wide meta-analyses identify multiple loci associated with smoking behavior. *Nature genetics* 2010; **42**(5): 441-447.
7. Hancock DB, Reginsson GW, Gaddis NC, Chen X, Saccone NL, Lutz SM *et al.* Genome-wide meta-analysis reveals common splice site acceptor variant in CHRNA4 associated with nicotine dependence. *Transl Psychiatry* 2015; **5**: e651.
8. Siedlinski M, Cho MH, Bakke P, Gulsvik A, Lomas DA, Anderson W *et al.* Genome-wide association study of smoking behaviours in patients with COPD. *Thorax* 2011; **66**(10): 894-902.
9. Thorgeirsson TE, Gudbjartsson DF, Surakka I, Vink JM, Amin N, Geller F *et al.* Sequence variants at CHRN3-CHRNA6 and CYP2A6 affect smoking behavior. *Nature genetics* 2010; **42**(5): 448-453.
10. Timofeeva MN, McKay JD, Smith GD, Johansson M, Byrnes GB, Chabrier A *et al.* Genetic polymorphisms in 15q25 and 19q13 loci, cotinine levels, and risk of lung cancer in EPIC. *Cancer Epidemiol Biomarkers Prev* 2011; **20**(10): 2250-2261.
11. Bloom AJ, Baker TB, Chen L-S, Breslau N, Hatsukami D, Bierut LJ *et al.* Variants in two adjacent genes, EGLN2 and CYP2A6, influence smoking behavior related to disease risk via different mechanisms. *Human Molecular Genetics* 2014; **23**(2): 555-561.

12. Thakur GA, Sengupta SM, Grizenko N, Choudhry Z, Joober R. Family-based association study of ADHD and genes increasing the risk for smoking behaviours. *Archives of disease in childhood* 2012; **97**(12): 1027.
13. Munafò MR, Flint J. The genetic architecture of psychophysiological phenotypes. *Psychophysiology* 2014; **51**(12): 1331-1332.
14. Keskitalo K, Broms U, Heliovaara M, Ripatti S, Surakka I, Perola M *et al.* Association of serum cotinine level with a cluster of three nicotinic acetylcholine receptor genes (CHRNA3/CHRNA5/CHRNA4) on chromosome 15. *Hum Mol Genet* 2009; **18**(20): 4007-4012.
15. Vink JM, Willemsen G, Boomsma DI. Heritability of smoking initiation and nicotine dependence. *Behav Genet* 2005; **35**(4): 397-406.
16. Carmelli D, Swan GE, Robinette D, Fabsitz R. Genetic Influence on Smoking — A Study of Male Twins. *New England Journal of Medicine* 1992; **327**(12): 829-833.
17. Kaprio J, Koskenvuo M, Sarna S. Cigarette smoking, use of alcohol, and leisure-time physical activity among same-sexed adult male twins. *Prog Clin Biol Res* 1981; **69 Pt C**: 37-46.
18. Liu DJ, Brazel DM, Turcot V, Zhan X, Gong J, Barnes DR *et al.* Exome chip meta-analysis elucidates the genetic architecture of rare coding variants in smoking and drinking behavior. *bioRxiv* 2017.
19. Centers for Disease Control and Prevention (CDC). Cigarette smoking among adults--United States, 2007. *MMWR Morbidity and mortality weekly report* 2008; **57**(45): 1221-1226.
20. Staley JR, Jones E, Kaptoge S, Butterworth AS, Sweeting MJ, Wood AM *et al.* A comparison of Cox and logistic regression for use in genome-wide association studies of cohort and case-cohort design. *European journal of human genetics : EJHG* 2017; **25**(7): 854-862.
21. Feng S, Liu D, Zhan X, Wing MK, Abecasis GR. RAREMETAL: fast and powerful meta-analysis for rare variants. *Bioinformatics* 2014; **30**(19): 2828-2829.
22. Zhan X, Hu Y, Li B, Abecasis GR, Liu DJ. RVTESTS: an efficient and comprehensive tool for rare variant association analysis using sequence data. *Bioinformatics* 2016; **32**(9): 1423-1426.
23. Willer CJ, Li Y, Abecasis GR. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* 2010; **26**(17): 2190-2191.
24. Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K *et al.* Genome-wide genetic data on ~500,000 UK Biobank participants. *bioRxiv* 2017.

25. Jiang B, Chen S, Jiang Y, Liu M, Iacono WG, Hewitt JK *et al.* Proper Conditional Analysis in the Presence of Missing Data Identified Novel Independently Associated Low Frequency Variants in Nicotine Dependence Genes. *bioRxiv* 2017.
26. Madsen BE, Browning SR. A Groupwise Association Test for Rare Mutations Using a Weighted Sum Statistic. *PLoS genetics* 2009; **5**(2): e1000384.
27. Morris AP, Zeggini E. An evaluation of statistical approaches to rare variant analysis in genetic association studies. *Genet Epidemiol* 2010; **34**(2): 188-193.
28. Wu MC. Rare variant association testing for sequencing data using the sequence kernel association test (SKAT). *Am J Hum Genet* 2011; **89**: 82-93.
29. Zhan X, Liu DJ. SEQMINER: An R-Package to Facilitate the Functional Interpretation of Sequence-Based Associations. *Genet Epidemiol* 2015; **39**(8): 619-623.
30. Hemani G, Zheng J, Elsworth B, Wade KH, Haberland V, Baird D *et al.* The MR-Base platform supports systematic causal inference across the human phenome. *Elife* 2018; **7**.
31. Pierce BL, Burgess S. Efficient design for Mendelian randomization studies: subsample and 2-sample instrumental variable estimators. *American journal of epidemiology* 2013; **178**(7): 1177-1184.
32. Bowden J, Davey Smith G, Burgess S. Mendelian randomization with invalid instruments: effect estimation and bias detection through Egger regression. *International journal of epidemiology* 2015; **44**(2): 512-525.
33. Rees JMB, Wood AM, Burgess S. Extending the MR-Egger method for multivariable Mendelian randomization to correct for both measured and unmeasured pleiotropy. *Stat Med* 2017; **36**(29): 4705-4718.
34. Bowden J, Davey Smith G, Haycock PC, Burgess S. Consistent Estimation in Mendelian Randomization with Some Invalid Instruments Using a Weighted Median Estimator. *Genet Epidemiol* 2016; **40**(4): 304-314.
35. Locke AE, Kahali B, Berndt SI, Justice AE, Pers TH, Day FR *et al.* Genetic studies of body mass index yield new insights for obesity biology. *Nature* 2015; **518**(7538): 197-206.
36. Schizophrenia Working Group of the Psychiatric Genomics Consortium, Ripke S, Neale BM, Corvin A, Walters JTR, Farh K-H *et al.* Biological insights from 108 schizophrenia-associated genetic loci. *Nature* 2014; **511**: 421.
37. Okbay A, Beauchamp JP, Fontana MA, Lee JJ, Pers TH, Rietveld CA *et al.* Genome-wide association study identifies 74 loci associated with educational attainment. *Nature* 2016; **533**(7604): 539-542.

38. MacArthur J, Bowler E, Cerezo M, Gil L, Hall P, Hastings E *et al.* The new NHGRI-EBI Catalog of published genome-wide association studies (GWAS Catalog). *Nucleic acids research* 2017; **45**(D1): D896-D901.
39. Battle A, Brown CD, Engelhardt BE, Montgomery SB. Genetic effects on gene expression across human tissues. *Nature* 2017; **550**(7675): 204-213.
40. Ng B, White CC, Klein H-U, Sieberts SK, McCabe C, Patrick E *et al.* An xQTL map integrates the genetic architecture of the human brain's transcriptome and epigenome. *Nat Neurosci* 2017; **advance online publication**.
41. Trabzuni D, Ryten M, Walker R, Smith C, Imran S, Ramasamy A *et al.* Quality control parameters on a large dataset of regionally dissected human control brains for whole genome expression studies. *Journal of Neurochemistry* 2011; **119**(2): 275-282.
42. Ongen H, Buil A, Brown AA, Dermitzakis ET, Delaneau O. Fast and efficient QTL mapper for thousands of molecular phenotypes. *Bioinformatics* 2016; **32**(10): 1479-1485.
43. Storey JD, Tibshirani R. Statistical significance for genomewide studies. *Proceedings of the National Academy of Sciences of the United States of America* 2003; **100**(16): 9440-9445.
44. Kamburov A, Wierling C, Lehrach H, Herwig R. ConsensusPathDB—a database for integrating human functional interaction networks. *Nucleic acids research* 2009; **37**(suppl_1): D623-D628.
45. Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P *et al.* A method and server for predicting damaging missense mutations. *Nat Meth* 2010; **7**(4): 248-249.
46. Shihab HA, Gough J, Cooper DN, Stenson PD, Barker GL, Edwards KJ *et al.* Predicting the functional, molecular, and phenotypic consequences of amino acid substitutions using hidden Markov models. *Human mutation* 2013; **34**(1): 57-65.
47. Bulik-Sullivan BK, Loh PR, Finucane HK, Ripke S, Yang J, Patterson N *et al.* LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nature genetics* 2015; **47**(3): 291-295.
48. Zheng J, Erzurumluoglu AM, Elsworth BL, Kemp JP, Howe L, Haycock PC *et al.* LD Hub: a centralized database and web interface to perform LD score regression that maximizes the potential of summary level GWAS data for SNP heritability and genetic correlation analysis. *Bioinformatics* 2017; **33**(2): 272-279.

49. Wang JC, Cruchaga C, Saccone NL, Bertelsen S, Liu P, Budde JP *et al.* Risk for nicotine dependence and lung cancer is conferred by mRNA expression levels and amino acid change in CHRNA5. *Hum Mol Genet* 2009; **18**(16): 3125-3135.
50. Gage SH, Jones HJ, Taylor AE, Burgess S, Zammit S, Munafò MR. Investigating causality in associations between smoking initiation and schizophrenia using Mendelian randomization. *Sci Rep* 2017; **7**: 40653.
51. Kelly C, McCreadie R. Cigarette smoking and schizophrenia. *Advances in Psychiatric Treatment* 2000; **6**(5): 327-331.
52. Speliotes EK, Willer CJ, Berndt SI, Monda KL, Thorleifsson G, Jackson AU *et al.* Association analyses of 249,796 individuals reveal 18 new loci associated with body mass index. *Nature genetics* 2010; **42**(11): 937-948.
53. Wheeler E, Huang N, Bochukova EG, Keogh JM, Lindsay S, Garg S *et al.* Genome-wide SNP and CNV analysis identifies common and low-frequency variants associated with severe early-onset obesity. *Nature genetics* 2013; **45**(5): 513-517.
54. Hecht SS. Tobacco Smoke Carcinogens and Lung Cancer. *JNCI: Journal of the National Cancer Institute* 1999; **91**(14): 1194-1210.
55. Ockene IS, Miller NH. Cigarette Smoking, Cardiovascular Disease, and Stroke. *A Statement for Healthcare Professionals From the American Heart Association* 1997; **96**(9): 3243-3247.
56. Uhlen M, Fagerberg L, Hallstrom BM, Lindskog C, Oksvold P, Mardinoglu A *et al.* Proteomics. Tissue-based map of the human proteome. *Science* 2015; **347**(6220): 1260419.
57. O'Leary NA, Wright MW, Brister JR, Ciufo S, Haddad D, McVeigh R *et al.* Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic acids research* 2016; **44**(D1): D733-745.
58. Ware JJ, Chen X, Vink J, Loukola A, Minica C, Pool R *et al.* Genome-Wide Meta-Analysis of Cotinine Levels in Cigarette Smokers Identifies Locus at 4q13.2. *Sci Rep* 2016; **6**: 20092.
59. Loukola A, Buchwald J, Gupta R, Palviainen T, Hallfors J, Tikkanen E *et al.* A Genome-Wide Association Study of a Biomarker of Nicotine Metabolism. *PLoS genetics* 2015; **11**(9): e1005498.
60. Turley P, Walters RK, Maghzian O, Okbay A, Lee JJ, Fontana MA *et al.* Multi-trait analysis of genome-wide association summary statistics using MTAG. *Nature genetics* 2018; **50**(2): 229-237.
61. Aberg KA, Liu Y, Bukszár J, *et al.* A comprehensive family-based replication study of schizophrenia genes. *JAMA Psychiatry* 2013; **70**(6): 573-581.

62. Zhang S, Chen H, Zhao X, Cao J, Tong J, Lu J *et al.* REV3L 3[prime]UTR 460 T>C polymorphism in microRNA target sites contributes to lung cancer susceptibility. *Oncogene* 2013; **32**(2): 242-250.
63. Demontis D, Rajagopal VM, Als TD, Grove J, Pallesen J, Hjorthoj C *et al.* Genome-wide association study implicates *CHRNA2* in cannabis use disorder. *bioRxiv* 2018.
64. Steenaard RV, Ligthart S, Stolk L, Peters MJ, van Meurs JB, Uitterlinden AG *et al.* Tobacco smoking is associated with methylation of genes related to coronary artery disease. *Clin Epigenetics* 2015; **7**: 54.
65. van den Berg ME, Warren HR, Cabrera CP, Verweij N, Mifsud B, Haessler J *et al.* Discovery of novel heart rate-associated loci using the Exome Chip. *Hum Mol Genet* 2017; **26**(12): 2346-2363.
66. Warren HR, Evangelou E, Cabrera CP, Gao H, Ren M, Mifsud B *et al.* Genome-wide association analysis identifies novel blood pressure loci and offers biological insights into cardiovascular risk. *Nature genetics* 2017; **49**(3): 403-415.
67. Hoffmann TJ, Ehret GB, Nandakumar P, Ranatunga D, Schaefer C, Kwok PY *et al.* Genome-wide association analyses using electronic health records identify new loci influencing blood pressure variation. *Nature genetics* 2017; **49**(1): 54-64.
68. Kong A, Thorleifsson G, Frigge ML, Vilhjalmsjon BJ, Young AI, Thorgeirsson TE *et al.* The nature of nurture: Effects of parental genotypes. *Science* 2018; **359**(6374): 424-428.
69. Kumasaka N, Aoki M, Okada Y, Takahashi A, Ozaki K, Mushiroda T *et al.* Haplotypes with Copy Number and Single Nucleotide Polymorphisms in CYP2A6 Locus Are Associated with Smoking Quantity in a Japanese Population. *PloS one* 2012; **7**(9): e44507.
70. Munafo MR, Tilling K, Taylor AE, Evans DM, Davey Smith G. Collider scope: when selection bias can substantially influence observed associations. *International journal of epidemiology* 2018; **47**(1): 226-235.