**Genome-wide haplotype-based association analysis of major depressive disorder in Generation Scotland and UK Biobank**

David M. Howard Ph.D.*[1], Lynsey S. Hall Ph.D. [1], Jonathan D. Hafferty M.D. [1], Yanni Zeng Ph.D. [1,2], Mark J. Adams Ph.D. [1], Toni-Kim Clarke Ph.D. [1], David J. Porteous Ph.D. [3], Reka Nagy BSc[2], Caroline Hayward Ph.D. [2,8], Blair H. Smith Ph.D. [4,8], Alison D. Murray Ph.D. [5,8], Niamh M. Ryan Ph.D. [3], Kathryn L. Evans Ph.D. [3,7], Chris S. Haley Ph.D. [2], Ian J. Deary Ph.D. [6,7,8], Pippa A. Thomson Ph.D. [3,7] and Andrew M. McIntosh M.D.[1,7,8]

Running Title: Haplotype-based association analysis of depression

Affiliations:
[1]Division of Psychiatry, University of Edinburgh, Royal Edinburgh Hospital, Edinburgh, UK
[2]Medical Research Council Human Genetics Unit, Institute of Genetics and Molecular Medicine, University of Edinburgh, Edinburgh, UK
[3]Centre for Genomic and Experimental Medicine, Institute of Genetics and Molecular Medicine, University of Edinburgh, Edinburgh, UK
[4]Division of Population Health Sciences, University of Dundee, Dundee, UK
[5]Aberdeen Biomedical Imaging Centre, University of Aberdeen, Aberdeen, UK
[6]Department of Psychology, The University of Edinburgh, Edinburgh, UK
[7]Centre for Cognitive Ageing and Cognitive Epidemiology, The University of Edinburgh, Edinburgh, UK
[8]Generation Scotland, Institute of Genetics and Molecular Medicine, University of Edinburgh, Edinburgh, UK

*Corresponding author: David M. Howard
Division of Psychiatry, University of Edinburgh, Royal Edinburgh Hospital, Edinburgh, UK
+44 131 537 6268
(e-mail: D.Howard@ed.ac.uk)

Keywords: Haplotype association analysis; Major Depressive Disorder; 6q21; Depression; Generation Scotland; UK Biobank

38 **ABSTRACT**

39 Genome-wide association studies using genotype data have had limited success in the identification of

40 variants associated with major depressive disorder (MDD). Haplotype data provide an alternative

41 method for detecting associations between variants in weak linkage disequilibrium with genotyped

42 variants and a given trait of interest. A genome-wide haplotype association study for MDD was

43 undertaken utilising a family-based population cohort, Generation Scotland: Scottish Family Health

44 Study (n = 18 773), as a discovery cohort with UK Biobank used as a population-based cohort

45 replication cohort (n = 25 035). Fine mapping of haplotype boundaries was used to account for

46 overlapping haplotypes potentially tagging the same causal variant. Within the discovery cohort, two

47 haplotypes exceeded genome-wide significance ($P < 5$ x $10^{-8}$) for an association with MDD. One of

48 these haplotypes was nominally significant in the replication cohort ($P < 0.05$) and was located in

49 6q21, a region which has been previously associated with bipolar disorder, a psychiatric disorder that

50 is phenotypically and genetically correlated with MDD. Several haplotypes with $P < 10^{-7}$ in the

51 discovery cohort were located within gene coding regions associated with diseases that are comorbid

52 with MDD. Using such haplotypes to highlight regions for sequencing may lead to the identification

53 of the underlying causal variants.

54 **INTRODUCTION**

55 Major depressive disorder (MDD) is a complex and clinically heterogeneous condition with core

56 symptoms of low mood and/or anhedonia over a period of at least two weeks. MDD is frequently

57 comorbid with other clinical conditions, such as cardiovascular disease,[1] cancer[2] and inflammatory

58 diseases.[3] This complexity and comorbidity suggests heterogeneity of aetiology and may explain why

59 there has been limited success in identifying causal genetic variants,[4-7] despite heritability estimates

60 ranging from 28% to 37%.[8, 9] Single nucleotide polymorphism (SNP)-based analyses are unlikely to

61 fully capture the variation in regions surrounding the genotyped markers, including untyped lower-

62 frequency variants and those that are in weak linkage disequilibrium (LD) with the common SNPs on

63 many genotyping arrays.

2

64    Haplotype-based analysis may help improve the detection of causal genetic variants as, unlike single

65    SNP-based analysis, it is possible to assign the strand of sequence variants and combine information

66    from multiple SNPs to identify rarer causal variants. A number of studies[10-12] have identified

67    haplotypes associated with MDD, albeit by focussing on particular regions of interest. In the current

68    study, a family and population-based cohort Generation Scotland: Scottish Family Health Study

69    (GS:SFHS) was utilised to ascertain genome-wide haplotypes in closely and distantly related

70    individuals.[13] A haplotype-based association analysis was conducted using MDD as a phenotype,

71    followed by additional fine-mapping of haplotype boundaries with a replication and meta-analysis

72    performed using the UK Biobank cohort.[14]

73    **MATERIALS AND METHODS**

74    Discovery cohort

75    The discovery phase of the study used the family and population-based Generation Scotland: Scottish

76    Family Health Study (GS:SFHS) cohort,[13] consisting of 23 960 individuals of whom 20 195 were

77    genotyped with the Illumina OmniExpress BeadChip (706 786 SNPs). Individuals with a genotype

78    call rate < 98% were removed, as well as those SNPs with a call rate < 98%, a minor allele frequency

79    (MAF) < 0.01 or those deviating from Hardy-Weinberg equilibrium ($P < 10^{-6}$). Individuals who were

80    identified as population outliers through principal component analyses of their genotypic information

81    were also removed.[15]

82    Following quality control there were 19 904 GS:SFHS individuals (11 731 females and 8 173 males)

83    that had genotypic information for 561 125 autosomal SNPs. These individuals ranged from 18-99

84    years of age with an average age of 47.4 years and a standard deviation of 15.0 years. There were 4

85    933 families that had at least two related individuals, this included 1 799 families with two members,

86    1 216 families with three members and 829 families with four members. The largest family group

87    consisted of 31 related individuals and there were 1 789 individuals that had no other family members

88    within GS:SFHS.

89    Replication cohort

3

90   The population-based UK Biobank[16] (provided as part of project #4844) was used as a replication

91   cohort to assess those haplotypes within GS:SFHS with $P < 10^{-6}$. The UK Biobank data consisted of

92   152 249 individuals with genomic data for 72 355 667 imputed variants.[17] The SNPs genotyped in

93   GS:SFHS were extracted from the UK Biobank data and those variants with an imputation infoscore

94   < 0.8 were removed, leaving 555 782 variants in common between the two cohorts. Those genotyped

95   individuals listed as non-white British and those that had also participated in GS:SFHS were removed

96   from within UK Biobank, leaving a total of 119 955 individuals.

97   **Genotype phasing and haplotype formation**

98   The genotype data for GS:SFHS and UK Biobank was phased using SHAPEIT v2.r837.[18] Genome-

99   wide phasing was conducted on the GS:SFHS cohort, whilst the phasing of UK Biobank was

100  conducted on a 50Mb window centred on those haplotypes identified within GS:SFHS with $P < 10^{-6}$.

101  The relatedness within GS:SFHS made it suitable for the application of the duoHMM method, which

102  improves phasing accuracy by also incorporating family information.[19] The default window size of

103  2Mb was used for UK Biobank and a 5Mb window was used for GS:SFHS as larger window sizes

104  have been demonstrated to be beneficial when there is increased identity by descent (IBD) in the

105  population.[18] The number of conditioning states per SNP was increased from the default of 100 states

106  to 200 states to improve phasing accuracy, with the default effective population size of 15 000 used.

107  To calculate the recombination rates between SNPs during phasing the HapMap phase II b37[20] was

108  used. This build was also used to partition the phased data into haplotypes.

109  Three window sizes (1cM, 0.5cM and 0.25cM) were used to establish the SNPs that formed each

110  haplotype.[21] Each window was then moved along the genome by a quarter of the respective window

111  size. There were a total of 97 333 windows with a mean number of SNPs per window of 157, 79 and

112  34 for the 1cM, 0.5cM and 0.25cM windows, respectively. Windows that were less 5 SNPs in length

113  were removed. Within each window, those haplotypes that had a minor allele frequency < 0.005 or

114  that deviating from Hardy-Weinberg equilibrium ($P < 10^{-6}$) were not tested for association. However,

115  they were included within the alternative haplotype when assessing the remaining 2 618 094

116    haplotypes. The reported haplotype positions relate to the outermost SNPs within each haplotype are

117    in base pair (bp) position according to GRCh37.

118    To approximate the number of independently segregating haplotypes the clump command within

119    Plink v1.90 [22] was applied. This provides an estimation of the Bonferroni correction required for

120    multiple testing. When applying an LD $r^2$ threshold of < 0.4 there were 1 070 216 independently

121    segregating haplotypes within GS:SFHS, equating to a $P$-value < 5 x $10^{-8}$ for genome-wide

122    significance. This threshold is also frequently applied to SNP-based and sequence-based association

123    studies to account for multiple testing.[23]

124    **Phenotype ascertainment and patient linkage**

125    Discovery cohort

126    Within GS:SFHS a diagnosis of MDD was made using initial screening questions and the Structured

127    Clinical Interview for the Diagnostic and Statistical Manual of Mental Disorders (SCID).[24] The SCID

128    is an internationally validated approach to identifying episodes of depression and was conducted by

129    clinical nurses trained in its administration. Further details regarding this diagnostic assessment have

130    been described previously.[25] In this study, MDD was defined by at least one instance of a major

131    depressive episode which initially identified 2 659 cases, 17 237 controls and 98 missing (phenotype

132    unknown) individuals.

133    In addition, the psychiatric history of cases and controls was examined using record linkage to the

134    Scottish Morbidity Record.[26] Within the control group, 1 072 participants were found to have attended

135    at least one psychiatry outpatient clinic and were excluded from the study. In addition, 47 of the MDD

136    cases were found to have additional diagnoses of either bipolar disorder or schizophrenia in

137    psychiatric inpatient records and were also excluded from the study. These participants had given

138    prior consent for anonymised record linkage to routine administrative clinical data.

139     In total there were 2 605 MDD cases and 16 168 controls following the removal of individuals based

140     on patient records and population stratification, equating to a prevalence of 13.9% for MDD in this

141     cohort.

142     Replication Cohort

143     Within the UK Biobank cohort, 25 035 participants completed a touchscreen assessment of depressive

144     symptoms and previous treatment. On the basis of their responses, diagnostic status was defined as

145     either 'probable single lifetime episode of major depression' or 'probable recurrent major depression

146     (moderate and severe)' and with control status defined as 'no mood disorder'. In total there were 8

147     508 cases and 16 527 controls, equating to a trait prevalence of 34.0% in this cohort, after the removal

148     of individuals with insufficient information or ambiguous phenotypes.[14]

149     **Statistical approach**

150     Discovery cohort

151     A mixed linear model was used to conduct an association analysis using GCTA v1.25.0 [27]:

152 $$y = X\beta + Z_1 u + Z_2 v + \varepsilon$$

153     where $y$ was the vector of binary observations for MDD. $\beta$ was the matrix of fixed effects, including

154     haplotype, sex, age and age$^2$. $u$ was fitted as a random effect taking into account the genomic

155     relationships (MVN $(0, G\sigma_u^2)$, where $G$ was a SNP-based genomic relationship matrix[28]). $v$ was a

156     random effect fitting a second genomic relationship matrix $G_t$ (MVN $(0, G_t\sigma_v^2)$ which modelled only

157     the more closely related individuals.[29] $G_t$ was equal to $G$ except that off-diagonal elements < 0.05

158     were set to 0. $X$, $Z_1$ and $Z_2$ were the corresponding incidence matrices. $\varepsilon$ was the vector of residual

159     effects and was assumed to be normally distributed, MVN $(0, I\sigma_\varepsilon^2)$.

160    The inclusion of the second genomic relationship matrix, $\boldsymbol{G_t}$, was deemed desirable as the fitting of

161    the single matrix $\boldsymbol{G}$ alone resulted in significant population stratification (intercept = 1.029 ± 0.003,

162    λGC = 1.026) following examination with LD score regression.[30] The fitting of both genomic

163    relationship matrices simultaneously produced no evidence of bias due to population stratification

164    (intercept = 1.002 ± 0.003, λGC = 1.005).

165    Replication cohort

166    A mixed linear model was used to assess the haplotypes in UK Biobank which were identified in the

167    discovery cohort with $P < 10^{-6}$ using GCTA v1.25.0 [27]:

168    $$y = X\beta + Z_1u + \varepsilon$$

169    where $y$ was the vector of binary observations for MDD. $\beta$ was the matrix of fixed effects, including

170    haplotype, sex, age, age$^2$, genotyping batch and recruitment centre. $u$ was fitted as a random effect

171    taking into account the SNP-based genomic relationships (MVN $(0, G\sigma_u^2)$. $X$ and $Z_1$ were the

172    corresponding incidence matrices and $\varepsilon$ was the vector of residual effects and was assumed to be

173    normally distributed, MVN $(0, I\sigma_\varepsilon^2)$. $\sigma_\varepsilon^2$). Replication success was judged on the statistical

174    significance of each haplotype using an inverse variance-weighted meta-analysis across both cohorts

175    conducted using Metal.[31]

176    Fine mapping

177    The method described above examines the effect of each haplotype against all other haplotypes in that

178    window. Therefore, a haplotype could be assessed against similar haplotypes containing the same

179    causal variant, limiting any observed phenotypic association. To investigate whether there were causal

180    variants located within directly overlapping haplotypes of the same window size, fine mapping of

181    haplotype boundaries was used. Where there were directly overlapping haplotypes, each with $P < 10^{-3}$

182    and with an effect in the same direction, i.e. both causal or both preventative, then any shared

183    consecutive regions formed a new haplotype that was assessed using the mixed model described

7

184    previously. This new haplotype was assessed using all individuals and was required to be at least 5

185    SNPs in length. A total of 47 new haplotypes were assessed from within 26 pairs of directly

186    overlapping haplotypes.

187    **RESULTS**

188    An association analysis for MDD was conducted using 2 618 094 haplotypes and 47 fine mapped

189    haplotypes within the discovery cohort, GS:SFHS. A genome-wide Manhattan plot of $-\log_{10} P$-values

190    for these haplotypes is provided in Figure 1 with a q-q plot provided in Supplementary Figure S1.

191    Within the discovery cohort, two haplotypes exceeded genome-wide significance ($P < 5 \times 10^{-8}$) for an

192    association with MDD, one located on chromosome 6 and the other located on chromosome 10. There

193    were 12 haplotypes with $P < 10^{-6}$ in the discovery cohort with replication sought for these haplotypes

194    using UK Biobank. Summary statistics from both cohorts and the meta-analysis for these 12

195    haplotypes are provided in Table 1. The protein coding genes which overlap these 12 haplotypes

196    along with the observed haplotype frequencies within the two cohorts are provided in Table 2. The

197    SNPs and alleles that constitute these 12 haplotypes are provided in Supplementary Table S1.

198    The two haplotypes on chromosome 6 (LD $r^2 = 0.74$) with $P < 10^{-6}$ in the discovery cohort both

199    achieved nominal significance ($P < 0.05$) in the replication cohort, with one reaching genome-wide

200    significance ($P < 5 \times 10^{-8}$) in the meta-analysis. A regional association plot of the region surrounding

201    these haplotypes within GS:SFHS is provided in Figure 2. Fine mapping was used to form the most

202    significant haplotype within the discovery cohort. Two directly overlapping 0.5cM haplotypes

203    consisting of 28 SNPs were identified between 108 335 345 and 108 454 437 bp (rs7749081 -

204    rs212829). These two haplotypes had $P$-values of $3.24 \times 10^{-5}$ and $5.57 \times 10^{-5}$, respectively and differed

205    at a single SNP (rs7749081). Exclusion of this single SNP defined a new 27 SNP haplotype that had a

206    genome-wide significant association with MDD ($P = 7.06 \times 10^{-9}$). Calculating the effect size at the

207    population level,[32] the estimates of the contribution of the two haplotypes to the total genetic variance

208    was $2.09 \times 10^{-4}$ and $2.38 \times 10^{-4}$, respectively, within GS:SFHS. None of the individual SNPs located

209    within either haplotype were associated with MDD in either cohort ($P \geq 0.05$).

8

210    A genome-wide significant haplotype ($P$ = 8.50 x $10^{-9}$) was identified on chromosome 10 within

211    GS:SFHS using a 0.5cM window. A regional association plot of the region surrounding this haplotype

212    is provided in Figure 3. This haplotype had an odds ratio (OR) of 2.33 (95% CI: 1.83 – 2.91) in the

213    discovery cohort and an OR of 1.15 (95% CI: 0.80 - 1.59) in the replication cohort. These were the

214    highest ORs observed in the respective cohorts. The estimate of the contribution of this haplotype to

215    the total genetic variance was 2.29 x $10^{-4}$ in the discovery cohort. Association analysis of the 92 SNPs

216    on this haplotype revealed that one SNP in GS:SFHS (rs17133585) and two SNPs in UK Biobank

217    (rs12413638 and rs10904290) were nominally significant ($P < 0.05$), although none had $P$-values <

218    0.001.

219    All 12 of the haplotypes with a $P$-value for association < $10^{-6}$ in the GS:SFHS discovery cohort were

220    risk factors for MDD (OR > 1) and within the replication cohort, 7 out of these 12 haplotypes had OR

221    > 1. None of the 95% confidence intervals for the replication ORs overlapped the 95% confidence

222    intervals of the discovery GS:SFHS cohort.

223    **DISCUSSION**

224    Twelve haplotypes were identified in the discovery cohort with $P < 10^{-6}$ of which two were significant

225    at the genome-wide level ($P < 5$ x $10^{-8}$) in the discovery cohort and one which was genome-wide

226    significant ($P < 5$ x $10^{-8}$) in the meta-analysis. A power analysis[33] was conducted using the genotype

227    relative risks observed in the discovery cohort, the sample sizes and haplotype frequencies in the

228    replication cohort and the prevalence of MDD reported for a structured clinical diagnosis of MDD in

229    other high income counties (14.6%).[34] There was sufficient power (> 0.99) to detect the twelve

230    haplotypes with $P < 10^{-6}$ identified in the discovery cohort within the replication cohort at a

231    significance threshold of 0.05%.

232    A complementary approach to replication is to identify the gene coding regions within haplotypes that

233    potentially provide a biologically informative explanation for an association with MDD. Those

234    haplotypes with $P < 10^{-7}$ in the discovery cohort and the gene coding regions that they overlap are

235    discussed below.

236   The two haplotypes on chromosome 6 overlapped with the Osteopetrosis Associated Transmembrane

237   Protein 1 (*OSTM1*) coding gene. *OSTM1* is associated with neurodegeneration[35, 36] and melanocyte

238   function,[37] and alpha-melanocyte stimulating hormone has been shown to have an effect on

239   depression-like symptoms.[38-40] This haplotype lies within the 6q21 region that has been associated

240   with bipolar disorder,[41-45] a disease that shares symptoms with MDD and has a correlated phenotypic

241   liability of 0.64.[46] This may indicate either a pleiotropic effect or clinical heterogeneity, whereby

242   patients may be misdiagnosed, i.e. patients may have MDD and transition to bipolar disorder in the

243   future or are sub-threshold for bipolar disorder and instead given a diagnosis of MDD.

244   The haplotype identified on chromosome 8 overlapped with the Interleukin 7 (*IL7*) protein coding

245   region. *IL7* is involved in maintaining T cell homeostasis[47] and proliferation,[48] which in turn

246   contributes to the immune response to pathogens. It has been proposed that impaired T cell function

247   may be a factor in the development of MDD,[49] with depressed subjects found to have elevated[50] or

248   depressed levels[51] of *IL7* serum. There is conjecture as to whether MDD causes inflammation or

249   represents a reaction to an increased inflammatory response,[52, 53] but it is most likely to be a

250   bidirectional relationship.[51]

251   The haplotype on chromosome 10 overlapped with two RNA genes: long intergenic non-protein

252   coding RNA 704 (*LINC00704*) and long intergenic non-protein coding RNA 705 (*LINC00705*). The

253   function of these non-protein coding genes is unreported. However, a study of cardiac neonatal lupus

254   which is a rare autoimmune disease demonstrated an association for a SNP (rs1391511) which is 15kb

255   from *LINC00705*.

256   Two Dutch studies[54, 55] have identified a variant (rs8023445) on chromosome 15 located within the

257   SRC (Src homology 2 domain containing) family, member 4 (*SHC4*) gene coding region that has a

258   moderate degree of association with MDD ($P = 1.64 \times 10^{-5}$ and $P = 9 \times 10^{-6}$, respectively). A variant

259   (rs10519201) within the *SHC4* coding region was also found to have an association ($P = 6.16 \times 10^{-6}$)

260   with Obsessive-Compulsive Personality Disorder in a UK-based study.[56] *SHC4* is expressed in

261   neurons[57] and regulates BDNF-induced MAPK activation,[58] which has been shown to be a key factor

262    in MDD pathophysiology.[59] The *SHC4* region overlaps with the haplotype on chromosome 15

263    identified in the discovery cohort (located at 49 206 902 – 49 260 601 bp) and therefore further

264    research to examine the association between the *SHC4* region and psychiatric disorders could be

265    warranted.

266    Haplotype-based analyses are capable of tagging variants due to the LD between the untyped variants

267    and the multiple flanking genotyped variants which make up the inherited haplotype. This approach

268    should provide greater power when there is comparatively higher IBD sharing, such as in GS:SFHS,

269    where there is a greater likelihood that a single haplotype is tagging the same causal variant across

270    that population. The UK Biobank was selected as replication cohort as it is a large population-based

271    sample that was expected to be genetically similar to the GS:SFHS discovery cohort. This was

272    confirmed by the similarity of the observed haplotype frequencies (Table 2) between the two cohorts.

273    The prevalence of MDD observed in the discovery cohort (13.7%) was comparable to that reported

274    (14.6%) within similar populations.[34] However, in the replication cohort, the trait prevalence was

275    notably higher (34.0%), most likely due to the differing methods of phenotypic ascertainment.

276    Additional work could seek to replicate the findings in further cohorts, as well as full meta-analysis of

277    all haplotypes within those cohorts. An additive model was used to analyse the haplotypes and

278    alternative approaches could implement a dominant model or an analysis of diplotypes (haplotype

279    pairs) for association with MDD.

280    **Conclusions**

281    This study identified two haplotypes within the discovery cohort that exceeded genome-wide

282    significance for association with a clinically diagnosed MDD phenotype. One of these haplotypes was

283    nominally significant in the replication cohort and was in LD with a haplotype that was genome-wide

284    significant in the meta-analysis. The genome-wide significant haplotype on chromosome 6 was

285    located on 6q21, which has been shown previously to be related to psychiatric disorders. There were a

286    number of haplotypes approaching genome-wide significance located within genic regions associated

287    with diseases that are comorbid with MDD and therefore these regions warrant further investigation.

288 The total genetic variance explained by the haplotypes identified was small, however these haplotypes

289 potentially represent biologically informative aetiological subtypes for MDD and merit further

290 analysis.

291 **ACKNOWLEDGEMENTS**

311 **CONFLICT OF INTEREST**

12

312 DJP and IJP are participants in UK Biobank. The authors report that no other financial interests or

313 potential conflicts of interest exist.

314 **REFERENCES**

315 1.    Huffman JC, Celano CM, Beach SR, Motiwala SR, Januzzi JL. Depression and cardiac
316       disease: epidemiology, mechanisms, and diagnosis. *Cardiovascular Psychiatry and*
317       *Neurology* 2013; **2013:** 14.

318
319 2.    Kang H-J, Kim S-Y, Bae K-Y, Kim S-W, Shin I-S, Yoon J-S*, et al.* Comorbidity of
320       depression with physical disorders: research and clinical implications. *Chonnam Medical*
321       *Journal* 2015; **51**(1)**:** 8-18.

322
323 3.    Raison CL, Capuron L, Miller AH. Cytokines sing the blues: inflammation and the
324       pathogenesis of depression. *Trends in Immunology* 2006; **27**(1)**:** 24-31.

325
326 4.    Major Depressive Disorder Working Group of the Psychiatric Gwas Consortium. A mega-
327       analysis of genome-wide association studies for major depressive disorder. *Molecular*
328       *Psychiatry* 2013; **18**(4)**:** 497-511.

329
330 5.    Converge Consortium. Sparse whole-genome sequencing identifies two loci for major
331       depressive disorder. *Nature* 2015; **523**(7562)**:** 588-591.

332
333 6.    Levinson DF, Mostafavi S, Milaneschi Y, Rivera M, Ripke S, Wray NR*, et al.* Genetic
334       studies of major depressive disorder: why are there no genome-wide association study
335       findings and what can we do about it? *Biological Psychiatry* 2014; **76**(7)**:** 510-512.

336
337 7.    Hyde CL, Nagle MW, Tian C, Chen X, Paciga SA, Wendland JR*, et al.* Identification of 15
338       genetic loci associated with risk of major depression in individuals of European descent.
339       *Nature Genetics* 2016; **advance online publication**.

340
341 8.    Lubke GH, Hottenga JJ, Walters R, Laurin C, de Geus EJC, Willemsen G*, et al.* Estimating
342       the genetic variance of major depressive disorder due to all single nucleotide polymorphisms.
343       *Biological Psychiatry* 2012; **72**(8)**:** 707-709.

344
345 9.    Sullivan PF, Neale MC, Kendler KS. Genetic epidemiology of major depression: review and
346       meta-analysis. *American Journal of Psychiatry* 2000; **157**(10)**:** 1552-1562.

347
348 10.   Zhang Z, Ni J, Zhang J, Tang W, Li X, Wu Z*, et al.* A haplotype in the 5'-upstream region of
349       the NDUFV2 gene is associated with major depressive disorder in Han Chinese. *Journal of*
350       *Affective Disorders* 2016; **190:** 329-332.

351
352 11.   Kim J-J, Mandelli L, Pae C-U, De Ronchi D, Jun T-Y, Lee C*, et al.* Is there protective
353       haplotype of dysbindin gene (DTNBP1) 3 polymorphisms for major depressive disorder.
354       *Progress in Neuro-Psychopharmacology and Biological Psychiatry* 2008; **32**(2)**:** 375-379.

355

356    12.    Klok MD, Giltay EJ, Van der Does AJW, Geleijnse JM, Antypa N, Penninx BWJH, *et al*. A
357           common and functional mineralocorticoid receptor haplotype enhances optimism and protects
358           against depression in females. *Translational Psychiatry* 2011; **1:** e62.

359
360    13.    Smith BH, Campbell A, Linksted P, Fitzpatrick B, Jackson C, Kerr SM, *et al*. Cohort profile:
361           Generation Scotland: Scottish Family Health Study (GS:SFHS). The study, its participants
362           and their potential for genetic research on health and illness. *International Journal of*
363           *Epidemiology* 2013; **42**(3)**:** 689-700.

364
365    14.    Smith DJ, Nicholl BI, Cullen B, Martin D, Ul-Haq Z, Evans J, *et al*. Prevalence and
366           characteristics of probable major depression and bipolar disorder within UK Biobank: cross-
367           sectional study of 172,751 participants. *PLoS ONE* 2013; **8**(11)**:** e75362.

368
369    15.    Amador C, Huffman J, Trochet H, Campbell A, Porteous D, Wilson JF, *et al*. Recent genomic
370           heritage in Scotland. *BMC Genomics* 2015; **16**(1)**:** 1-17.

371
372    16.    Allen NE, Sudlow C, Peakman T, Collins R. UK biobank data: come and get it. *Science*
373           *Translational Medicine* 2014; **6**(224)**:** 224ed224.

374
375    17.    Marchini J (2015). UK Biobank phasing and imputation documentation. Version 1.2:
376           http://biobank.ctsu.ox.ac.uk/crystal/docs/impute_ukb_v1.pdf.

377
378    18.    Delaneau O, Zagury J-F, Marchini J. Improved whole-chromosome phasing for disease and
379           population genetic studies. *Nature Methods* 2013; **10**(1)**:** 5-6.

380
381    19.    O'Connell J, Gurdasani D, Delaneau O, Pirastu N, Ulivi S, Cocca M, *et al*. A general
382           approach for haplotype phasing across the full spectrum of relatedness. *PLoS Genetics* 2014;
383           **10**(4)**:** e1004234.

384
385    20.    The International HapMap Consortium. A second generation human haplotype map of over
386           3.1 million SNPs. *Nature* 2007; **449**(7164)**:** 851-861.

387
388    21.    Browning BL, Browning SR. Improving the accuracy and efficiency of identity-by-descent
389           detection in population data. *Genetics* 2013; **194**(2)**:** 459-471.

390
391    22.    Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira Manuel A, Bender D, *et al*. PLINK: a
392           tool set for whole-genome association and population-based linkage analyses. *American*
393           *Journal of Human Genetics* 2007; **81**(3)**:** 559-575.

394
395    23.    Sham PC, Purcell SM. Statistical power and significance testing in large-scale genetic studies.
396           *Nat Rev Genet* 2014; **15**(5)**:** 335-346.

397
398    24.    First MB, Spitzer RL, Gibbon Miriam., Williams JBW (2002). Structured Clinical Interview
399           for DSM-IV-TR Axis I Disorders, Research Version, Patient Edition. (SCID-I/P)

25. Fernandez-Pujals AM, Adams MJ, Thomson P, McKechanie AG, Blackwood DHR, Smith BH, *et al*. Epidemiology and heritability of major depressive disorder, stratified by age of onset, sex, and illness course in generation scotland: scottish family health study (GS:SFHS). *PLoS ONE* 2015; **10**(11)**:** e0142197.

26. Information Services Division (2016). SMR Data Manual: http://www.ndc.scot.nhs.uk/Data-Dictionary/SMR-Datasets.

27. Yang J, Zaitlen NA, Goddard ME, Visscher PM, Price AL. Advantages and pitfalls in the application of mixed-model association methods. *Nature Genetics* 2014; **46**(2)**:** 100-106.

28. Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, *et al*. Common SNPs explain a large proportion of the heritability for human height. *Nature Genetics* 2010; **42**(7)**:** 565-569.

29. Zaitlen N, Kraft P, Patterson N, Pasaniuc B, Bhatia G, Pollack S, *et al*. Using extended genealogy to estimate components of heritability for 23 quantitative and dichotomous traits. *PLoS Genetics* 2013; **9**(5)**:** e1003520.

30. Bulik-Sullivan BK, Loh P-R, Finucane HK, Ripke S, Yang J, Schizophrenia Working Group of the Psychiatric Genomics Consortium, *et al*. LD score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nature Genetics* 2015; **47**(3)**:** 291-295.

31. Willer CJ, Li Y, Abecasis GR. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics* 2010; **26**(17)**:** 2190-2191.

32. Park J-H, Wacholder S, Gail MH, Peters U, Jacobs KB, Chanock SJ, *et al*. Estimation of effect size distribution from genome-wide association studies and implications for future discoveries. *Nature Genetics* 2010; **42**(7)**:** 570-575.

33. Purcell S, Cherny SS, Sham PC. Genetic power calculator: design of linkage and association genetic mapping studies of complex traits. *Bioinformatics* 2003; **19**(1)**:** 149-150.

34. Bromet E, Andrade LH, Hwang I, Sampson NA, Alonso J, de Girolamo G, *et al*. Cross-national epidemiology of DSM-IV major depressive episode. *BMC Medicine* 2011; **9**(1)**:** 1-16.

35. Kasper D, Planells-Cases R, Fuhrmann JC, Scheel O, Zeitz O, Ruether K, *et al*. Loss of the chloride channel ClC-7 leads to lysosomal storage disease and neurodegeneration. *The EMBO Journal* 2005; **24**(5)**:** 1079-1091.

36. Pandruvada SNM, Beauregard J, Benjannet S, Pata M, Lazure C, Seidah NG, *et al*. Role of ostm1 cytosolic complex with kinesin 5B in intracellular dispersion and trafficking. *Molecular and Cellular Biology* 2016; **36**(3)**:** 507-521.

446
447    37.    Hoek KS, Schlegel NC, Eichhoff OM, Widmer DS, Praetorius C, Einarsson SO, *et al*. Novel
448           MITF targets identified using a two-step DNA microarray strategy. *Pigment Cell &*
449           *Melanoma Research* 2008; **21**(6)**:** 665-676.

450
451    38.    Maes M, DeJonckheere C, Vandervorst C, Schotte C, Cosyns P, Raus J, *et al*. Abnormal
452           pituitary function during melancholia: Reduced α-melanocyte-stimulating hormone secretion
453           and increased intact ACTH non-suppression. *Journal of Affective Disorders* 1991; **22**(3)**:** 149-
454           157.

455
456    39.    Goyal SN, Kokare DM, Chopde CT, Subhedar NK. Alpha-melanocyte stimulating hormone
457           antagonizes antidepressant-like effect of neuropeptide Y in Porsolt's test in rats.
458           *Pharmacology Biochemistry and Behavior* 2006; **85**(2)**:** 369-377.

459
460    40.    Kokare DM, Singru PS, Dandekar MP, Chopde CT, Subhedar NK. Involvement of alpha-
461           melanocyte stimulating hormone (α-MSH) in differential ethanol exposure and withdrawal
462           related depression in rat: Neuroanatomical–behavioral correlates. *Brain Research* 2008; **1216:**
463           53-67.

464
465    41.    Knight J, Rochberg NS, Saccone SF, Nurnberger JI, Rice JP. An investigation of candidate
466           regions for association with bipolar disorder. *American Journal of Medical Genetics Part B:*
467           *Neuropsychiatric Genetics* 2010; **153B**(7)**:** 1292-1297.

468
469    42.    Dick DM, Foroud T, Flury L, Bowman ES, Miller MJ, Rau NL, *et al*. Genomewide linkage
470           analyses of bipolar disorder: a new sample of 250 pedigrees from the national institute of
471           mental health genetics initiative. *American Journal of Human Genetics* 2003; **73**(1)**:** 107-114.

472
473    43.    Park N, Juo SH, Cheng R, Liu J, Loth JE, Lilliston B, *et al*. Linkage analysis of psychosis in
474           bipolar pedigrees suggests novel putative loci for bipolar disorder and shared susceptibility
475           with schizophrenia. *Molecular Psychiatry* 2004; **9**(12)**:** 1091-1099.

476
477    44.    Pato CN, Pato MT, Kirby A, Petryshen TL, Medeiros H, Carvalho C, *et al*. Genome-wide
478           scan in Portuguese Island families implicates multiple loci in bipolar disorder: Fine mapping
479           adds support on chromosomes 6 and 11. *American Journal of Medical Genetics Part B:*
480           *Neuropsychiatric Genetics* 2004; **127B**(1)**:** 30-34.

481
482    45.    Fabbri C, Serretti A. Genetics of long-term treatment outcome in bipolar disorder. *Progress in*
483           *Neuro-Psychopharmacology and Biological Psychiatry* 2016; **65:** 17-24.

484
485    46.    McGuffin P, Rijsdijk F, Andrew M, Sham P, Katz R, Cardno A. The heritability of bipolar
486           affective disorder and the genetic relationship to unipolar depression. *Archives of General*
487           *Psychiatry* 2003; **60**(5)**:** 497-502.

488
489    47.    Surh CD, Sprent J. Homeostasis of Naive and Memory T Cells. *Immunity* 2008; **29**(6)**:** 848-
490           862.

491

492    48.    Kittipatarin C, Khaled AR. Interlinking interleukin-7. *Cytokine* 2007; **39**(1)**:** 75-83.

493
494    49.    Miller AH. Depression and immunity: A role for T cells? *Brain, Behavior, and Immunity*
495           2010; **24**(1)**:** 1-8.

496
497    50.    Simon NM, McNamara K, Chow CW, Maser RS, Papakostas GI, Pollack MH, *et al*. A
498           detailed examination of cytokine abnormalities in major depressive disorder. *European*
499           *Neuropsychopharmacology* 2008; **18**(3)**:** 230-233.

500
501    51.    Lehto SM, Huotari A, Niskanen L, Herzig K-H, Tolmunen T, Viinamäki H, *et al*. Serum IL-7
502           and G-CSF in major depressive disorder. *Progress in Neuro-Psychopharmacology and*
503           *Biological Psychiatry* 2010; **34**(6)**:** 846-851.

504
505    52.    Stewart JC, Rand KL, Muldoon MF, Kamarck TW. A prospective evaluation of the
506           directionality of the depression-inflammation relationship. *Brain, Behavior, and Immunity*
507           2009; **23**(7)**:** 936-944.

508
509    53.    Irwin MR, Miller AH. Depressive disorders and immunity: 20 years of progress and
510           discovery. *Brain, Behavior, and Immunity* 2007; **21**(4)**:** 374-383.

511
512    54.    Aragam N, Wang K-S, Pan Y. Genome-wide association analysis of gender differences in
513           major depressive disorder in the Netherlands NESDA and NTR population-based samples.
514           *Journal of Affective Disorders* 2011; **133**(3)**:** 516-521.

515
516    55.    Sullivan PF, de Geus EJC, Willemsen G, James MR, Smit JH, Zandbelt T, *et al*. Genome-
517           wide association for major depressive disorder: a possible role for the presynaptic protein
518           piccolo. *Molecular Psychiatry* 2008; **14**(4)**:** 359-375.

519
520    56.    Boraska V, Davis OSP, Cherkas LF, Helder SG, Harris J, Krug I, *et al*. Genome-wide
521           association analysis of eating disorder-related symptoms, behaviors, and personality traits.
522           *American Journal of Medical Genetics* 2012; **159B**(7)**:** 803-811.

523
524    57.    Hawley SP, Wills MKB, Rabalski AJ, Bendall AJ, Jones N. Expression patterns of ShcD and
525           Shc family adaptor proteins during mouse embryonic development. *Developmental Dynamics*
526           2011; **240**(1)**:** 221-231.

527
528    58.    You Y, Li W, Gong Y, Yin B, Qiang B, Yuan J, *et al*. ShcD interacts with TrkB via its PTB
529           and SH2 domains and regulates BDNF-induced MAPK activation. *BMB Rep* 2010; **43**(7)**:**
530           485-490.

531
532    59.    Duric V, Banasr M, Licznerski P, Schmidt HD, Stockmeier CA, Simen AA, *et al*. A negative
533           regulator of MAP kinase causes depressive behavior. *Nature Medicine* 2010; **16**(11)**:** 1328-
534           1332.

535

536    **Figure 1.** Manhattan plot representing the –$\log_{10}$ *P*-values for an association between each assessed
537    haplotype in the Generation Scotland: Scottish Family Health Study cohort and Major Depressive
538    Disorder
539
540    **Figure 2.** Regional association plot representing the –$\log_{10}$ *P*-values for an association between
541    haplotypes in the Generation Scotland: Scottish Family Health Study cohort and Major Depressive
542    Disorder within the 107.4 – 107.6 Mb region on chromosome 6. The start and end position (using
543    build GRCh37) of haplotypes represent the outermost SNP positions within the windows examined.
544    The warmth of colour represents the $r^2$ with the genome-wide significant haplotype located between
545    108 338 267 and 108 454 437 bp.
546
547    **Figure 3.** Region association plot representing the –$\log_{10}$ *P*-values for an association between
548    haplotypes in the Generation Scotland: Scottish Family Health Study cohort and Major Depressive
549    Disorder within the 3.6 – 5.8 Mb region on chromosome 10. The start and end position (using build
550    GRCh37) of haplotypes represent the outermost SNP positions within the windows examined. The
551    warmth of colour represents the $r^2$ with the genome-wide significant haplotype located between 4 588
552    261 and 4 822 210 bp.
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588

18

**Table 1.** The genetic association between Major Depressive Disorder and 12 haplotypes in the Generation Scotland: Scottish Family Health Study (GS:SFHS) discovery cohort (where $P < 10^{-6}$), the replication cohort (UK Biobank) and a meta-analysis.

| Haplotype | | | GS:SFHS | | UK Biobank | | Meta-analysis | |
|---|---|---|---|---|---|---|---|---|
| Chr. | Position (bp) | Window Size (cM) | Odds Ratio (95% CI) | *P*-value | Odds Ratio (95% CI) | *P*-value | Odds Ratio (95% CI) | *P*-value |
| { 6 * | 108338267 - 108454437 | 0.34 | **1.83 (1.53 - 2.16)** | **7.06 x 10⁻⁹** | **1.11 (1.01 - 1.22)** | **3.62 x 10⁻²** | 1.26 (1.16 - 1.37) | 3.14 x 10⁻⁷ |
| 6 | 108407662 - 108454437 | 0.25 | 1.68 (1.42 - 1.96) | 8.17 x 10⁻⁸ | **1.14 (1.04 - 1.24)** | **4.47 x 10⁻³** | **1.25 (1.16 - 1.35)** | **4.38 x 10⁻⁸** |
| 7 | 139682412 - 139708901 | 0.25 | 2.17 (1.67 - 2.73) | 4.37 x 10⁻⁷ | 0.87 (0.68 - 1.08) | 2.20 x 10⁻¹ | 1.28 (1.08 - 1.49) | 4.67 x 10⁻³ |
| { 8 | 79700362 - 80387861 | 0.5 | 1.98 (1.56 - 2.46) | 9.02 x 10⁻⁷ | 1.06 (0.86 - 1.28) | 5.93 x 10⁻¹ | 1.36 (1.18 - 1.56) | 6.29 x 10⁻⁵ |
| 8 | 79759499 - 80156474 | 0.25 | 1.77 (1.47 - 2.10) | 7.90 x 10⁻⁸ | 1.05 (0.91 - 1.21) | 5.06 x 10⁻¹ | 1.28 (1.15 - 1.42) | 1.14 x 10⁻⁵ |
| 10 | 4588261 - 4822210 | 0.5 | **2.33 (1.83 - 2.91)** | **8.50 x 10⁻⁹** | 1.15 (0.80 - 1.59) | 4.39 x 10⁻¹ | 1.67 (1.40 - 1.98) | 7.92 x 10⁻⁸ |
| 11 * | 2260854 - 2437425 | 0.41 | 1.64 (1.38 - 1.91) | 2.86 x 10⁻⁷ | 1.00 (0.87 - 1.34) | 9.91 x 10⁻¹ | 1.26 (1.10 - 1.34) | 1.32 x 10⁻⁴ |
| 12 | 48159721 - 48263828 | 0.25 | 2.00 (1.58 - 2.47) | 4.78 x 10⁻⁷ | 0.97 (0.79 - 1.17) | 7.36 x 10⁻¹ | 1.29 (1.12 - 1.48) | 6.51 x 10⁻⁴ |
| 12 | 116904503 - 117062860 | 0.25 | 2.13 (1.64 - 2.69) | 9.90 x 10⁻⁷ | 1.04 (0.79 - 1.34) | 7.79 x 10⁻¹ | 1.45 (1.22 - 1.71) | 5.37 x 10⁻⁵ |
| 15 | 49206902 - 49260601 | 0.25 | 2.03 (1.62 - 2.48) | 9.21 x 10⁻⁸ | 1.09 (0.88 - 1.32) | 4.04 x 10⁻¹ | 1.41 (1.22 - 1.61) | 4.39 x 10⁻⁶ |
| { 15 | 93806447 - 93851224 | 0.5 | 1.58 (1.34 - 1.83) | 4.47 x 10⁻⁷ | 0.93 (0.81 - 1.05) | 2.38 x 10⁻¹ | 1.16 (1.05 - 1.27) | 2.50 x 10⁻³ |
| 15 | 93821340 - 93845622 | 0.25 | 1.52 (1.31 - 1.75) | 8.67 x 10⁻⁷ | 0.91 (0.81 - 1.03) | 1.37 x 10⁻¹ | 1.13 (1.03 - 1.23) | 6.97 x 10⁻³ |

Bold values indicate genome-wide statistical significance ($P < 5 \times 10^{-8}$) was achieved in the GS:SFHS cohort or the meta-analysis, or that nominal statistical significance ($P < 0.05$) was achieved in the UK Biobank. Base pair (bp) positions are based on build GRCh37. * indicates haplotype boundaries defined by the fine mapping approach. { indicates linkage disequilibrium ($r^2$) > 0.5 between haplotypes in the GS:SFHS cohort.

19

615  **Table 2.** Protein coding genes located overlapping with the 12 haplotypes with $P < 10^{-6}$ in the
616  Generation Scotland: Scottish Family Health Study (GS:SFHS) discovery cohort and the frequencies
617  of those haplotypes in GS:SFHS and UK Biobank.

| Chr. | Position (bp) | Protein coding genes | Haplotype frequency | |
|---|---|---|---|---|
| | | | GS:SFHS | UK Biobank |
| 6 | 108338267 - 108454437 | OSTM1 | 0.0152 | 0.0197 |
| 6 | 108407662 - 108454437 | OSTM1 | 0.0193 | 0.0241 |
| 7 | 139682412 - 139708901 | TBXAS1 | 0.0066 | 0.0069 |
| 8 | 79700362 - 80387861 | IL7 | 0.0076 | 0.0081 |
| 8 | 79759499 - 80156474 | IL7 | 0.0147 | 0.0157 |
| 10 | 4588261 - 4822210 | | 0.0064 | 0.0027 |
| 11 | 2260854 - 2437425 | ASCL2, CLorf21, TSPAN32, CD81, TSSC4, TRPM5 | 0.0196 | 0.0187 |
| 12 | 48159721 - 48263828 | SLC48A1, RAPGEF3, HDAC7, VDR | 0.0078 | 0.0090 |
| 12 | 116904503 - 117062860 | MAP1LC3B2 | 0.0057 | 0.0045 |
| 15 | 49206902 - 49260601 | SHC4 | 0.0082 | 0.0080 |
| 15 | 93806447 - 93851224 | | 0.0224 | 0.0206 |
| 15 | 93821340 - 93845622 | | 0.0265 | 0.0243 |

618  Base pair (bp) positions are based on build GRCh37 with protein coding regions obtained from
619  Ensembl, GRCh37.p13. Haplotype frequencies were calculated using unrelated individuals and
620  excluding UK Biobank participants recruited in Glasgow or Edinburgh. { indicates a linkage
621  disequilibrium ($r^2$) > 0.5 between haplotypes in the GS:SFHS cohort.

20