

1 **The genetic and phenotypic correlates of neonatal Complement Component 3 and 4**
2 **protein concentrations with a focus on psychiatric and autoimmune disorders**

3

4 Nis Borbye-Lorenzen, PhD*equal first author
5 Center for Neonatal Screening, Department of Congenital Disorders, Statens Serum Institute,
6 Copenhagen, Denmark

7

8 Zhihong Zhu, PhD*equal first author
9 National Centre for Register-Based Research, Aarhus University, 8210 Aarhus V, Denmark

10

11

12 Esben Agerbo, DrMedSc
13 The Lundbeck Foundation Initiative for Integrative Psychiatric Research, iPSYCH, 8210 Aarhus V,
14 Denmark
15 Centre for Integrated Register-based Research, Aarhus University, CIRRAU, 8210 Aarhus V, Denmark
16 National Centre for Register-Based Research, Aarhus University, 8210 Aarhus V, Denmark

17

18 Clara Albiñana, MSc
19 National Centre for Register-Based Research, Aarhus University, 8210 Aarhus V, Denmark
20 The Lundbeck Foundation Initiative for Integrative Psychiatric Research, iPSYCH, 8210 Aarhus V,
21 Denmark

22

23 Michael E. Benros, MD, PhD
24 Copenhagen Research Centre for Mental Health, Mental Health Centre Copenhagen, Copenhagen
25 University Hospital, Hellerup, Denmark
26 Department of Immunology and Microbiology, Faculty of Health and Medical Sciences,
27 University of Copenhagen, Copenhagen, Denmark

28

29 Beilei Bian, MSc
30 Institute for Molecular Bioscience, The University of Queensland, Brisbane, Queensland, Australia

31

32 Anders D Børghlum, MD, PhD
33 Department of Biomedicine and the iSEQ Center, Aarhus University, Aarhus, Denmark
34 The Lundbeck Foundation Initiative for Integrative Psychiatric Research, iPSYCH, 8210 Aarhus V,
35 Denmark
36 Center for Genomics and Personalized Medicine, CGPM, Aarhus, Denmark

37

38 Cynthia M. Bulik, PhD
39 Department of Psychiatry, University of North Carolina at Chapel Hill, Chapel Hill, USA
40 Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden
41 Department of Nutrition, University of North Carolina at Chapel Hill, Chapel Hill, USA

42

43 Jean-Christophe Philippe Goldtsche Debost, MD, PhD
44 Aarhus University Hospital Skejby, Department of Psychosis, Aarhus Nord, Denmark
45 National Centre for Register-Based Research, Aarhus University, 8210 Aarhus V, Denmark

46

47 Jakob Grove, PhD
48 Department of Biomedicine (Human Genetics), Aarhus University, Aarhus, Denmark
49 The Lundbeck Foundation Initiative for Integrative Psychiatric Research, iPSYCH, 8210 Aarhus V,
50 Denmark

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

51 Center for Genomics and Personalized Medicine, Aarhus, Denmark
52 Bioinformatics Research Centre, Aarhus University, 8000 Aarhus C, Denmark
53
54 David M. Hougaard, DMSci
55 Department for Congenital Disorders, Statens Serum Institut, 2300 Copenhagen S, Denmark.
56 The Lundbeck Foundation Initiative for Integrative Psychiatric Research, iPSYCH, 8210 Aarhus V,
57 Denmark
58
59 Allan F McRae, PhD
60 Institute for Molecular Bioscience, The University of Queensland, Brisbane, Queensland 4072,
61 Australia
62
63 Ole Mors, PhD
64 Psychosis Research Unit, Aarhus University Hospital – Psychiatry, Denmark
65 The Lundbeck Foundation Initiative for Integrative Psychiatric Research, iPSYCH, 8210 Aarhus V,
66 Denmark
67
68 Preben Bo Mortensen, MD, PhD
69 National Centre for Register-Based Research, Aarhus University, 8210 Aarhus V, Denmark
70 The Lundbeck Foundation Initiative for Integrative Psychiatric Research, iPSYCH, 8210 Aarhus V,
71 Denmark
72 Centre for Integrated Register-based Research, Aarhus University, CIRRAU, 8210 Aarhus V, Denmark
73
74 Katherine L. Musliner, PhD
75 Department of Affective Disorders, Aarhus University and Aarhus University Hospital-Psychiatry,
76 Aarhus Denmark
77
78 Merete Nordentoft, MD, DrMedSc
79 The Lundbeck Foundation Initiative for Integrative Psychiatric Research, iPSYCH, 8210 Aarhus V,
80 Denmark
81 Mental Health Services in the Capital Region of Denmark, Mental Health Center Copenhagen,
82 University of Copenhagen, 2100 Copenhagen, Denmark
83 Department of Clinical Medicine, University of Copenhagen, 2200 Copenhagen N, Denmark
84
85 Liselotte V. Petersen, PhD.
86 National Centre for Register-Based Research, Aarhus University, 8210 Aarhus V, Denmark
87
88 Florian Privé, PhD
89 National Centre for Register-Based Research, Aarhus University, 8210 Aarhus V, Denmark
90
91 Julia Sidorenko, PhD
92 Institute for Molecular Bioscience, The University of Queensland, Brisbane, Queensland 4072,
93 Australia
94
95 Kristin Skogstrand, PhD
96 Center for Neonatal Screening, Department of Congenital Disorders, Statens Serum Institute,
97 Copenhagen, Denmark
98
99 Thomas Werge, PhD
100 Institute of Biological Psychiatry, Mental Health Services, Copenhagen University Hospital.
101 Department of Clinical Medicine, University of Copenhagen, 2200 Copenhagen N, Denmark

102 Lundbeck Center for Geogenetics, GLOBE Institute, University of Copenhagen, Copenhagen,
103 Denmark.
104 The Lundbeck Foundation Initiative for Integrative Psychiatric Research, iPSYCH, 8210 Aarhus V,
105 Denmark
106
107 Naomi R Wray, PhD
108 Institute for Molecular Bioscience, The University of Queensland, Brisbane, Queensland 4072,
109 Australia
110 Queensland Brain Institute, The University of Queensland, Brisbane, Queensland 4072, Australia
111
112
113 Bjarni J. Vilhjálmsson, PhD
114 National Centre for Register-Based Research, Aarhus University, 8210 Aarhus V, Denmark
115 Bioinformatics Research Centre, Aarhus University, 8000 Aarhus C, Denmark
116 The Lundbeck Foundation Initiative for Integrative Psychiatric Research, iPSYCH, 8210 Aarhus V,
117 Denmark
118
119 John J. McGrath, MD, PhD
120 National Centre for Register-Based Research, Aarhus University, 8210 Aarhus V, Denmark
121 Queensland Centre for Mental Health Research, The Park Centre for Mental Health, Brisbane,
122 Queensland 4076, Australia
123 Queensland Brain Institute, University of Queensland, Brisbane, Queensland 4072, Australia
124
125
126

127 **Abstract**

128 The complement system, including complement components 3 and 4 (C3, C4), traditionally has been
129 linked to innate immunity. More recently, complement components have also been implicated in
130 brain development and the risk of schizophrenia. Based on a large, population-based case-cohort
131 study, we measured the blood concentrations of C3 and C4 in 68,768 neonates. We found a strong
132 correlation between the concentrations of C3 and C4 (phenotypic correlation = 0.65, P -value <
133 1.0×10^{-100} , genetic correlation = 0.38, P -value = 1.9×10^{-35}). A genome-wide association study (GWAS)
134 for C4 protein concentration identified 36 independent loci, 30 of which were in or near the major
135 histocompatibility complex on chromosome 6 (which includes the *C4* gene), while six loci were found
136 on six other chromosomes. A GWAS for C3 identified 15 independent loci, seven of which were
137 located in the *C3* gene on chromosome 19, and eight loci on five other chromosomes. We found no
138 association between (a) measured neonatal C3 and C4 concentrations, imputed C4 haplotypes, or
139 predicted *C4* gene expression, with (b) schizophrenia (SCZ), bipolar disorder (BIP), depression (DEP),
140 autism spectrum disorder, attention deficit hyperactivity disorder or anorexia nervosa diagnosed in
141 later life. Mendelian randomisation (MR) suggested a small positive association between higher C4
142 protein concentration and an increased risk of SCZ, BIP, and DEP, but these findings did not persist in
143 more stringent analyses. Evidence from MR supported causal relationships between C4
144 concentration and several autoimmune disorders: systemic lupus erythematosus (SLE, OR and 95%
145 confidence interval, 0.37, 0.34 – 0.42); type-1 diabetes (T1D, 0.54, 0.50 - 0.58); multiple sclerosis
146 (MS, 0.68, 0.63 - 0.74); rheumatoid arthritis (0.85, 0.80 - 0.91); and Crohn's disease (1.26, 1.19 -
147 1.34). A phenome-wide association study (PheWAS) in UK Biobank confirmed that the genetic
148 correlates of C4 concentration were associated a range of autoimmune disorders including coeliac
149 disease, thyrotoxicosis, hypothyroidism, T1D, sarcoidosis, psoriasis, SLE and ankylosing spondylitis.
150 We found no evidence of associations between C3 versus mental or autoimmune disorders based on
151 either MR or PheWAS. In general, our results do not support the hypothesis that C4 is causally
152 associated with the risk of SCZ (nor several other mental disorders). We provide new evidence to
153 support the hypothesis that higher C4 concentration is associated with lower risks of autoimmune
154 disorders.

155 **Keywords:** Complement Component, C3, C4, schizophrenia, psychiatric disorders, autoimmune
156 disorders, genome-wide association study, Mendelian randomization, phenome-wide association
157 study, dried blood spots.

158

159 Introduction

160 The complement systems are an integral part of the innate immune response¹⁻⁴. These
161 phylogenetically-ancient systems involve complex and interlinked amplification cascades, which can
162 be triggered to protect the body from pathogens. Elements of the system are also involved in a
163 range of additional physiological functions. For example, a growing body of evidence links elements
164 of the complement systems (e.g. Complement Component 4; C4) to brain development and
165 psychiatric disorders⁵⁻⁷.

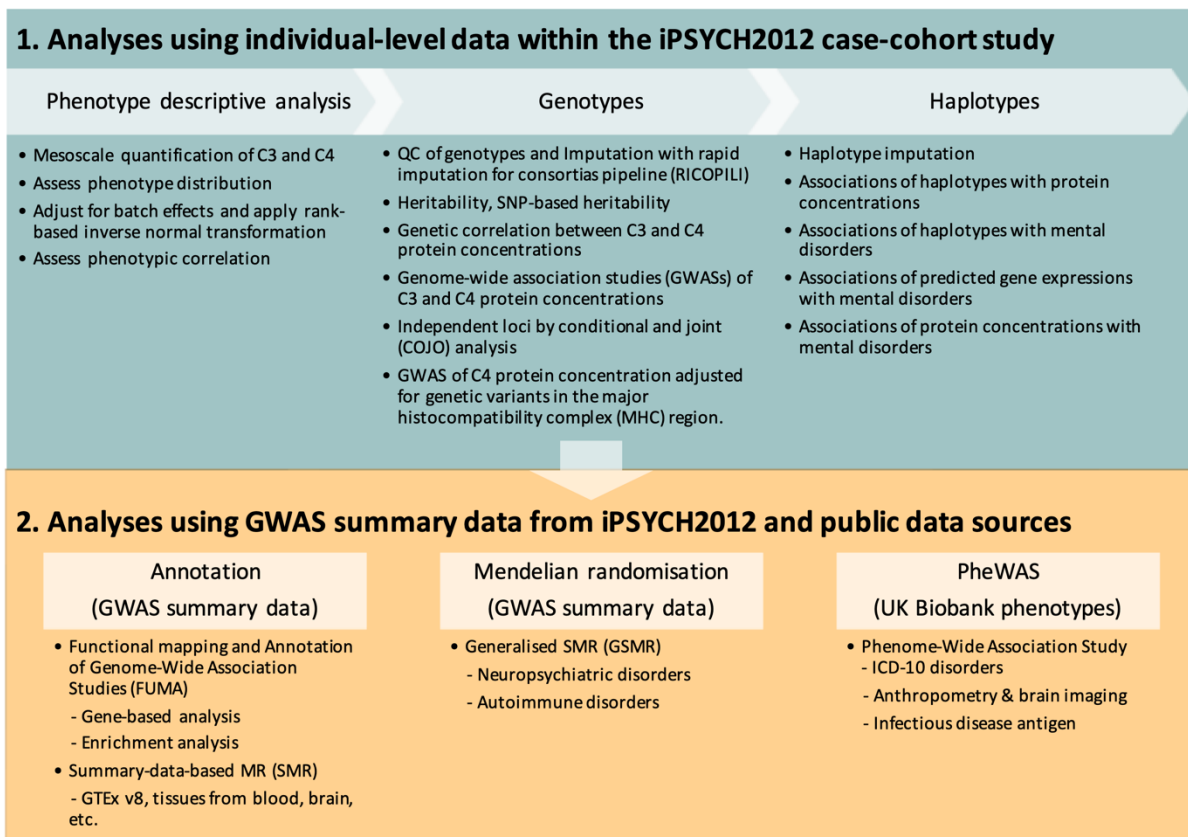
166 Several of the genes that encode components of the complement system, including *C4*, are located
167 within the major histocompatibility complex (MHC). Linking disease phenotypes with loci within the
168 MHC is difficult because of the long-range linkage disequilibrium (LD) in this region. Furthermore,
169 the *C4* gene has two homologous isoforms (*C4A* and *C4B*), each of which can vary according to an
170 insertion of a human endogenous retrovirus (*HERV*) transposon, and which can vary between one to
171 three genocopies per haplotype⁸. Sekar and colleagues⁹ proposed that genetic variants involving the
172 *C4* gene could account for the strong signal over this region detected in genome-wide association
173 studies (GWAS) of schizophrenia¹⁰. They found that the expression of mRNA transcripts coding for
174 the *C4A*-related isoform was increased in post-mortem schizophrenia brain samples (cases = 35,
175 controls = 70), and also reported that the *C4A* copy number was associated with both increased *C4A*
176 expression in the brain and increased risk of schizophrenia. These findings are of interest to
177 neurodevelopmental disorders, given evidence that *C4* and related members of the complement
178 systems are involved in synaptic pruning during early brain development^{9,11-13}. Apart from the links
179 with schizophrenia, increased *C4A* copy number has also been associated with a *decreased* risk of
180 autoimmune disorders^{14,15}. In light of the shared genetic architecture between different types of
181 mental disorders¹⁶, and the links between *C4A* alleles and risk of autoimmune disorders, there is a
182 need to explore if the putative risk haplotypes are associated with a wider range of both mental- and
183 autoimmune disorders.

184 As more copies of the *C4A* gene are associated with increased expression of *C4A*-related transcripts
185 in the brain^{9,17}, it is reasonable to assume that an increased copy number of the *C4A* gene would also
186 be reflected in increased expression of the C4 protein. C4 is an abundant circulating protein,
187 produced mainly in the liver, and evidence from transgenic mouse experiments¹² and human
188 observational studies^{18,19} confirms a dose-response relationship between increased *C4* copy number
189 and increased concentration of C4 protein in the peripheral circulation. The concentration of the C4
190 protein is correlated with other members of the complement family, including complement
191 component C3 (encoded by the *C3* gene)²⁰. A recent systematic review and meta-analysis found that
192 the serum concentrations of C4 and C3 did not differ in those with schizophrenia versus controls²¹.
193 However, the neurodevelopmental hypothesis of schizophrenia suggests that early life disruption of
194 brain development may underpin the subsequent adult-onset disorder^{22,23}. Additional evidence
195 suggests that complement-related synaptic pruning may be most prominent during early life (e.g. C3
196 protein expression peaks in the early post-natal period and decreases with age)^{7,24}. We are aware of
197 one study that examined neonatal C4 concentration in a schizophrenia case-control study (cases =
198 75, controls = 644)²⁵. This study found evidence that an increased concentration of one of the two
199 measured peptides within the protein encoded by *C4A* was associated with an increased risk of
200 subsequent schizophrenia. There is a need to explore the association between complement-related
201 protein concentrations and later mental disorders based on larger samples of neonatal blood
202 samples.

203 We examined the association between neonatal C3 and C4 concentrations and later health
204 outcomes based on a population-based case-cohort study that had access to archived neonatal dried

205 blood spots^{26,27}. Building on the only published GWAS of serum complement components C3 and C4
 206 (studies based on 3,495 Han Chinese men¹⁸), we completed GWASs for C3 and C4 neonatal
 207 concentrations based on 68,768 samples and estimated the heritability of these phenotypes. We
 208 were also able to assess the association between imputed C4 haplotypes and observed neonatal C4
 209 protein concentration. We examined the association of both (a) imputed C4 haplotypes and (b)
 210 observed C3 and C4 protein concentration in neonatal dried blood spots versus the risk of a range of
 211 clinically observed mental disorders in the case-cohort study (i.e. schizophrenia, bipolar disorder,
 212 depression, autism spectrum disorder, attention deficit hyperactivity disorder and anorexia
 213 nervosa). Based on the results of the GWASs of C3 and C4, we used: (a) bioinformatics tools to
 214 explore gene properties of the genome-wide significant loci (summary-data-based MR [SMR], gene-
 215 based analyses and gene set analyses); (b) Mendelian randomization analyses to explore the
 216 associations between C3 and C4 protein concentrations versus a range of mental disorders (as listed
 217 above) and autoimmune disorders (multiple sclerosis [MS], type-1 diabetes [T1D], Crohn’s disease
 218 [CD], ulcerative colitis [UC], rheumatoid arthritis [RA], and systemic lupus erythematosus [SLE]), and
 219 finally, (c) phenome-wide association studies (PheWAS)²⁸ to examine the relationship between the
 220 genetic correlates of C3 and C4 protein concentrations versus mental and autoimmune disorders, as
 221 well as a wide range of other health phenotypes. A summary of the overall methods is shown in
 222 Figure 1.

223



224

225 **Figure 1.** Methods figure.

226

227 **METHODS**

228

229 ***The iPSYCH2012 study***

230

231 Key elements of this study were based on the Lundbeck Foundation Initiative for Integrative
232 Psychiatric Research (iPSYCH) sample²⁶, a population-based case-cohort designed to study the
233 genetic and environmental factors of schizophrenia (SCZ), bipolar disorder (BIP), depression (DEP),
234 autism spectrum disorder (ASD) and attention deficit hyperactivity disorder (ADHD). The original
235 iPSYCH sample (known as iPSYCH2012) included information on case status complete through 31
236 December 2012. We also included 4,791 anorexia nervosa cases (AN; ANGI-DK) from the Anorexia
237 Nervosa Genetics Initiative (ANGI)²⁹, which had the same design as iPSYCH2012. Henceforth, we
238 refer to iPSYCH2012 as the combined dataset with the ANGI samples. The iPSYCH2012 sample is
239 nested within the entire Danish population born between 1981 and 2005 ($n=1,472,762$). Diagnoses
240 were identified in the Danish Central Psychiatric Research Register^{30,31}, which includes all inpatient
241 contacts in Danish psychiatric hospitals since 1969 and all outpatient and emergency contacts since
242 1995. The ICD-10 codes used to classify the psychiatric disorder cases can be found in

243 **Supplementary Table 1.** The phenotype information for the iPSYCH2012 participants was updated
244 for the target mental disorders until December 2016. The case-cohort sample includes a population-
245 based random sub-cohort³² ($N = 30,000$) with an inclusion probability of 2.04% of the study base
246 ($30,000 / 1,472,762$). This sub-cohort also includes some participants with the target mental
247 disorders of interest. The genotypes and C3 and C4 protein concentrations were measured in
248 neonatal dried bloodspots (DBSs) taken as part of routine screening at birth from all babies born in
249 Denmark since 1981 and stored in the Danish Neonatal Screening Biobank³³. Dried blood spot
250 samples have been collected from practically all neonates born in Denmark since 1st May 1981 and
251 stored at -20°C . Samples are collected 4–7 days after birth. After the dried blood spots were
252 retrieved from the biobank, samples were extracted in a PBS buffer and stored for further use at -
253 80°C . Subsequently, DNA was extracted according to previously published methods³⁴. After storage
254 the protein extracts were assayed for C3 and C4 concentrations. Thus, all genotypes and C3/C4
255 protein concentration data originated from a single DBS extraction. Additional details related to
256 blood spot extraction and storage are provided in **Supplementary Methods 1**.

257

258 ***Ethical framework***

259 Material from the Danish Neonatal Screening Biobank has been used primarily for screening for
260 congenital disorders, but are also stored for follow-up diagnostics, screening, quality control and
261 research. According to Danish legislation, material from The Danish Neonatal Screening Biobank can
262 be used for research after approval from the Biobank, and the relevant Scientific Ethical Committee.
263 There is also a mechanism in place ensuring that one can opt out of having the stored material used
264 for research. The Danish Data Protection Agency and the Danish Health Data Authority approved this
265 study. According to Danish law, informed consent is not required for register-based studies. All data
266 accessed were deidentified.

267 ***C3 and C4 protein concentrations***

268 These methods have been described in a related study³⁵. Two 3.2 mm discs of DBS were punched
269 into 96 well polymerase chain reaction plates (72.1981.202, Sarstedt). The extracts were analyzed
270 with a multiplex immunoassay (also measuring vitamin D binding protein³⁵) using U-plex plates
271 (Meso-Scale Diagnostics (MSD), Maryland, US) employing antibodies specific for complement C3
272 (HYB030-07 and HYB030-06) and complement C4 (MA1-72520 (ThermoFisher Scientific) and
273 HYB162-04). The antibodies were purchased from SSI Antibodies (Copenhagen, Denmark) except if
274 otherwise stated. Extracts were analyzed diluted 1:70 in diluent 101 (#R51AD, MSD). Capture
275 antibodies (used at $10\ \mu\text{g}/\text{mL}$ as input concentration) were biotinylated in-house using EZ-Link Sulfo-
276 NHS-LC-Biotin (#21327, Thermo Fisher Scientific) and detection antibodies were SULFO-tagged

277 (R91AO, MSD), both at a challenge ratio of 20:1. As calibrators, we used complement components
278 purified from human: C3: #PSP-109 (Nordic Biosite, Copenhagen, DK), C4: abx060108 (Abbexa,
279 Cambridge, UK). Calibrators were diluted in diluent 101, detection antibodies (used at 1 µg/mL) were
280 diluted in diluent 3 (#R50AP, MSD). Controls were made in-house from part of the calibrator solution
281 in one batch, aliquoted in portions for each plate, and stored at -20°C until use. The samples were
282 prepared on the plates as recommended by the manufacturer and were read on the QuickPlex SQ
283 120 (MSD) 4 min after adding 2x Read buffer T (#R92TC, MSD). Analyte concentrations were
284 calculated from the calibrator curves on each plate using 4PL logistic regression using the MSD
285 Workbench software.

286 Intra-assay variations were calculated from 38 measurements analyzed on the same plate of a pool
287 of extract made from 304 samples. Inter-assay variations were calculated from controls analyzed in
288 duplicate on each plate during the sample analysis, 1022 plates in total. Lower limits of detections
289 were calculated as 2.5 standard deviations from 40 replicate measurements of the zero calibrator.
290 The higher detection limit was defined as the highest calibrator concentration. The lower and upper
291 detection limits for: (a) C3 were 95.4 and 7.98×10^4 µg/L respectively, and (b) C4 were 55.2 and
292 7.98×10^4 µg/L respectively. The intra- and inter-assay coefficient of variation (CV) for (a) C3 were
293 5.2% and 18.1% respectively; and for (b) C4 were 3.9% and 8.5% respectively. To validate the
294 stability of the samples during storage, we randomly selected 15-16 samples from five years (1984,
295 1992, 2000, 2008, and 2016; a total of 76 samples). After extracting the samples and adding them to
296 an MSD plate, the rest of the extracts were frozen for 2 months, thawed and measured as described
297 above to imitate the freeze-thaw cycle of the samples in the study. The oldest samples (from 1984)
298 recorded higher concentrations, most probably due to a change in the type of filter paper after
299 1989. In light of this artefact, we adjusted all DBP values by plate (the sequence of testing followed
300 the date of birth of the sample). This is described below. Additional details related to pre-analytic
301 variation are provided in **Supplementary Methods 1**.

302 ***Imputation of genotypes***

303 DNA genotyping was conducted at the Broad Institute (Boston, MA, USA) using the Infinium
304 PsychChip v1.0 array (Illumina, San Diego, CA, USA)³⁶. We restricted the genotyped SNPs to 252,339
305 high quality, common SNPs based on build hg19 (the same human genome reference build was used
306 throughout this study). Details of the filtering can be found elsewhere (Schork et al., 2019). Briefly,
307 we excluded SNPs with minor allele frequency (MAF) < 0.01, Hardy Weinberg Equilibrium (HWE) *p*-
308 value < 1.0×10^{-6} or non-SNP alleles (i.e., insertions and deletions, INDELS). 245,328 autosomal and
309 7,011 X-chromosome (chrX) SNPs were retained and used to impute SNPs using the Ricopili
310 pipeline³⁷ with the Haplotype Reference Consortium (HRC)³⁸ as the imputation reference panel
311 (accession number: EGAD00001002729). 6,743,499 autosomal SNPs, 227,371 chrX SNPs for males
312 and 184,517 chrX SNPs for females were retained with missing rate < 0.02 and genotype call
313 probability > 0.8. We further excluded the imputed SNPs with imputation info score < 0.8, MAF <
314 0.01 or HWE *p*-value < 1.0×10^{-6} . 5,201,724 SNPs were retained in autosomes and 126,109 SNPs were
315 retained on chrX. We then used the common SNPs to infer the genetic ancestries of 80,873
316 participants in the iPSYCH2012 study, 75,764 individuals of European ancestry and 5,109 individuals
317 of non-European ancestry. Details are provided in **Supplementary Method 2**.

318 ***Imputation of C4 haplotypes***

319 C4 haplotypes were imputed from reference data^{9,14} using the genotyped SNPs in the iPSYCH2012
320 sample. The human C4 haplotypes have various copy numbers, including two isotypic
321 polymorphisms, C4A (A) and C4B (B). Each isotype has two length-polymorphisms due to a human
322 endogenous retroviral (HERV) insertion, long form (L, with HERV insertion) and short form (S,
323 without HERV insertion). The isotypic and length polymorphisms lead to four alleles in a C4 copy, AL,
324 AS, BL and BS. Using the genotyped SNPs, the C4 haplotype reference was used to impute the C4

325 alleles and the number of *C4* copies (with a maximum copy number of 4). The *C4* haplotype
326 imputation panel comprised whole genome sequencing data from 1,234 individuals of multiple
327 ancestries, which enabled us to identify *C4* alleles with high accuracy. We used Beagle software³⁹ for
328 the imputation with the *C4* haplotype reference. The imputation results provided the counts of
329 alleles, but were unable to confidently distinguish all combinations of variants, for example,
330 between the haplotypes *AS-BL* and *AL-BS*. We counted the two *C4* alleles (*C4A* and *C4B*) with
331 combination of *HERV* using a subset of the imputed result, where combinations can be confidently
332 distinguished (details are provided in **Supplementary Method 3**). Both counts of *C4* allele
333 combinations and reported studies⁴⁰ indicated that the *C4A* gene is more likely to carry *HERV*
334 insertion than the *C4B* gene. Therefore, the *C4* haplotype is assumed to be *AL-BS* rather than *AS-BL*,
335 consistent with methods described by Sekar et al⁹. The imputed counts were converted to the *C4*
336 haplotypes. Eight common *C4* haplotypes (allele frequencies ≥ 0.01) were imputed in the
337 iPSYCH2012 study (**Supplementary Table 3**). The allele frequencies of the 8 haplotypes were
338 consistent with other studies^{14,41}. Given the common *C4* haplotypes, we counted the copy numbers
339 of the *C4* alleles (**Supplementary Figure 2**) for each participant. 28 individuals (0.04%) carried 4
340 copies of *C4B* and 35 individuals (0.05%) carried 6 copies of *HERV* insertion. Therefore, we excluded
341 these individuals with very rare copy numbers. The *C4A* copy number is strongly correlated with *C4B*
342 and *HERV* copy numbers (Pearson correlation between *C4A* and *C4B* = -0.52; between *C4A* and *HERV*
343 = 0.73).

344 **Quality control of the C3 and C4 protein concentrations**

345 The C3 and C4 protein concentrations were measured in 78,268 iPSYCH2012 participants of multiple
346 ancestries. We focused on 68,768 individuals of European ancestry with measures of C3 and C4
347 protein concentrations. The protein assay plates captured a substantial amount of variance (C3 =
348 49.4%, C4 = 45.3%). Therefore, we used a linear mixed model (LMM)⁴² approach to adjust protein
349 concentrations, $y = \sum z_{\text{plate}} u_{\text{plate}} + e$, where y represents the C3/4 protein concentration; z_{plate}
350 represents protein assay plate, a random variable; u_{plate} represents the random effect of protein
351 assay plate; and e represents residual. The mixed model regression was conducted by the R package
352 of lme4 (Bates et al., 2015). The rank-based inverse normal transformation (RINT)⁴³ was applied to
353 the residuals to have mean 0 and variance 1. The standard deviations (SDs) adjusted for variance
354 captured by protein assay plate were used for the interpretation of results of C3 and C4 protein
355 concentrations, for C3 protein concentration, 1 SD unit = 2.56 $\mu\text{g/L}$ ($3.60 \mu\text{g/L} \times \sqrt{1 - 0.49}$), and for
356 C4 protein concentration, 1 SD unit = 2.46 $\mu\text{g/L}$ ($3.33 \mu\text{g/L} \times \sqrt{1 - 0.45}$).

357 **Heritability and SNP-based heritability of the C3 and C4 protein concentrations**

358 The iPSYCH2012 cohort had 75,764 participants of European ancestry. 19,113 participants who
359 shared a genetic relatedness (entry of genetic relationship matrix, $r_{\text{GRM}} \geq 0.05$ with at least one
360 other individual were considered as relatives; 3,253 first degree ($r_{\text{GRM}} \geq 0.4$), 2,077 second degree
361 relatives ($0.2 \leq r_{\text{GRM}} < 0.4$) and 13,783 third degree relatives ($0.05 \leq r_{\text{GRM}} < 0.2$). We jointly estimated
362 both the heritability (h^2) and the SNP-based h^2 (h^2_{SNP}) of the C3 and C4 protein concentrations by
363 using the method proposed by Zaitlen et al⁴⁴. This method assumes a normal distribution of SNP
364 effect sizes. The GWAS studies of protein concentrations⁴⁵⁻⁴⁷ observed that SNPs in or near the
365 coding genes would capture more phenotypic variance than the remaining SNPs. Therefore, we used
366 two approaches to further explore the h^2_{SNP} , 1) estimating it using all common SNPs by BayesR⁴⁸,
367 which assumes a mixture distribution of SNP effect sizes and can be used to test number of SNPs
368 with nil, small, median and large genetic variance (small $R^2 < 0.01\%$, median $0.01\% \leq R^2 < 0.1\%$, and
369 large $0.1\% \leq R^2 < 1\%$) in addition to estimation of h^2_{SNP} , and 2) partitioning h^2_{SNP} into (a) $h^2_{\text{cis-chr}}$,
370 explained by SNPs on the chromosome where the coding gene (cis-chr SNPs) was positioned, and (b)

371 $h^2_{\text{trans-chr}}$, explained by the remaining SNPs (trans-chr SNPs). To partition h^2_{SNP} , we divided SNPs into
372 two subsets, 1) cis-chr and 2) trans-chr SNPs. The $h^2_{\text{cis-chr}}$ and $h^2_{\text{trans-chr}}$ were estimated by using GCTA-
373 GREML⁴⁹ and BayesR. Only genetically unrelated participants were included in the two analyses. The
374 genetic relationship matrix used in the Zaitlen and GREML analyses were estimated from 5,201,724
375 common SNPs. Only the subset of HapMap phase 3 (HM3) SNPs were included in the BayesR
376 analyses because of the computation complexity (853,129 HM3 SNP in total; for C3 protein
377 concentration, 132,239 cis-chr SNPs and 839,891 trans-chr SNPs; for C4 protein concentration,
378 60,631 cis-SNPs and 792,498 trans-SNPs). The Zaitlen method and GREML were implemented in
379 Genome-wide Complex Trait Analysis (GCTA)⁵⁰. BayesR was implemented in Genome-wide Complex
380 Trait Bayesian analysis (GCTB). The URLs for these programs are provided below.

381 We estimated the genetic correlation between C3 and C4 concentrations by BOLT-REML⁵¹. To further
382 examine if the genetic correlation was primarily driven by the two protein-coding genes (i.e. C3 and
383 C4), we conducted the BOLT-REML and Haseman-Elston regression⁵² analyses using the trans-chr
384 SNPs. The Haseman-Elston regression was implemented in GCTA.

385 ***GWAS of C3 and C4 protein concentrations***

386 We performed the GWAS analysis of the C3 and C4 protein concentrations by fastGWA⁵³. The
387 fastGWA is a LMM method which can include all individuals of European ancestry regardless of
388 relatedness. 5,201,724 imputed SNPs were analysed in the GWAS. In addition, fastGWA can include
389 candidate markers which optimizes power if particular SNPs capture a large proportion of the total
390 variance⁵⁴. Therefore, we excluded the SNPs in and near the coding gene for the required GRM in
391 the GWAS, C3: chr19, 4.67Mb – 8.74Mb, C4: chr6, 24.8Mb – 33.9Mb. For each GWAS, we fitted
392 birthyear, sex, wave (i.e., genotyping batch) and the first 20 PCs as covariates in the model. The PCs
393 were estimated by FastPCA⁵⁵, excluding the same SNPs as we did for the required GRM. We
394 conducted the GWASs using all SNPs on autosomal and sex chromosomes. SNPs on X chromosome
395 for males (coded as 0/2) were tested as diploid, assuming X chromosome of males has half dosage
396 compensation⁵⁶. We used GCTA-COJO⁵⁰ to identify the SNPs which were independently associated
397 with the two concentrations. We randomly sampled 10,000 participants from the population-based
398 sub-cohort of iPSYCH2012 as the LD reference cohort. The GWAS significance threshold was 5.0×10^{-8} .
399

400 To explore if the enrichment of mental disorder cases in the iPSYCH2012 case-cohort could induce
401 bias within the GWASs, we conducted simulations with ascertained individuals and performed
402 GWASs in the population-based sub-cohort (**Supplementary Method 4**).

403 ***Associations between C4 haplotypes and protein concentrations***

404 We examined the associations between the imputed C4 haplotypes and the two observed C3 and C4
405 protein concentrations. We first examined the associations of C4 copy numbers using a LMM
406 approach, $y_{\text{protein}} = x_{\text{copy}}b_{\text{copy}} + \sum x_c b_c + \sum z_{\text{-MHC}} u_{\text{-MHC}} + e$, where y_{protein} was C3/4 protein concentration;
407 x_{copy} was copy number of C4 allele, either C4A, C4B or HERV; b_{copy} represents effect of copy number;
408 x_c was covariate with b_c being its effect; Both b_{copy} and b_c were fixed effects. The covariates in the
409 model were the same as those fitted in the GWAS of C4 protein concentration. Because of the strong
410 linkage disequilibrium in the MHC region, fitting SNPs in the region is likely to underestimate the
411 effects of C4 allele count and reduce power. Therefore, we fitted the SNPs outside the MHC region
412 ($z_{\text{-MHC}}$) in the model with $u_{\text{-MHC}}$ being their random effects. In practice, the effect of copy number
413 from the linear mixed model could be estimated by generalized least squares (GLS) method, $b = (X^T V^{-1} X)^{-1} X^T V^{-1} y_{\text{protein}}$, where $X = \{x_{\text{copy}}, x_c\}$ and V was phenotypic covariance matrix of C3/4 protein
414

415 concentration. It was implemented by GCTA-GREML. All the individuals of European ancestry were
416 included in the analysis. The three *C4* allele counts were correlated. Therefore, we estimated the
417 joint effects using the same LMM approach as above, $y_{\text{protein}} = x_{C4A}b_{C4A} + x_{C4B}b_{C4B} + x_{\text{HERV}}b_{\text{HERV}} + \sum x_c b_c +$
418 $\sum Z_{\text{MHC}}u_{\text{MHC}} + e$. In the model, x_{C4A} , x_{C4B} and x_{HERV} represent respectively counts of *C4A*, *C4B* and *HERV*.
419 The remaining variables were defined as above. Secondly, we further examined the associations of
420 the imputed *C4* haplotypes using the LMM approach. Previous studies have reported a strong effect
421 for the *C4A* gene⁹, while the effect of *C4B* remains unclear⁵⁷. Therefore, we used 'BS' as the
422 reference haplotype to estimate joint effects of the remaining haplotypes. The regression model can
423 be expressed as, $y_{\text{protein}} = \sum x_{\text{allele}}b_{\text{allele}} + \sum x_c b_c + \sum Z_{\text{MHC}}u_{\text{MHC}} + e$, where x_{allele} represents *C4* haplotype.
424 Seven *C4* haplotypes were included in the model, except for BS. The remaining parameters were
425 defined as above. The estimated effect can be interpreted as the effect of the *C4* haplotype
426 compared to BS. All the European participants were included in the analysis. The significance
427 threshold for these analyses was the same as the main GWAS significance threshold (i.e., 5.0×10^{-8}).

428 ***FUMA and SMR***

429 We conducted gene-based analysis by Functional Mapping and Annotation of Genome-Wide
430 Association Studies (FUMA)⁵⁸. There were 18,305 genes available for the gene-based analysis, thus
431 the Bonferroni corrected threshold was 1.4×10^{-6} ($= 0.05 / (18,305 \times 2)$). We conducted Summary-
432 data-based Mendelian Randomisation (SMR)⁵⁹ to identify pleiotropic genes for C3 and C4
433 concentrations. For the SMR analysis, the eQTL data (i.e., summary statistics from associations of
434 gene expressions), was Genotype-Tissue Expression version 8 (GTEx v8)⁶⁰. The LD reference sample
435 with 10,000 participants was the same as the above GCTA-COJO analysis. 22,338 gene-tagged probes
436 within 49 tissues (200,144 probes in total) which had significant SNPs were included in the SMR
437 analysis. The Bonferroni significance threshold was 1.2×10^{-6} ($= 0.05 / (200,144 \times 2)$).

438 ***Associations between C4 haplotypes and mental disorders observed within the iPSYCH2012 case-*** 439 ***cohort study***

440 Based on the associations with protein concentrations, we conducted the associations between *C4*
441 haplotypes and 6 iPSYCH2012 disorders (SCZ, BIP, DEP, ASD, ADHD and AN). We used three
442 approaches to examine the relationships, 1) associations with *C4* allele counts, 2) associations with
443 imputed *C4* haplotypes, 3) associations with predicted *C4* gene expression in the brain. Because the
444 iPSYCH2012 case-cohort study has person-level data on the age-at-first contact with psychiatric
445 services, we were able to assess the risk of mental disorders within the more informative time-to-
446 event framework, using Cox proportional hazards regression (Cox PH) to analyse the hazards of *C4*
447 allele counts and haplotypes with respect to the mental disorder of interest. For *C4* allele count, we
448 examined the joint effects due to their correlations, $h(t) = h_0(t)\exp(x_{C4A}b_{C4A} + x_{C4B}b_{C4B} + x_{\text{HERV}}b_{\text{HERV}} +$
449 $\sum x_c b_c)$. In the model, $h_0(t)$ represents the baseline hazard while $h(t)$ represents the hazard at time t
450 between baseline and December 2016. The remaining variables were defined as above. For *C4*
451 haplotypes, we examined the joint effects using the Cox PH model, $h(t) = h_0(t)\exp(\sum x_{\text{allele}}b_{\text{allele}} +$
452 $\sum x_c b_c)$. All the variables were defined as above. In addition to *C4* haplotypes, we used the predicted
453 *C4A* and *C4B* gene expressions as outlined in the post-mortem brain study of Sekar et al.⁹ The
454 association was conducted with a Cox PH model, $h(t) = h_0(t)\exp(x_{C4A_predicted}b_{C4A_predicted} + x_{C4B_predicted}$
455 $b_{C4B_predicted} + \sum x_c b_c)$, where $x_{C4A_predicted}$ and $x_{C4B_predicted}$ represent the predicted *C4A* and *C4B* gene
456 expressions, respectively. We included only unrelated individuals of European ancestry in the three
457 analyses.

458 In the time-to-event analysis, the cases were the diagnosed participants by December 2016, and the
459 non-cases are defined as the entire cohort excluding those individuals with the disorder of interest.

460 Therefore, we defined six psychiatric-disorder samples for the time-to-event analyses. The sample
461 sizes for cases and non-cases are shown in **Supplementary Table 1**.

462 ***Associations between protein concentrations and mental disorders observed within the*** 463 ***iPSYCH2012 case-cohort study***

464 Based on the associations between *C4* haplotypes and 1) the two protein concentrations (*C3* and *C4*)
465 and 2) six mental disorders, we explored the associations between *C3* and *C4* protein concentrations
466 and mental disorders observed in the iPSYCH2012 case-cohort study, using Cox PH models. We first
467 tested the marginal effects of two concentrations, using the model $h(t) = h_0(t)\exp(x_{\text{protein}}b_{\text{protein}} +$
468 $\sum x_c b_c)$. The variables were defined as above. Due to the high correlation, we then fitted both
469 concentrations jointly, $h(t) = h_0(t)\exp(x_{C3_protein}b_{C3_protein} + x_{C4_protein}b_{C4_protein} + \sum x_c b_c)$, where $x_{C3_protein}$
470 and $x_{C4_protein}$ represent *C3* and *C4* concentration, respectively. $b_{C3_protein}$ and $b_{C4_protein}$, effects of two
471 protein concentrations, were fixed effects. In the three analyses, we included only unrelated
472 individuals of European ancestry.

473 ***Mendelian Randomisation analysis based on summary statistics***

474 We explored the relationships between protein concentrations and mental and autoimmune
475 disorders using the generalised summary-data-based Mendelian Randomisation (GSMR) method⁶¹.
476 Because of the possible link between *C3* and *C4* and brain function⁷, we also included two
477 neurodegenerative disorders—Alzheimer’s disease, and amyotrophic lateral sclerosis in these
478 analyses. Thus, there were 8 broadly-defined neuropsychiatric disorders (i.e., SCZ⁶², DEP⁶³, BIP⁶⁴,
479 ASD⁶⁵, ADHD⁶⁶, AN⁶⁷, Alzheimer’s disease⁶⁸, amyotrophic lateral sclerosis⁶⁹), and 6 autoimmune
480 disorders (i.e., multiple sclerosis⁷⁰, type-1 diabetes⁷¹, Crohn’s disease⁷², ulcerative colitis⁷²,
481 rheumatoid arthritis⁷³, systemic lupus erythematosus⁷⁴). The GWAS summary statistics for these
482 disorders were publicly available (additional details provided in **Supplementary Table 19**).
483 Unfortunately, detailed GWAS summary statistics for Sjögren's syndrome were not available. The
484 GSMR method was implemented in GCTA. The GSMR method includes options to exclude potentially
485 pleiotropic loci (HEIDI filtering). Similar to the methods for pleiotropy exclusion, it assumes fewer
486 pleiotropic SNPs than valid variants and MR estimates may fluctuate with the inclusion of those SNP
487 outliers⁶¹. Given that a substantial proportion of the variance in *C4* concentration was associated
488 with SNPs in the MHC region (including the *C4* gene), genetic variants used in the GSMR analysis may
489 be dominated by those positioned in or near the region. Therefore, the interpretation of results
490 based on the large effect loci within the MHC can be misleading. As a planned analyses, we ran the
491 GSMR method again based on summary statistics from the GWAS of *C4* concentration adjusted for
492 COJO SNPs within MHC region (as a covariate). In addition, in the presence of bidirectional GSMR
493 findings (i.e., evidence of both protein concentration impacting on phenotype of interest, and vice
494 versa; forward and reverse direction respectively), symmetrical (i.e., equivalent) effect sizes may
495 reflect the presence of pleiotropy. It is of note that effects of *C3/4* protein concentration on these
496 mental and autoimmune disorders by GSMR were equivalent to log odds ratio (OR) from logistic
497 regression if we have all the phenotypes in a cohort⁶¹. We, therefore, reported the logOR and 95% CI
498 in the results. The effects of mental and autoimmune disorders on *C3/4* protein concentration may
499 be required to be transformed for better interpretation⁷⁵. However, we only used the reserve GSMR
500 results to examine the presence of reverse causation and pleiotropy. Therefore, we reported the
501 GSMR estimates directly. The LD reference sample which included 10,000 participants was the same
502 as for the GCTA-COJO analysis. We conducted both forward (from *C3/4* concentration to disorder)
503 and reverse (from disorder to *C3/4* concentration) analyses. We set HEIDI-outlier threshold at 0.01
504 to filter horizontal pleiotropy. The GSMR Bonferroni corrected significance threshold was 1.9×10^{-3} (= $0.05 / (2 \times 13)$).
505

506
507

508 **PheWAS based on UK Biobank phenotypes**

509 Based on the GSMR analysis results, we conducted phenome-wide association studies (PheWASs) to
510 explore the relationships with disorder outcomes. The PheWASs were conducted in the UK Biobank
511 (UKB) cohort⁷⁶, a large population cohort with 487,409 participants of multiple ancestries. The
512 genotypes were imputed to the HRC³⁸ and UK10K⁷⁷ reference panels by the UKB group. The quality
513 controls were described in detail elsewhere⁷⁸, including generic ancestry determination, quality
514 controls of imputed SNPs, and estimation of principal components. In the study, we included
515 1,130,559 HM3 SNPs on autosomal chromosomes, with MAF ≥ 0.01 , HWE P -value $\geq 1.0 \times 10^{-6}$,
516 because only effects of HM3 SNPs were predicted by BayesR. The genetic relationship matrix was
517 estimated by GCTA. 347,769 unrelated participants of European ancestry were retained with genetic
518 relationship < 0.05 . In the PheWAS analysis, we included 1,148 UKB phenotypes, 1) 1,027 disorders
519 which were classified by ICD-10 codes, 2) 51 anthropometric measurements and brain imaging traits,
520 and 3) 70 infectious disease antigens. The quantitative traits were standardised by RINT to have
521 mean 0 and variance 1. We then used the model to test the associations, for quantitative traits, $y =$
522 $x_{\text{protein_prs}}b_{\text{protein_prs}} + \sum x_c b_c + e$, where y represents quantitative trait in UKB; $x_{\text{protein_prs}}$ represents
523 polygenic scores for neonatal C3 and C4 protein concentrations predicted by BayesR; x_c represent
524 the covariate variables including birth year, sex and 20 PCs. For dichotomous traits, $\text{logit}(y) =$
525 $x_{\text{protein_prs}}b_{\text{protein_prs}} + \sum x_c b_c + e$, where y represents the dichotomous trait and definitions of the
526 remaining variables were the same as above. In addition, we conducted the PheWAS analyses for
527 males and females separately using the same approach. Polygenic scores were predicted using
528 GWASs in both sexes. The Bonferroni corrected significance threshold was 7.3×10^{-6} ($= 0.05 / (1148 \times$
529 $3 \times 2)$).

530

531 **RESULTS**

532

533 The iPSYCH2012 study, a population-based case-cohort, was designed to study the genetic and
534 environmental factors of 6 mental disorders (**Supplementary Table 1**), SCZ, BIP, DEP, ASD, ADHD and
535 AN. The study included 80,873 individuals of multiple ancestries. 75,764 European individuals were
536 retained by principal components (PCs) projection (**Supplementary Figure 1**). The following analyses
537 in our study were based on the European individuals. We imputed C4 haplotypes using the reference
538 data^{9,14}. The imputation process predicted the counts of three C4 alleles (C4A, C4B and HERV), but is
539 not able to confidently distinguish all combinations. From the subset of imputation result without
540 ambiguity, C4A allele is more likely to carry a HERV than C4B allele (~ 1.5 higher) while C4A allele is
541 much less likely to not carry a HERV (**Supplementary Table 2**, for C4A, 0.2%, for C4B, 21.2%), which is
542 consistent with previous studies⁴⁰. Based on that, eight common C4 haplotypes were imputed with
543 allele frequency (AF) ≥ 0.01 (**Supplementary Table 3**). Their frequencies were respectively BS (12%),
544 AL (4%), AL-BS (23%), AL-BL (43%), AL-BS-BS (2%), AL-AL (11%), AL-AL-BS (3%), and AL-AL-BL (2%),
545 consistent with the published studies^{14,41}. We counted the copy numbers of the three types of C4
546 alleles (C4A, C4B and HERV, **Supplementary Figure 2**). The copy numbers (i.e., count) for the
547 different types of C4 alleles were correlated. C4A count was negatively correlated with C4B count (r
548 $= -0.52$, P -value $< 1.0 \times 10^{-100}$). HERV count was positively correlated with C4A count ($r = 0.73$, P -value
549 $< 1.0 \times 10^{-100}$), but negatively correlated with C4B count ($r = -0.17$, P -value $< 1.0 \times 10^{-100}$).

550

551 There were 68,768 participants of European ancestry with measures of C3 and C4 protein
552 concentrations. The distributions of the observed neonatal C3 and C4 protein concentrations were
553 right skewed, with mean, median, SD, and interquartile range being respectively 7.1, 6.7, 3.6, and
554 5.1 - 9.2 $\mu\text{g/L}$ for C3 protein concentration, and 6.9, 6.5, 3.3 and 4.9 - 9.0 $\mu\text{g/L}$ for C4 protein
555 concentration (**Supplementary Table 4**). Significant differences were observed in the protein
556 concentrations between males and females (C3 protein concentration: difference = -0.32, SE = 0.03,
557 P -value = 6.5×10^{-32} ; C4 protein concentration: difference = -0.30, SE = 0.03, P -value = 2.7×10^{-33}).

558 While the variance captured by sex was small ($R^2 = 0.19\%$ for C3 protein concentration and 0.20% for
559 C4 protein concentration), we fitted sex as a covariate in the following analyses. To account for the
560 influence of duration of storage (**Supplementary Figure 3**) and between-protein assay plate
561 variation, we regressed the concentrations of the plates using a linear mixed model (LMM)
562 approach. The residuals were standardised (mean 0, variance 1) using rank-based inverse normal
563 transformation (RINT). After standardisation, the concentrations of C3 and C4 were positively
564 correlated $r_p = 0.65$ (P -value $< 1 \times 10^{-100}$, **Supplementary Figure 4**).

565

566 **Heritability of C3 and C4 protein concentrations**

567 The h^2 of C4 by Zaitlen's method⁴⁴ was 40% (SE = 0.03, P -value = 2.7×10^{-44} , **Supplementary Table 5**)
568 while the h^2_{SNP} was 26% (SE = 0.006, P -value $< 1.0 \times 10^{-100}$). For C3, h^2 was 21% (SE = 0.03, P -value =
569 1.1×10^{-11}) and the h^2_{SNP} was 4% (SE = 0.005, P -value = 3.2×10^{-14}). The high genetic variance of C4
570 concentration was confirmed by BayesR⁴⁸ — the h^2_{SNP} for C4 was 24% (SE = 0.004, P -value 1.0×10^{-100})
571 and 6% (SE = 0.005, P -value = 4.4×10^{-39}) for C3. Moreover, ~ 50 HM3 SNPs captured substantial
572 genetic variance ($R^2 > 0.1\%$), for C3 protein concentration 39 SNPs and for C4 protein concentration
573 62 SNPs. The observation was consistent with our expectation that SNPs in or near the salient
574 encoding genes would capture a substantial proportion of phenotypic variance for these two
575 circulating proteins. Thus, we partitioned SNPs into two subsets; (a) those on the chromosome
576 where the coding gene is located (*cis*-chr SNPs), and (b) those on the remaining chromosomes
577 (*trans*-chr SNPs), and jointly estimating SNP-based h^2 at the two subsets of SNPs, SNP-based $h^2_{\text{cis-chr}}$
578 and $h^2_{\text{trans-chr}}$, using GREML⁴⁹. For C4, $h^2_{\text{cis-chr}} = 14\%$ (SE = 0.005, P -value $< 1.0 \times 10^{-100}$) and $h^2_{\text{trans-chr}} = 4\%$
579 (SE = 0.006, P -value = 9.4×10^{-13}). For C3, $h^2_{\text{cis-chr}}$ was 0.4% (SE = 0.001, P -value = 1.1×10^{-3}) and $h^2_{\text{trans-chr}}$
580 = 3% (SE = 0.006, P -value = 1.4×10^{-7}). The *cis*-chr SNPs of C4 concentration captured more genetic
581 variance than the *trans*-chr SNPs. A higher genetic variance was found at *cis*-chr SNPs by BayesR. For
582 C4, $h^2_{\text{cis-chr}} = 22\%$ (SE = 0.003, P -value $< 1.0 \times 10^{-100}$) and $h^2_{\text{trans-chr}} = 5\%$ (SE = 0.004, P -value = 6.2×10^{-30}).
583 For C3, $h^2_{\text{cis-chr}} = 2\%$ (SE = 0.001, P -value = 3.9×10^{-80}) and $h^2_{\text{trans-chr}} = 5\%$ (SE = 0.004, P -value = 5.9×10^{-26}).
584 Both the Zaitlen method and GREML assume genetic variance at each SNP follows a normal
585 distribution while BayesR assumes a mixture distribution. The BayesR estimates are expected to
586 better model the true underlying genetic architecture. We only used individual-level-data-based
587 methods (such as Zaitlen's method, GREML and BayesR) in the analysis as genotypes and
588 phenotypes were all available, since summary-level-data-based methods modelling LD from
589 reference cohorts are less accurate. In summary, both C3 and C4 concentrations were heritable
590 traits. SNPs positioned in the C4 gene accounted for a substantial proportion of the genetic variance
591 of C4 concentration.

592

593 The genetic correlation (r_g) between the two concentrations based on BOLT-REML⁵¹ was 0.38
594 (**Supplementary Table 6**, SE = 0.03, P -value = 1.9×10^{-35}), smaller than the phenotypic correlation
595 (0.65). Given the large genetic effects at the coding genes for C3 and C4 we then estimated r_g using
596 SNPs other than chromosomes 6 or 19 (related to the location of C4 and C3 genes, respectively) to
597 further investigate if the correlation was driven by *cis*-chr SNPs, $r_g = 0.82$ (SE = 0.05, P -value =
598 4.8×10^{-65}). The high correlation was confirmed by Haseman-Elston regression⁵² ($r_g = 0.78$, SE = 0.19,
599 P -value = 4.1×10^{-5}), using *trans*-chr SNPs. These results indicated that C3 and C4 were genetically
600 correlated, and this genetic correlation was not driven by SNPs in or near their respective encoding
601 genes.

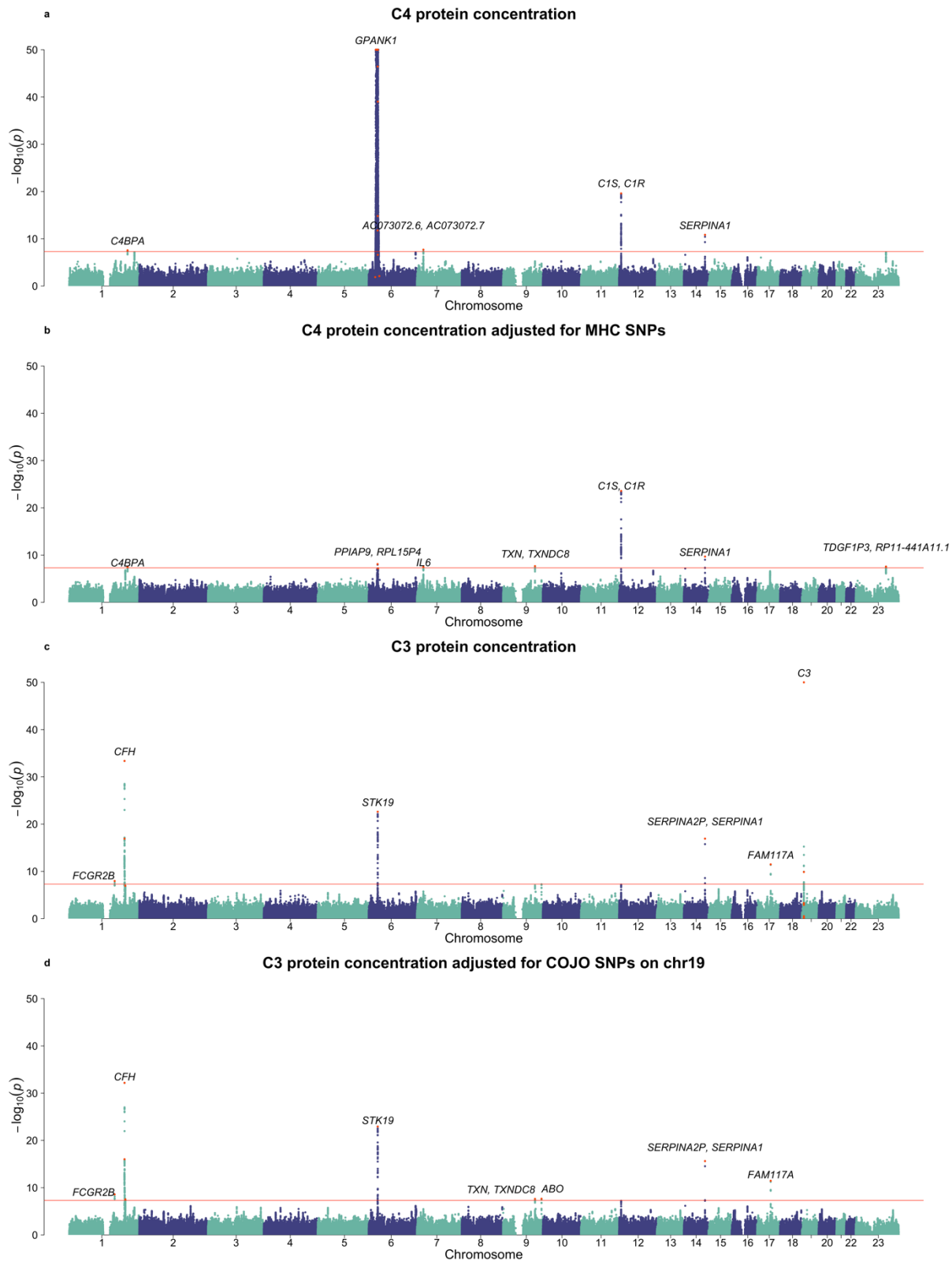
602

603 **GWAS of C3 and C4 protein concentrations**

604 We used fastGWA⁵³ to conduct the GWAS analysis based on 75,764 participants of European
605 ancestry and 5,327,833 common SNPs, 5,201,724 in autosomes and 126,109 on the X chromosome
606 (**Figure 2**). We conducted a GCTA-COJO⁷⁹ analysis to help identify putative independent SNPs. For C4
607 protein concentration, 34 autosomal SNPs were identified as genome-wide significant

608 (Supplementary Table 7), and all were autosomal. For C3 protein concentration, 14 significant SNPs
609 were identified, again all were autosomal (Supplementary Table 8).

610
611
612
613



614
615

616 **Figure 2** GWASs of neonatal C4 and C3 protein concentrations. a) unadjusted C4 protein
617 concentration, b) C4 protein concentration adjusted for COJO SNPs in the MHC region (fitted as
618 covariates in the regression model), c) unadjusted C3 protein concentration and d) C3 protein
619 concentration adjusted for COJO SNPs on chromosome 19. The COJO SNPs fitted as covariates in
620 GWAS of adjusted protein concentration (panels b and d) were identified from GCTA-COJO analysis
621 of unadjusted protein concentration. The COJO SNPs were highlighted with red colour. The top-
622 associated SNPs were annotated with their overlapped or nearest genes. The GWAS threshold was
623 5.0×10^{-8} .

624
625 Of the 34 SNPs significantly associated with C4 protein concentration, 30 (88.2%, 30/34) were found
626 on chromosome 6. Of these 29 were in the MHC region and 27 (79.4%, 27/34) SNPs were positioned
627 within 2Mb of the *C4* gene (chr6, 31.9 Mb). These 27 SNPs explained 16.7% of phenotypic variance in
628 C4 concentration, which is consistent with the estimated $h^2_{\text{cis-chr}}$. SNP rs113720465 (32,005,355bp,
629 ~1Kb away from *C4B-AS1* [32,000-32,004Kb]) had the largest effect size (the A allele was associated
630 with an increase of 0.76 standard deviation units of C4 protein concentration), however SNP
631 rs3117579 had the smallest *P*-value (within an exon of *GPANK1*). Given this large effect size, it is
632 possible that SNPs in LD at $r^2 < 0.01$ (the COJO threshold of independence) could also be reported as
633 genome-wide significant through correlation. Thus, we conducted a GWAS fitting the COJO SNPs in
634 and near MHC region as fixed effects (**Figure 2**). We identified 8 significant loci by COJO, 6 of which
635 were significant from GWAS of unadjusted C4 protein concentration. The 2 additional loci were on
636 chromosomes 9 (rs6477754) and X (rs12012736). Interestingly, nearly all the 8 COJO SNPs were
637 annotated to the genes biologically related to complement-related pathways (**Supplementary Figure**
638 **5**). For example, *C4BPA* (rs12057769) encodes a binding protein of C4. The *IL6* gene (rs2066992)
639 encodes a cytokine stimulated in response to infections and injuries. *C1S* (7.1Mb on chr12) and *C1R*
640 (7.2Mb on chr12), the nearest genes of rs11064501, are the protein-coding genes of two C1 sub-
641 components.

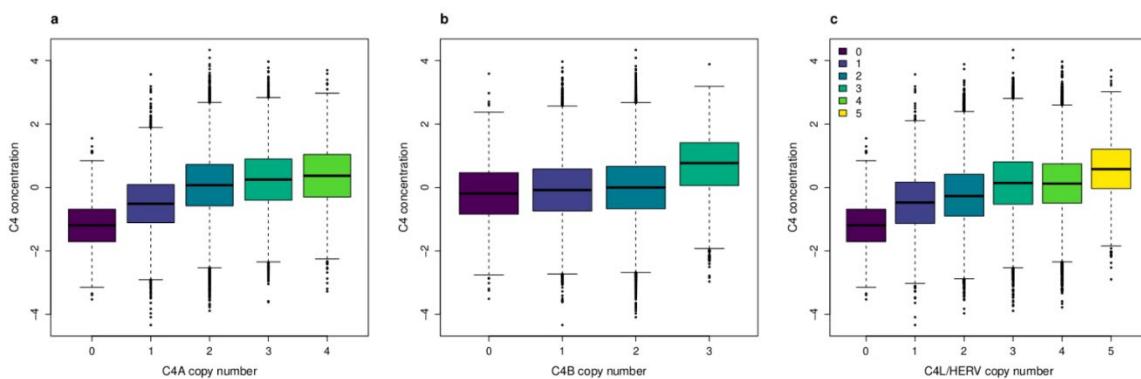
642
643 With respect to C3 protein concentration, 7 COJO SNPs were positioned within 2Mb of the *C3* gene
644 (chr19, 6.7Mb) — these loci explained 3% of phenotypic variance in C3 concentration. After fitting
645 these 7 COJO SNPs as covariates, 8 significant COJO SNPs were identified (**Supplementary Figure 6**).
646 We found a SNP within the *ABO* gene, which has recently been identified as a ‘master regulator’ of
647 plasma protein concentration^{46,80}. The gene annotations of the remaining SNPs encode proteins
648 which involve immune- and/or C3-related pathways: (a) *FCGR2B* (rs844), which encodes an
649 inhibitory receptor for the Fc region of immunoglobulin gamma (IgG), (b) *CFH* (rs558103 and
650 rs11580821) which encodes Complement Factor H, a key factor that inhibits the alternative pathway
651 and the amplification loop downstream from C3, (c) *STK19* gene (rs114492815) which is close to the
652 *C4A* gene, and (d) *FAM117A* (rs12949906), which has enhanced gene expressions in dendritic cells
653 (i.e., antigen-presenting cells involved in the immune system⁸¹).

654
655 For both GWASs of C3 and C4 protein concentration, we found no evidence of potential
656 ascertainment bias related to the enrichment of cases with mental disorders in the iPSYCH2012
657 case-cohort study (**Supplementary Method 4, Supplementary Figures 8-10**). Therefore, the
658 following post-GWAS analyses were based on the results from the full iPSYCH2012 sample.

659 660 **Associations between C4 haplotypes and C3/4 protein concentration**

661 Both h^2_{SNP} and GWAS results indicated strong effects of the SNPs in the MHC region for both C4 and
662 C3 protein concentrations. Due to the complex LD structure in this region, we used the imputed C4
663 haplotypes to investigate phenotypic associations of these genetic variants. We first examined the
664 associations between the imputed C4 haplotypes with the observed C4 protein concentration, using
665 an LMM approach. As expected, more copies of *C4* allele (either *C4A*, *C4B*, with or without *HERV*)
666 were strongly associated with higher C4 protein concentration (**Figure 3**). The *C4A* copy number

667 (b_{C4A}) had greater effect than $C4B$ (b_{C4B}) and $HERV$ (b_{HERV}), $b_{C4A} = 0.3$ (Supplementary Table 9, SE =
668 0.01, P -value $< 1.0 \times 10^{-100}$), $b_{C4B} = 0.2$ (SE = 0.01, P -value $< 1.0 \times 10^{-100}$), and $b_{HERV} = 0.2$ (SE = 0.004, P -
669 value $< 1.0 \times 10^{-100}$). The $C4$ copy numbers were correlated. Therefore, we fitted all 3 gene copy
670 numbers in a regression model to estimate the joint effects. The $C4A$ copy number had nearly
671 identical effect to $C4B$ copy number, $b_{C4A} = 0.6$ (SE = 0.01, P -value $< 1.0 \times 10^{-100}$), $b_{C4B} = 0.6$ (SE = 0.01,
672 P -value $< 1.0 \times 10^{-100}$). The beta estimates associated with the $HERV$ copy number was less than the
673 comparable estimates for $C4A$ and $C4B$, and was negatively associated with $C4$ protein
674 concentration, $b_{HERV} = -0.08$ (SE = 0.005, P -value = 5.0×10^{-51}). This may reflect the strong correlation
675 with $C4A$ ($r = 0.73$) and negative correlation with $C4B$ ($r = -0.17$). The result suggested 1 more copy of
676 $C4A$ or $C4B$ is likely to have 1.6 $\mu\text{g/L}$ ($\sim 0.6 \times \text{SD}$ unit) higher $C4$ protein concentration given the same
677 amount of $HERV$. We calculated the captured variance ($R^2 = s^2 b^2$) that were comparable between
678 the $C4$ copy numbers. In the formula, s^2 was variance of $C4$ copy number, analogous to variance of
679 allele count. Of interest, the s^2 of $C4A$ count was greater than $C4B$ count (Supplementary Table 3, s^2
680 = 0.55 for $C4A$ and 0.31 for $C4B$). Therefore, $C4A$ count had a larger contribution to $C4$ protein
681 concentration than the $C4B$ count ($R^2 = 23\%$ for $C4A$ and 11% for $C4B$). In total, both counts captured
682 17.3% of variance in $C4$ protein concentration, accounting for the negative correlation between the
683 two allele counts ($r = -0.52$). The captured genetic variance was in line with $h^2_{\text{cis-chr}}$, and the genetic
684 variance at the MHC SNPs. In summary, the imputed counts of both $C4A$ and $C4B$ were associated
685 with the observed $C4$ protein concentration and $C4A$ count had a greater contribution than $C4B$
686 count.
687
688



689
690

691 **Figure 3** Plot of $C4$ copy number versus $C4$ protein concentration. There were three $C4$ alleles, (a)
692 $C4A$, (b) $C4B$ and (c) $C4L/HERV$. The colours represent $C4$ allele counts.

693

694

695 Based on the effects of $C4$ allele counts, we then examined the association between the commonly
696 observed $C4$ haplotypes and $C4$ protein concentration. For this analysis, we used the BS haplotype as
697 the reference category (because of [1] the positive association between $C4$ allele count and $C4$
698 protein concentration, and [2] the greater contribution of $C4A$ count to $C4$ protein concentration).
699 All the remaining 7 common haplotypes were associated with increased $C4$ protein concentration.
700 Due to their higher frequencies, $AL-BS$ and $AL-BL$ haplotypes captured greater variance of $C4$ protein
701 concentration (Supplementary Table 10, $R^2 = 7.6\%$ for $AL-BS$ and 7.1% for $AL-BL$).

702

703 We then examined $C3$ protein concentration. In keeping with expectations, we did not identify any
704 significant associations between either $C4$ copy number or $C4$ haplotype, versus $C3$ protein
705 concentration. However, we found the $AL-BS$ haplotype was nominally significantly associated with
706 $C3$ protein concentration ($b_{AL-BS} = 0.23$, SE = 0.03, P -value = 5.4×10^{-3}). From the GWAS of $C3$

707 concentration, there was a COJO SNP positioned within the MHC region (rs114492815). While this
708 SNP was in very weak association with each of the *C4* allele counts ($R^2 < 0.005$ for *C4A* count and *C4B*
709 count, 0.01 for *C4L/HERV* count), there was a moderate association with *AL-B5* ($R^2 = 0.11$).
710 Therefore, we ran the analysis again, fitting rs114492815 as an additional covariate. After this
711 adjustment, none of the haplotypes were associated with the C3 concentration (**Supplementary**
712 **Table 10**). In general, C3 concentration was independent of C4 alleles.

713

714 **Functional mapping of GWAS**

715 Having found the significant SNPs from the GWASs, we explored the genes associated with both
716 concentrations. For C4 concentration, we identified 263 significant genes by Functional Mapping and
717 Annotation of Genome-Wide association Studies (FUMA), 257 (98%) on chromosome 6 and five (2%)
718 on the remaining chromosomes (**Supplementary Table 11**). These findings are consistent with the
719 high LD between loci in the MHC region and the high gene density in this region. Some
720 differentiation between genes was achieved by using SMR which integrates the trait associations
721 with significant eQTL associations. For SMR, using the eQTL summary data GTEx version 8, we
722 identified 56 pleiotropic genes, 55 on chromosome 6 and one on chromosome 1 (**Supplementary**
723 **Table 12**). We noted that the number of identified genes by SMR was smaller than by FUMA. Many
724 genes on chromosome 6 were significant on the SMR test but failed to pass HEIDI test (i.e., there
725 was evidence of pleiotropy). This was because of the complex LD and likely multiple causal alleles.
726 Interestingly, SMR analysis found that *C4A*, *C4B* and *C4BPA* were all significantly associated with
727 neonatal C4 protein concentration. These are three of the major genes involved in regulation of C4
728 protein concentration. The genetic correlates of neonatal C4 protein concentration were associated
729 with higher *C4A* gene expressions in 8 brain tissues, amygdala ($b_{XY} = 0.70$, SE = 0.11, P -value =
730 3.7×10^{-10}), anterior cingulate cortex ($b_{XY} = 0.74$, SE = 0.13, P -value = 6.8×10^{-9}), caudate basal ganglia
731 ($b_{XY} = 0.70$, SE = 0.09, P -value = 2.6×10^{-16}), cerebella hemisphere ($b_{XY} = 0.44$, SE = 0.05, P -value =
732 8.2×10^{-18}), brain cerebellum ($b_{XY} = 0.43$, SE = 0.04, P -value = 5.1×10^{-22}), hippocampus ($b_{XY} = 0.68$, SE =
733 0.10, P -value = 8.0×10^{-12}), hypothalamus ($b_{XY} = 0.79$, SE = 0.11, P -value = 4.4×10^{-12}), and putamen
734 basal ganglia ($b_{XY} = 0.76$, SE = 0.12, P -value = 8.6×10^{-10}). Except for cerebellar hemisphere and
735 cerebellum, the effect sizes of these associations were comparable with a mean of $b_{XY} = 0.73$.
736 Overall, the findings from FUMA and SMR indicate strong associations between genes in the MHC
737 region and C4 protein concentration, and the *C4A* gene was likely to have causal effect on C4 protein
738 concentration in brain tissues. Interestingly, these significant genes were enriched with the Kyoto
739 Encyclopedia of Genes and Genomes (KEGG) gene-sets of systemic lupus erythematosus (SLE, P -
740 value = 1.1×10^{-70}) and complement systems (P -value = 7.2×10^{-4}) (**Supplementary Figure 7**). For C3
741 protein concentration, we identified 19 genes by FUMA (**Supplementary Table 13**). Within this set of
742 genes, the *DXO* gene (chr6: 31.9Mb), which is positioned close to rs114492815, a significant SNP
743 from the C3 GWAS, passed both SMR and HEIDI (a test of pleiotropy) tests (**Supplementary Table**
744 **14**). All these genes by FUMA and SMR were enriched with KEGG gene-sets of SLE (P -value = 3.5×10^{-6}),
745 leishmania infection (P -value = 5.0×10^{-6}), complement systems (P -value = 9.1×10^{-6}) and Fc gamma
746 R-mediated phagocytosis (P -value = 9.5×10^{-4}) (**Supplementary Figure 7**). In general, the identified
747 genes (especially those within MHC region) suggest an association between SLE and the complement
748 systems.

749

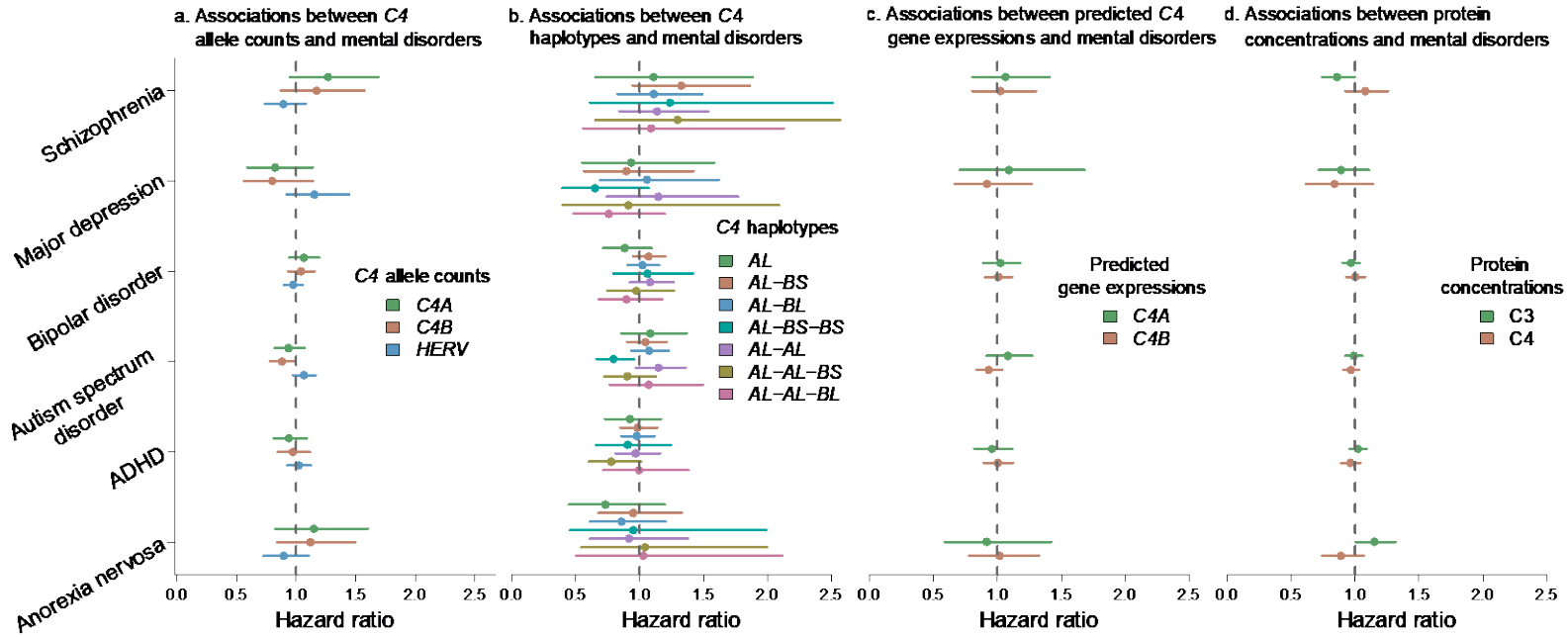
750

751 **Associations with mental disorders within the iPSYCH2012 case-cohort study**

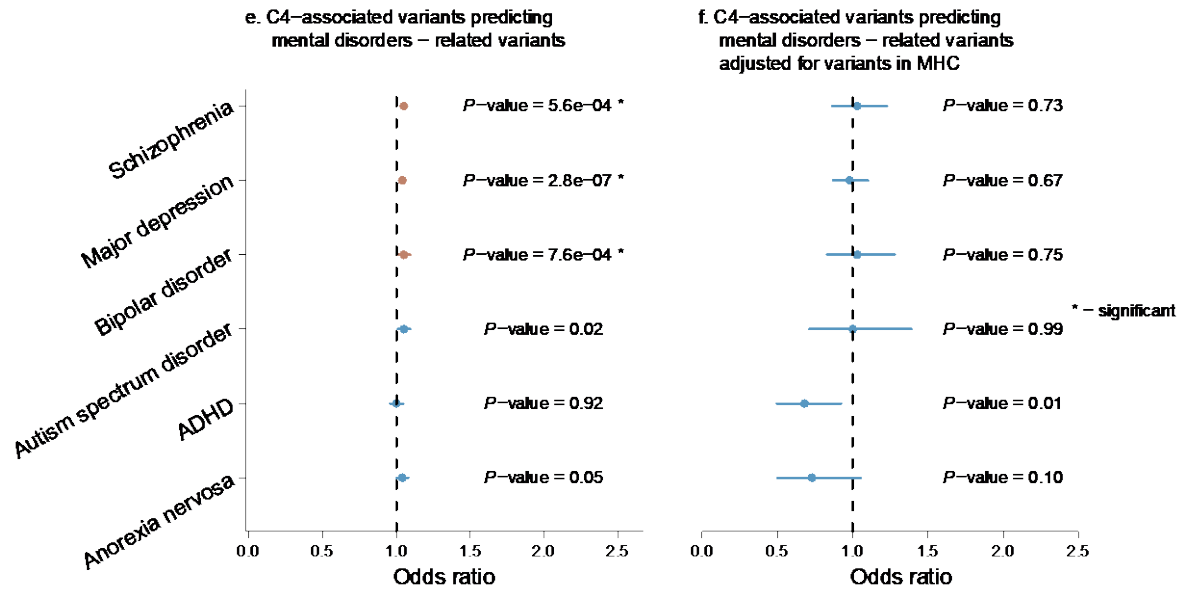
752 We did not identify significant associations between *C4A*, *C4B*, or *HERV* copy numbers
753 (**Supplementary Table 15** or C4-related haplotypes (**Supplementary Table 16**) and any of the six
754 mental disorders. Based on the formula between imputed C4 haplotypes and observed C4 gene
755 expression (i.e., RNA concentration) in post-mortem brain tissue⁹, we found no significant
756 associations between these estimates and any of the 6 mental disorders (**Supplementary Table 17**).
757 Importantly, we did not find any significant associations between observed neonatal C4 protein

758 concentration and any of the six mental disorders. In models that accounted for the strong
759 correlation between C3 and C4 concentration, we found no significant association between C3
760 concentration and any of the six mental disorders. (**Supplementary Table 18**).

Cox PH regression within the iPSYCH2012 case-cohort study



Mendelian randomisation based on summary statistics



762 **Figure 4** Association between C4-related measures and mental disorders, and GSMR analyses with
763 mental disorders. There were 6 mental disorders in the analyses: SCZ, DEP, BIP, ASD, ADHD and AN.
764 The results shown in the top row were from time-to-event analyses between mental disorders and
765 C4-related genotypes and phenotypes, including a) C4 allele counts, b) imputed C4 haplotypes, c)
766 predicted C4 gene expressions and d) C3 and C4 protein concentrations. The analyses were
767 conducted in the iPSYCH2012 cohort. The results from Mendelian randomisation analyses
768 (conducted by GSMR) were shown in the bottom row. The GSMR analysis using GWAS summary
769 statistics predicted relationships between C4 protein concentration and mental disorders, e) using
770 genetic variants from GWAS of C4 protein concentration, f) using genetic variants from GWAS of C4
771 protein concentration adjusted for COJO SNPs in the MHC region. The Bonferroni corrected
772 thresholds were provided in Methods. Bars represent 95% confidence interval. Significant results
773 were highlighted with “*”.

774

775 **GSMR relationships with candidate neuropsychiatric and autoimmune disorders**

776

777 We conducted Mendelian randomization analyses to examine relationships between the two protein
778 concentrations (C3 and C4) and neuropsychiatric and autoimmune disorders (**Supplementary Table**
779 **19**). The results are shown in **Figures 4 and 5**. In the unadjusted analysis (i.e., all loci including the
780 MHC region; with and without HEIDI filtering), higher C4 protein concentration was found to be
781 associated with three mental disorders (**Figure 4 and Supplementary Figure 11**, SCZ, DEP, BIP). The
782 odds ratios for these three findings were small (1.05 or less). We found that these GSMR results
783 were strongly dependent on SNP instruments that were in and near the MHC region (e.g., for
784 schizophrenia and bipolar disorder 126 out of 130 SNPs, and for major depression 103 out of 107
785 SNPs). These three findings did not persist in the analyses adjusted for the MHC region SNPs. Overall,
786 these analyses do not lend weight to the hypothesis that C4 is causally related to the risk of the
787 psychiatric disorders included in the analyses.

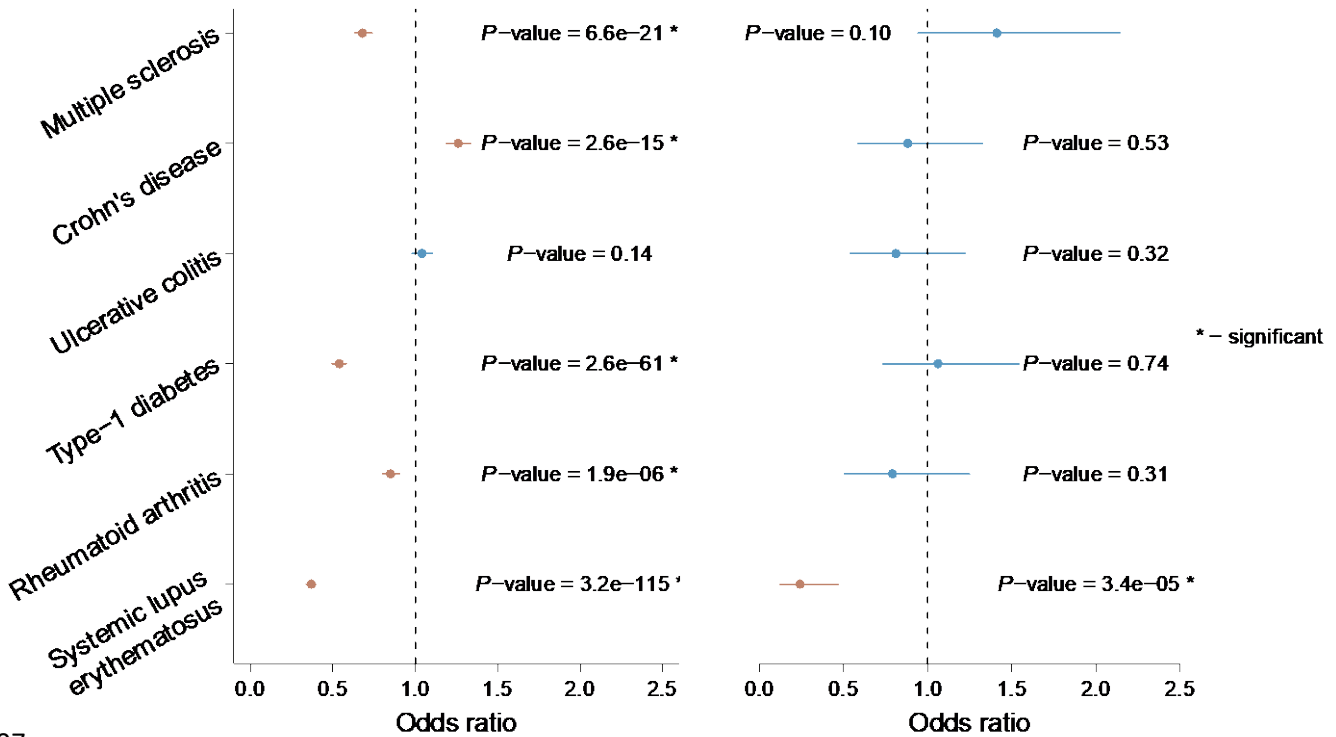
788

789 In contrast, we found strong, protective effects of C4 concentration for several autoimmune
790 disorders. Higher C4 concentration was associated with lower risks of multiple sclerosis, type-1
791 diabetes, rheumatoid arthritis and systemic lupus erythematosus (**Figure 5**). The effects were very
792 large, especially for type-1 diabetes (OR = 0.54, 95% confidence interval [CI] = 0.50 – 0.58, $N_{\text{SNP}} = 47$)
793 and systemic lupus erythematosus (OR = 0.37, 95% CI = 0.34 – 0.42, $N_{\text{SNP}} = 103$) (**Supplementary**
794 **Table 20**). We identified that higher C4 concentration increased the risk of Crohn’s disease (OR =
795 1.26, 95% CI = 1.19 – 1.34, $N_{\text{SNP}} = 86$). The strong effect of neonatal C4 protein concentration on
796 these disorders were not caused by reverse causation (**Supplementary Table 21 and Supplementary**
797 **Figure 12**). After removing the pleiotropic SNPs, genetic variants associated with all autoimmune
798 disorders (except for type-1 diabetes) were not associated with C4 protein concentration in the
799 reverse GSMR analysis (from autoimmune disorder to neonatal C4 protein concentration). When we
800 examined the relationships adjusted for the MHC region SNPs, the significant association with SLE
801 persisted. The effect size was comparable to that found using unadjusted C4 GWAS (with
802 adjustment, OR = 0.24, 95% CI = 0.12 – 0.47, $N_{\text{SNP}} = 7$; without adjustment, OR = 0.37, 95% CI = 0.34 –
803 0.42, $N_{\text{SNP}} = 103$). Overall, these findings further support the hypothesis that higher C4 protein
804 concentration is causally related to a reduced risk of systemic lupus erythematosus—it is predicted
805 that an increase of 2.46 $\mu\text{g/L}$ (1 SD unit) of C4 concentration would be associated with a 76%
806 reduced risk (1 – 0.24) of systemic lupus erythematosus.

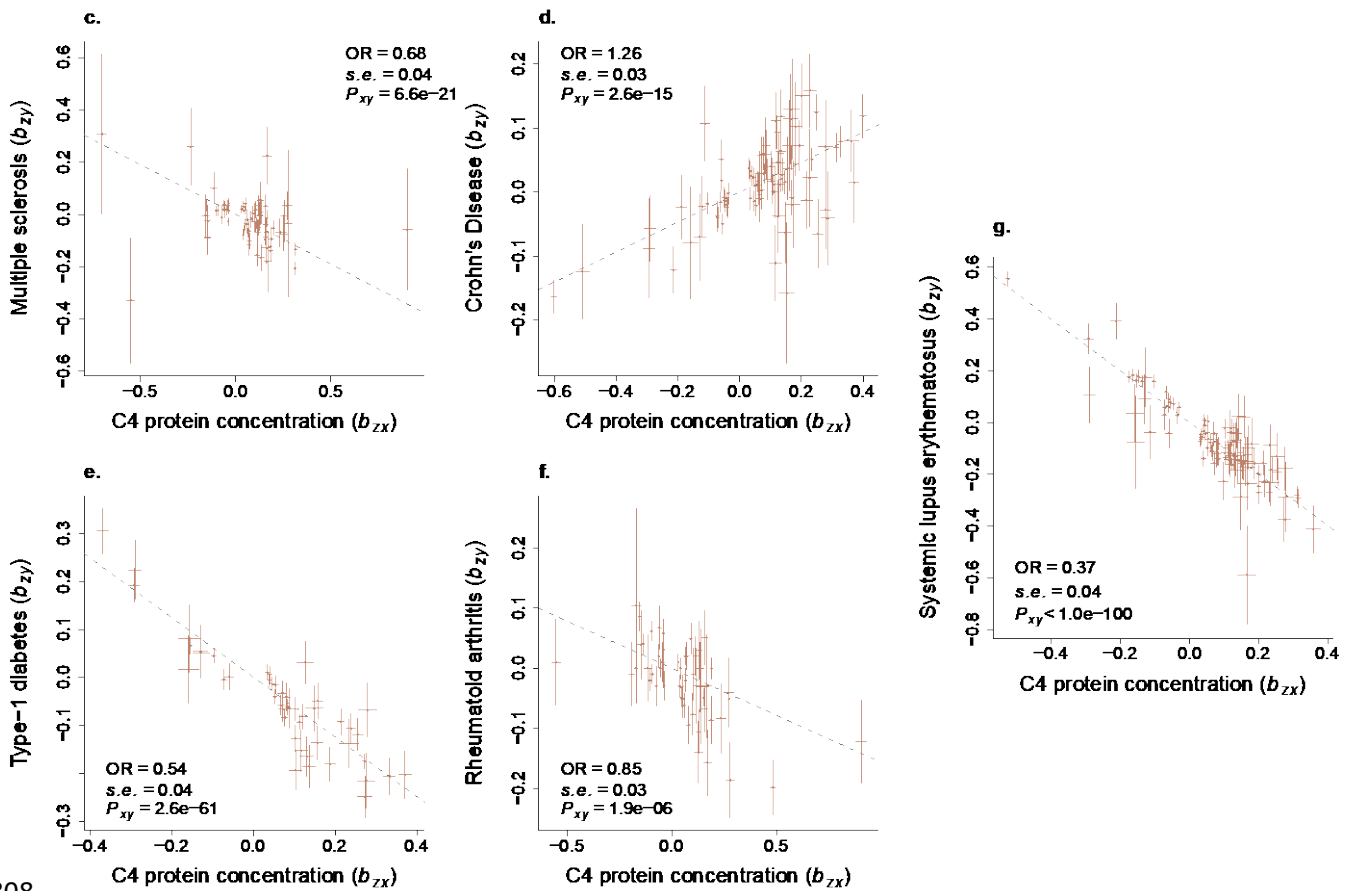
Mendelian randomisation based on summary statistics

a. C4-associated variants predicting autoimmune disorders – related variants

b. C4-associated variants predicting autoimmune disorders – related variants adjusted for variants in MHC



807



808

809 **Figure 5** GSMR analyses of C4 protein concentration against autoimmune disorders. The top two
810 panels showed the Mendelian randomisation results a) using genetic variants from GWAS of C4
811 protein concentration, b) using genetic variants from GWAS of C4 protein concentration adjusted for
812 COJO SNPs in the MHC region. The Bonferroni corrected threshold was 1.8×10^{-3} . Bars in the top 2
813 panels represented 95% confidence interval. Significant results were highlighted with “*”. The lower
814 five panels show the effects of genetic variants for C4 protein concentration without adjustment
815 against effects for five autoimmune disorders. Significant findings were identified for multiple
816 sclerosis, Crohn's disease, type-1 diabetes, rheumatoid arthritis and systemic lupus erythematosus.
817 The estimates in the panels (OR, standard error and *P*-value) were estimated from GSMR. Slope of
818 dash line represented logOR. Bars in the bottom 5 panels represented standard errors of SNPs from
819 GWAS. Potential pleiotropic SNPs were excluded.
820

821 We then explored the relationships between C3 concentration and neuropsychiatric and
822 autoimmune disorders by bidirectional GSMR. Mindful that analyses based on fewer instruments
823 may be underpowered to detect small effects, no significant associations were identified with
824 pleiotropic SNPs removed (**Supplementary Table 20**). Our findings provide no support for the
825 hypothesis that C3 protein concentration is related to the risk of the neuropsychiatric nor
826 autoimmune disorders examined in this study.

827 **C3 and C4 phenome-wide association studies in the UK Biobank**

828 With respect to C4 concentration, the PheWAS study in the UK Biobank identified significant
829 associations with 35 phenotypes (**Supplementary Table 22 and Supplementary Figure 13**). Many of
830 these were related to autoimmunity. Of the top 8 disease associations ranked by *P*-value, higher C4
831 concentration was associated with a reduced risk of six disorders, two were associated with an
832 increased risk. The top 8 were intestinal malabsorption (which includes coeliac disease; ICD10 = K90,
833 OR = 0.54, 95% confidence interval [CI] = 0.53 - 0.56); thyrotoxicosis [hyperthyroidism] (ICD10 = E05,
834 OR = 0.77, 95% CI = 0.75 - 0.79); hypothyroidism (ICD10= E03, OR= 0.92, 95% CI = 0.90 – 0.93);
835 insulin-dependent diabetes mellitus (ICD10 = E10, OR = 0.80, 95% CIs = 0.78 – 0.83); sarcoidosis
836 (ICD10 = D86; OR = 0.79, 95% CIs = 0.75 – 0.83); psoriasis (ICD10 = L40; OR = 1.08, 95% CIs = 1.06 –
837 1.10); systemic lupus erythematosus (ICD10 = M32, OR = 0.74, 95% CI = 0.69 - 0.80) and ankylosing
838 spondylitis (ICD10 = M45, OR = 1.22, 95% CIs = 1.16- 1.28). We also found a significant result for
839 multiple sclerosis (ICD-10 = G35, OR = 0.88, 95% CI = 0.84 - 0.92). The attenuated results for the
840 autoimmune disorders that were included in GMSR may be related to the low prevalence of these
841 disorders (and the smaller number of cases in the UKB). In addition to the disorders which met the
842 Bonferroni-corrected *P*-value threshold (*P*-value < 7.3×10^{-6}), we note that Sjögren's syndrome (ICD10
843 = M35; OR = 0.95, 95% CI = 0.92 - 0.97) was nominally significant. There were no significant findings
844 between C4 and any neuropsychiatric disorders (ICD10 F codes). No significant difference was found
845 between males and females. Overall, these findings lend weight to the hypotheses that neonatal C4
846 protein concentration is (a) not associated with the risk of neuropsychiatric disorders, but (b) is
847 associated with reduced risks of several autoimmune disorders. There were no significant
848 associations between C3 and any of the 1,148 phenotypes, which was in line with our GSMR findings
849 (**Supplementary Table 23 and Supplementary Figure 14**).

850 851 **DISCUSSION**

852 Our findings provide new insights into the genetic architecture of C3 and C4. Reassuringly, we found
853 a robust association between C4-related haplotypes (including copy number) and neonatal C4
854 protein concentration. The C3 and C4 protein concentrations were phenotypically and genetically

855 correlated ($r_p = 0.65$, P -value $< 1 \times 10^{-100}$; $r_g = 0.35$, P -value = 1.9×10^{-35} , using all SNPs; and $r_g = 0.78$, P -
856 value = 4×10^{-5} using trans-chr SNPs). The C3 and C4 GWAS findings identified variants in genes that
857 encode important proteins within the inter-connected complement pathways. In contrast to a
858 previous study⁹, we found that neither a higher imputed *C4* haplotype count nor a higher observed
859 *C4* protein concentration was associated with an increased risk of schizophrenia or any other mental
860 disorder diagnosed in later life. We did, however, find evidence from Mendelian randomization
861 studies that support the hypothesis that *C4* protein concentration is associated with a range of
862 autoimmune disorders. In models that incorporate the correlation between *C4* and *C3* protein
863 concentrations, we found no association between *C3* and an altered risk of any mental disorder nor
864 autoimmune disorder. This following discussion will focus on 5 key findings.

865 First, with respect to *C4* protein concentration, we found a stronger contribution of the *C4A* count
866 compared to the *C4B* count. In keeping with prior smaller studies^{18,19}, we confirmed that the copy
867 numbers of *C4A* and *C4B* were robustly positively associated with the concentration of the *C4*
868 protein (*C4A* count; $b = 0.33$, $SE = 0.005$, P -value $< 1 \times 10^{-100}$; *C4B* count; $b = 0.18$, $SE = 0.007$, P -value $<$
869 1×10^{-100}). In joint analyses that accounted for the pattern of correlations between the different types
870 of *C4* allele counts, the *C4A* count had twice the contribution to overall *C4* concentration compared
871 with *C4B* count.

872 Second, the protein concentrations of *C3* and *C4*, which are key components of the complement
873 initiation pathways⁸², were highly heritable. Both pedigree- and SNP-based h^2 estimates (standard
874 error) were appreciable for *C4* (0.40 (0.03) and 0.26 (0.006) respectively). The same estimates for *C3*
875 were smaller; 0.21 (0.03) and 0.04 (0.005) respectively. As expected, SNPs within and near the
876 respective coding genes (*C4*, *C3*) contributed to more than half of the genetic variance of their
877 related proteins.

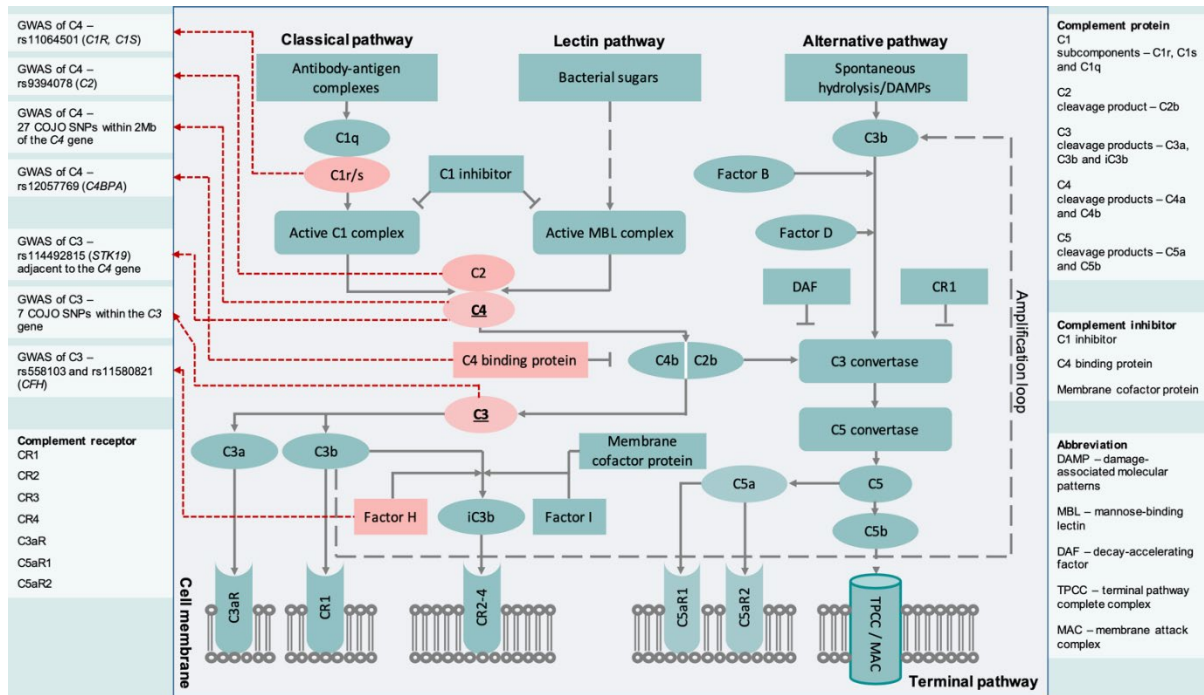
878 Third, our sample sizes for *C3* and *C4* concentrations GWASs were nearly twenty times larger than
879 the only published GWAS for these proteins¹⁸. With respect to *C4* protein concentration, our series
880 of linked GWASs provide important information about the genetic correlates of *C4* protein
881 concentration that lie outside of the MHC complex. In the GWAS, 30 quasi-independent hits were on
882 chromosome six, within the MHC region. Six additional loci were found on chromosomes 1, 7, 9, 12,
883 14 and X. We identified a locus on chromosome 1 within *C4BPA*, which encodes *C4* binding protein
884 (closely involved in *C4* protein regulation). The locus on chromosome 12 (rs11064501) is adjacent to
885 two genes that encode proteins involved in complement cascade initiation (*C1s*, *C1R*). Interestingly,
886 a locus (rs12012736) was identified on the X chromosome. This locus may be one of the factors that
887 contributed to the small sex differences found for the *C3* and *C4* protein concentrations and to the
888 known sex differences in the risk of autoimmune disorders⁸³.

889 With respect to *C3* protein concentration, apart from loci within the *C3* gene (seven quasi-
890 independent loci within this gene on chromosome 19), we found a locus within *FCGR2B* (Fc gamma
891 receptor IIb), which encodes a receptor for the Fc region of immunoglobulin gamma complexes
892 (IgG). The IgG complex forms part of the machinery required for the phagocytosis of immune
893 complexes. One locus in the MHC complex was identified, which is adjacent to the *C4A* gene. In
894 keeping with the prior GWAS¹⁸, we identified two loci within *CFH*, the gene that encodes
895 complement factor H. This protein is involved in complement regulation and has been linked to
896 several disease phenotypes (most notably, with age-related macular degeneration)⁸⁴. *CHF*
897 specifically regulates *C3*, which slows the downstream complement activation. We also found a locus
898 within *ABO*, which was identified as having associations with over 50 other protein
899 concentrations^{46,80}, thus variants in this gene could directly or indirectly influence generic protein

900 metabolic pathways (e.g., upstream metabolic steps, downstream protein degradation and
 901 excretion). In summary, our study has highlighted how genetic variants within several components
 902 of the complement cascades (i.e., at the systems level) could influence the concentration of key
 903 circulation proteins such as C3 and C4. We have summarized these findings in **Figure 6**.

904

905



906

907 **Figure 6** Summary of the results from GWASs of neonatal C3 and C4 protein concentrations
 908 displayed within the complement cascade. For significant loci identified from COJO, proteins
 909 encoded by annotated genes were highlighted with red colour.

910

911 Fourth, convergent evidence found no association between several *C4*-related measures and risk of
 912 SCZ. Our study measured *C4* protein concentration in 68,768 neonates, an age more proximal to the
 913 period of brain development consistent with the impact of *C4* expression and synaptic pruning^{6,24}.
 914 The strong association between *C4*-related copy number and measured *C4* protein concentration,
 915 and the biologically-plausible loci/genes identified in the *C3* and *C4* GWASs lend weight to the
 916 validity of our measures. We found that none of the following variables were associated with an
 917 altered risk of SCZ: (a) observed neonatal *C4* concentration, (b) copy numbers of either *C4A*, *C4B*, or
 918 *HERV*, (c) major *C4*-related haplotypes, nor (d) imputed brain *C4A* RNA expression. Furthermore,
 919 there were no associations between these *C4*-related variables and any of the other 5 iPSYCH target
 920 psychiatric disorders. We also note that the PheWAS study found no significant associations
 921 between the summary statistics of *C4* protein concentration and the UK Biobank-measured brain
 922 volumes (n = 28,613). Reassuringly, we note that the SMR analyses identified (a) *C4A* gene
 923 expression was strongly linked to *C4* protein concentration, and (b) higher *C4* neonatal protein
 924 concentration was associated with increased *C4A* gene expression in brain tissue (including
 925 amygdala, anterior cingulate cortex, caudate basal ganglia, cerebellum, hippocampus, hypothalamus
 926 and putamen basal ganglia). These particular findings support the hypothesis that the loci we
 927 observed in the GWAS not only influence circulating *C4* protein concentration (i.e., as measured in

928 the neonatal dried blood spots), but may also influence the expression of the *C4* gene in the brain. In
929 summary, we found no evidence to support the hypothesis that *C4*-related variables were causally
930 related to the risk of SCZ, nor the other mental disorders included in the iPSYCH case-cohort study.
931 Our findings allow us to refine future directions for schizophrenia research. We note that the most
932 strongly associated SNPs in the two most recent SCZ GWASes (rs1233578 [28,712,247bp]¹⁰ and
933 rs140365013 [27,523,869bp]⁶²) are both over 3Mb upstream from the *C4A* gene (31,95 – 31.97Mb).
934 There are 598 known genes^{87,88} (including 417 protein-coding genes) annotated between
935 rs140365013 and the *C4A* transcription start site, so this region should provide fertile grounds for
936 the generation of new candidate genes that might explain the top hits in recent SCZ GWASes.

937 Fifth, we found convergent evidence linking higher *C4* protein concentration and a reduced risk of
938 several autoimmune disorders, and an increased risk for Crohn's disease. Based on Mendelian
939 randomization analyses, there was robust evidence with respect to a lower risk of systemic lupus
940 erythematosus. The protective effect remained significant when we adjusted *C4* concentration for
941 MHC SNPs (from *C4* with MHC SNPs to SLE: OR = 0.37, 95% CI = 0.34 – 0.42; from *C4* adjusted for
942 MHC SNPs to SLE: OR = 0.24, 95% CI = 0.12 - 0.47). Evidence also emerged for a lower risk of type-1
943 diabetes, multiple sclerosis, and rheumatoid arthritis. Reassuringly, the UKB-based PheWAS found
944 variants associated with increased neonatal *C4* protein concentration were associated with (a)
945 reduced risks of a wide range of disorders (including coeliac disease, thyrotoxicosis, hypothyroidism,
946 type 1 diabetes, sarcoidosis, SLE, nephrotic syndrome, and multiple sclerosis; Sjögren's syndrome
947 was nominally significant), and (b) increased risks of several disorders (including psoriasis, ankylosing
948 spondylitis, iridocyclitis; Crohn's disease was nominally significant). Our findings are consistent with
949 a meta-analysis based on 16 case-control studies, where low *C4* gene copy number (<4) was
950 associated with an increased risk of any type of autoimmune disorder, including systemic lupus
951 erythematosus¹⁵. A study based on large-scale genetic and transcriptomic datasets by Kim et al.¹⁷
952 suggested that *C4A*-related gene expression was not associated with risk of schizophrenia-related
953 synaptic gene expression, but was associated with disorders including inflammatory bowel disease,
954 rheumatoid arthritis, and lupus. Our findings support these conclusions. Recently, it was reported
955 that variants in *C4A* and *C4B*, which were thought to increase the risk for SCZ, are protective for two
956 autoimmune disorders (systemic lupus erythematosus, Sjögren's syndrome)¹⁴. We also observed
957 that higher *C4* protein concentration was associated with increased risks of several autoimmune
958 disorders. The mechanisms of action underpinning the pattern of increased and decreased risk of
959 different autoimmune disorders remain poorly understood^{87,88}. In summary, our findings provide
960 convergent evidence to support the hypothesis that *C4* protein concentration is associated with the
961 risk of a range of autoimmune disorders.

962 Many GWASs have found links between loci in the MHC region and risk of autoimmune disorders⁸⁹⁻
963 ⁹⁴. Until recently, this has been interpreted as a connection between HLA genes and autoimmune
964 disease. Recently, Kamitaki et al. have shown that the link between the MHC locus and SLE and
965 Sjögren's may be explained by *C4A-C4B* allelic variance¹⁴ within the MHC region, thus expanding on
966 smaller studies linking low copy number of *C4*^{15,95} and *C4A*⁹⁶ with higher risk of systemic lupus
967 erythematosus. This raises the possibility that the correlation between the MHC region and the other
968 above-mentioned autoimmune diseases are also explained by the *C4A-C4B* allelic variance. Bian et
969 al.⁴¹ observed strong linkage disequilibrium with *HLA* alleles and BS (one of the *C4* alleles).
970 Regardless of these speculations, we found no evidence of causal relationships between (a) *C4* copy
971 number and *C4* haplotypes, (b) predicted *C4A* and *C4B* gene expressions, and (c) *C4* protein
972 concentration versus SCZ and a range of mental disorders. We identified pleiotropy between *C4*
973 concentration and three mental disorders (SCZ, DEP, BIP) from GSMR analyses. The complex linkage
974 disequilibrium between the *C4* gene and other genes in MHC region (including *HLA* genes) suggests

975 that we should be cautious when interpreting genotype to phenotype associations for loci within the
976 MHC region.

977 **Strengths and Limitations of the study**

978 Our study has several strengths. Our sample was nearly 20 times larger than the only other
979 published GWAS of C3 and C4¹⁸. With respect to the hypothesis linking complement to brain
980 development, our complement assays were collected from neonatal samples (versus adult samples).
981 Because the onset of mental disorders such as SCZ is often in the second and third decade of life, our
982 samples are unlikely to be impacted by reverse causation (e.g., smoking may be linked to
983 complement gene expression in the brain¹⁷) and medication effects may impact on post-mortem
984 gene expression studies⁹⁷. With respect to limitations, because our samples were based on neonatal
985 C3 and C4 concentrations, it remains to be seen if the genetic correlates we identified are stable
986 across the lifespan. Also, we used an antibody that has been demonstrated to measure total C4 (i.e.,
987 both C4A and C4B), thus we are unable to isolate the concentrations of the two isoforms. The C3 and
988 C4 concentrations in our study were derived from circulating plasma proteins, whereas the
989 concentration of these proteins may vary between organs/tissues and also in response to local tissue
990 activation pathways.

991
992 The study by Sekar et al.⁹, which was based on C4 haplotypes imputed from 28,799 schizophrenia
993 cases and 35,986 controls, found that the *AL-AL* haplotype was associated with an odds ratio of 1.27
994 compared with the *BS* allele (i.e. a 27% increased odds). Our study, based on 2,517 cases and 51,799
995 non-cases, and using a more informative imputation training set⁹, found no significant association
996 for this comparison (HR = 1.14, 95% CI = 0.84 - 1.54; **Supplementary Table 16**). Because the Sekar et
997 al. study included more schizophrenia cases in their analyses, it is feasible that our study was
998 underpowered to confidently detect an association between C4 haplotypes and SCZ⁹. However, we
999 had access to the concentrations of protein product of these haplotypes from a very large sample
1000 size (n = 68,768), and thus we could estimate the variance of neonatal C4 protein concentration.
1001 Based on this observed variance, our study had 80% power to confidently detect a 25% increased
1002 risk of SCZ (OR = 1.25) by 2.46 µg/L C4 protein concentration (1 standard deviation unit). Thus, our
1003 study had sufficient power to detect the effect size previously identified by Sekar et al.⁹.

1004 1005 **Conclusions**

1006 Our study provides new insights into the genetic and phenotypic correlates of C3 and C4 protein
1007 concentration and helps unravel the contribution of different C4-related copy numbers and
1008 haplotypes to C4 protein concentration. Based on convergent evidence, we found no evidence to
1009 support an association between C4-related measures and risk of SCZ, nor other psychiatric disorders.
1010 Mindful of the pitfalls of linking genotypes with phenotypes within the MHC region, we encourage
1011 the research community to continue to actively explore additional candidate loci for schizophrenia
1012 within this region. In contrast to our findings regarding mental disorders, convergent evidence
1013 emerged supporting an association between C4 protein concentration and risk of autoimmune
1014 disorders.

1015

1016 **URLs**

1017 PLINK2: <https://www.cog-genomics.org/plink/2.0/>

1018 GCTA: <https://yanglab.westlake.edu.cn/software/gcta/#Overview/>

1019 BayesR: <https://cnsgenomics.com/software/gctb/#Overview/>

1020 BOLT-REML: https://alkesgroup.broadinstitute.org/BOLT-LMM/BOLT-LMM_manual.html

1021 FUMA: <https://fuma.ctglab.nl/>

1022 SMR: <https://yanglab.westlake.edu.cn/software/smr/>

1023 GTEx version 8 (SMR format): <https://yanglab.westlake.edu.cn/software/smr/#DataResource/>

1024 Human Protein Atlas: <https://www.proteinatlas.org/>

1025 UCSC: <https://genome.ucsc.edu/>

1026

1027 **Acknowledgements**

1028 The Genotype-Tissue Expression (GTEx) Project was supported by the Common Fund of the Office of
1029 the Director of the National Institutes of Health, and by NCI, NHGRI, NHLBI, NIDA, NIMH, and NINDS.
1030 The *cis*-eQTL data of GTEx version 8 was summarized into SMR format.

1031 This research has been conducted using the data resource under the database of Genotypes and
1032 Phenotypes (dbGaP) study accession phs001992, and using the UK Biobank Resource under
1033 Application Number 12505. The authors thank GenomeDK and Aarhus University for providing
1034 computational resources and support that contributed to these research results.

1035 We thank Dr. Joana A. Revez for her insightful comments on BayesR and R tools.

1036 **Availability of data and materials**

1037 Owing to the sensitive nature of these data, individual level data can be accessed only through
1038 secure servers where download of individual level information is prohibited. Each scientific project
1039 must be approved before initiation, and approval is granted to a specific Danish research institution.
1040 International researchers may gain data access through collaboration with a Danish research
1041 institution. More information about getting access to the iPSYCH data can be obtained at
1042 <https://ipsych.au.dk/about-ipsych/>. The summary statistics from the GWAS for C3 and C4 and
1043 related analyses will be made available via the GWAS Catalogue <https://www.ebi.ac.uk/gwas/>
1044 (Accession numbers to follow).

1045

1046 **Funding**

1047 This study was supported by the Danish National Research Foundation, via a Niels Bohr
1048 Professorship to John McGrath. Bjarni Vilhjalmsson was also supported by a Lundbeck Foundation
1049 Fellowship (R335-2019-2339).

1050 This research was conducted using the Danish National Biobank resource, supported by the Novo
1051 Nordisk Foundation. The iPSYCH team was supported by grants from the Lundbeck Foundation
1052 (R102-A9118, R155-2014-1724, and R248-2017-2003), NIMH (1R01MH124851-01 to A.D.B.) and the
1053 Universities and University Hospitals of Aarhus and Copenhagen. High-performance computer
1054 capacity for handling and statistical analysis of iPSYCH data on the GenomeDK HPC facility was
1055 provided by the Center for Genomics and Personalized Medicine and the Centre for Integrative
1056 Sequencing, iSEQ, Aarhus University, Denmark (grant to ADB). The Anorexia Nervosa Genetics
1057 Initiative (ANGI) was an initiative of the Klarman Family Foundation. Genotyping of the Anorexia
1058 patients were funded by the Klarman Family Foundation.

1059 MEB was supported by the Independent Research Fund Denmark (grant number, 7025-00078B) and
1060 by an unrestricted grant from The Lundbeck Foundation (grant number, R268-2016-3925); JCD was
1061 supported by a grant from the Danish Council for Independent Research (grant number, 0134-
1062 00227B); AFM was supported by an ARC Future Fellowship (FT200100837); KLM was supported by
1063 grants from The Lundbeck Foundation and the Brain & Behavior Research Foundation; NRW was
1064 supported by NHMRC 1173790 and 1113400.

1065 CMB Is supported by NIMH (R56MH129437; R01MH120170; R01MH124871; R01MH119084;
1066 R01MH118278; R01 MH124871); Brain and Behavior Research Foundation Distinguished Investigator
1067 Grant; Swedish Research Council (Vetenskapsrådet, award: 538-2013-8864); Lundbeck Foundation
1068 (Grant no. R276-2018-4581).

1069 **Disclosures**

1070 CM Bulik reports: Shire (grant recipient, Scientific Advisory Board member); Lundbeckfonden (grant
1071 recipient); Pearson (author, royalty recipient); Equip Health Inc. (Clinical Advisory Board); Other
1072 authors have nothing to disclose.

1073

1074

1075 **References**

- 1076 1 Minton, K. Innate immunity: The inside story on complement activation. *Nat Rev Immunol*
1077 **14**, 61 (2014). <https://doi.org/10.1038/nri3603>
- 1078 2 Merle, N. S., Noe, R., Halbwachs-Mecarelli, L., Fremeaux-Bacchi, V. & Roumenina, L. T.
1079 Complement System Part II: Role in Immunity. *Front Immunol* **6**, 257 (2015).
1080 <https://doi.org/10.3389/fimmu.2015.00257>
- 1081 3 Merle, N. S., Church, S. E., Fremeaux-Bacchi, V. & Roumenina, L. T. Complement System Part
1082 I - Molecular Mechanisms of Activation and Regulation. *Front Immunol* **6**, 262 (2015).
1083 <https://doi.org/10.3389/fimmu.2015.00262>
- 1084 4 Reis, E. S., Mastellos, D. C., Hajishengallis, G. & Lambris, J. D. New insights into the immune
1085 functions of complement. *Nat Rev Immunol* **19**, 503-516 (2019).
1086 <https://doi.org/10.1038/s41577-019-0168-x>
- 1087 5 Magdalon, J. *et al.* Complement System in Brain Architecture and Neurodevelopmental
1088 Disorders. *Front. Neurosci.* **14**, 23 (2020). <https://doi.org/10.3389/fnins.2020.00023>
- 1089 6 Stephan, A. H., Barres, B. A. & Stevens, B. The complement system: an unexpected role in
1090 synaptic pruning during development and disease. *Annu. Rev. Neurosci.* **35**, 369-389 (2012).
1091 <https://doi.org/10.1146/annurev-neuro-061010-113810>
- 1092 7 Presumey, J., Bialas, A. R. & Carroll, M. C. Complement System in Neural Synapse Elimination
1093 in Development and Disease. *Adv. Immunol.* **135**, 53-79 (2017).
1094 <https://doi.org/10.1016/bs.ai.2017.06.004>
- 1095 8 Blanchong, C. A. *et al.* Genetic, structural and functional diversities of human complement
1096 components C4A and C4B and their mouse homologues, Slp and C4. *Int Immunopharmacol*
1097 **1**, 365-392 (2001). [https://doi.org/10.1016/s1567-5769\(01\)00019-4](https://doi.org/10.1016/s1567-5769(01)00019-4)
- 1098 9 Sekar, A. *et al.* Schizophrenia risk from complex variation of complement component 4.
1099 *Nature* **530**, 177-183 (2016). <https://doi.org/10.1038/nature16549>
- 1100 10 Schizophrenia Working Group of the Psychiatric Genomics, C. Biological insights from 108
1101 schizophrenia-associated genetic loci. *Nature* **511**, 421-427 (2014).
1102 <https://doi.org/10.1038/nature13595>
- 1103 11 Sellgren, C. M. *et al.* Increased synapse elimination by microglia in schizophrenia patient-
1104 derived models of synaptic pruning. *Nat Neurosci* **22**, 374-385 (2019).
1105 <https://doi.org/10.1038/s41593-018-0334-7>
- 1106 12 Yilmaz, M. *et al.* Overexpression of schizophrenia susceptibility factor human complement
1107 C4A promotes excessive synaptic loss and behavioral changes in mice. *Nat. Neurosci.* **24**,
1108 214-224 (2021). <https://doi.org/10.1038/s41593-020-00763-8>
- 1109 13 O'Connell, K. S. *et al.* Association between complement component 4A expression, cognitive
1110 performance and brain imaging measures in UK Biobank. *Psychol. Med.*, 1-11 (2021).
1111 <https://doi.org/10.1017/s0033291721000179>
- 1112 14 Kamitaki, N. *et al.* Complement genes contribute sex-biased vulnerability in diverse
1113 disorders. *Nature* **582**, 577-581 (2020). <https://doi.org/10.1038/s41586-020-2277-x>
- 1114 15 Li, N. *et al.* Association between C4, C4A, and C4B copy number variations and susceptibility
1115 to autoimmune diseases: a meta-analysis. *Sci. Rep.* **7**, 42628 (2017).
1116 <https://doi.org/10.1038/srep42628>
- 1117 16 Cross-Disorder Group of the Psychiatric Genomics, C. *et al.* Genetic relationship between five
1118 psychiatric disorders estimated from genome-wide SNPs. *Nat. Genet.* **45**, 984-994 (2013).
1119 <https://doi.org/10.1038/ng.2711>
- 1120 17 Kim, M. *et al.* Brain gene co-expression networks link complement signaling with convergent
1121 synaptic pathology in schizophrenia. *Nat. Neurosci.* **24**, 799-809 (2021).
1122 <https://doi.org/10.1038/s41593-021-00847-z>

- 1123 18 Yang, X. *et al.* Genome-wide association study for serum complement C3 and C4 levels in
1124 healthy Chinese subjects. *PLoS Genet.* **8**, e1002916 (2012).
1125 <https://doi.org/10.1371/journal.pgen.1002916>
- 1126 19 Yang, Y. *et al.* Diversity in intrinsic strengths of the human complement system: serum C4
1127 protein concentrations correlate with C4 gene size and polygenic variations, hemolytic
1128 activities, and body mass index. *J. Immunol.* **171**, 2734-2745 (2003).
1129 <https://doi.org/10.4049/jimmunol.171.5.2734>
- 1130 20 Hebert, L. A., Cosio, F. G. & Neff, J. C. Diagnostic significance of hypocomplementemia.
1131 *Kidney Int.* **39**, 811-821 (1991). <https://doi.org/10.1038/ki.1991.102>
- 1132 21 Mongan, D. *et al.* Peripheral complement proteins in schizophrenia: A systematic review and
1133 meta-analysis of serological studies. *Schizophr. Res.* **222**, 58-72 (2020).
1134 <https://doi.org/10.1016/j.schres.2020.05.036>
- 1135 22 Murray, R. M. & Lewis, S. W. Is schizophrenia a neurodevelopmental disorder? *Br. Med. J.*
1136 *(Clin. Res. Ed)* **295**, 681-682 (1987).
- 1137 23 Weinberger, D. R. Implications of normal brain development for the pathogenesis of
1138 schizophrenia. *Arch. Gen. Psychiatry* **44**, 660-669 (1987).
- 1139 24 Schafer, D. P. *et al.* Microglia sculpt postnatal neural circuits in an activity and complement-
1140 dependent manner. *Neuron* **74**, 691-705 (2012).
1141 <https://doi.org/10.1016/j.neuron.2012.03.026>
- 1142 25 Cooper, J. D. *et al.* Schizophrenia-risk and urban birth are associated with proteomic changes
1143 in neonatal dried blood spots. *Transl Psychiatry* **7**, 1290 (2017).
1144 <https://doi.org/10.1038/s41398-017-0027-0>
- 1145 26 Pedersen, C. B. *et al.* The iPSYCH2012 case-cohort sample: new directions for unravelling
1146 genetic and environmental architectures of severe mental disorders. *Mol. Psychiatry* **23**, 6-
1147 14 (2018). <https://doi.org/10.1038/mp.2017.196>
- 1148 27 Bybjerg-Grauholm, J. *et al.* The iPSYCH2015 Case-Cohort sample: updated directions for
1149 unravelling genetic and environmental architectures of severe mental disorders. *medRxiv*,
1150 2020.2011.2030.20237768 (2020). <https://doi.org/10.1101/2020.11.30.20237768>
- 1151 28 Pendergrass, S. A. *et al.* The use of phenome-wide association studies (PheWAS) for
1152 exploration of novel genotype-phenotype relationships and pleiotropy discovery. *Genet.*
1153 *Epidemiol.* **35**, 410-422 (2011). <https://doi.org/10.1002/gepi.20589>
- 1154 29 Thornton, L. M. *et al.* The Anorexia Nervosa Genetics Initiative (ANGI): Overview and
1155 methods. *Contemp. Clin. Trials* **74**, 61-69 (2018). <https://doi.org/10.1016/j.cct.2018.09.015>
- 1156 30 Munk-Jorgensen, P. & Mortensen, P. B. The Danish Psychiatric Central Register. *Dan. Med.*
1157 *Bull.* **44**, 82-84 (1997).
- 1158 31 Mors, O., Perto, G. P. & Mortensen, P. B. The Danish Psychiatric Central Research Register.
1159 *Scand J Public Health* **39**, 54-57 (2011). <https://doi.org/10.1177/1403494810395825>
- 1160 32 Borgan, O., Langholz, B., Samuelsen, S. O., Goldstein, L. & Pogoda, J. Exposure stratified case-
1161 cohort designs. *Lifetime Data Anal* **6**, 39-58 (2000).
1162 <https://doi.org/10.1023/a:1009661900674>
- 1163 33 Norgaard-Pedersen, B. & Hougaard, D. M. Storage policies and use of the Danish Newborn
1164 Screening Biobank. *J. Inherit. Metab. Dis.* **30**, 530-536 (2007).
1165 <https://doi.org/10.1007/s10545-007-0631-x>
- 1166 34 Hollegaard, M. V. *et al.* Whole genome amplification and genetic analysis after extraction of
1167 proteins from dried blood spots. *Clin. Chem.* **53**, 1161-1162 (2007).
1168 <https://doi.org/10.1373/clinchem.2006.082313>
- 1169 35 Albiñana, C. *et al.* Genetic correlates of vitamin D-binding protein and 25 hydroxyvitamin D
1170 in neonatal dried blood spots. *medRxiv*, 2022.2006.2008.22276164 (2022).
1171 <https://doi.org/10.1101/2022.06.08.22276164>
- 1172 36 Gunderson, K. L. *et al.* Whole-genome genotyping. *Methods Enzymol.* **410**, 359-376 (2006).
1173 [https://doi.org/10.1016/S0076-6879\(06\)10017-8](https://doi.org/10.1016/S0076-6879(06)10017-8)

- 1174 37 Lam, M. *et al.* RICOPIIL: Rapid Imputation for COnsortias PipeLine. *Bioinformatics* **36**, 930-
1175 933 (2020). <https://doi.org/10.1093/bioinformatics/btz633>
- 1176 38 McCarthy, S. *et al.* A reference panel of 64,976 haplotypes for genotype imputation. *Nat.*
1177 *Genet.* **48**, 1279-1283 (2016). <https://doi.org/10.1038/ng.3643>
- 1178 39 Browning, B. L., Zhou, Y. & Browning, S. R. A One-Penny Imputed Genome from Next-
1179 Generation Reference Panels. *Am. J. Hum. Genet.* **103**, 338-348 (2018).
1180 <https://doi.org/10.1016/j.ajhg.2018.07.015>
- 1181 40 Wouters, D. *et al.* High-throughput analysis of the C4 polymorphism by a combination of
1182 MLPA and isotype-specific ELISA's. *Mol. Immunol.* **46**, 592-600 (2009).
1183 <https://doi.org/10.1016/j.molimm.2008.07.028>
- 1184 41 Bian, B., Couvy-Duchesne, B., Wray, N. R. & McRae, A. F. The role of critical immune genes in
1185 brain disorders: insights from neuroimaging immunogenetics. *Brain Commun* **4**, fcac078
1186 (2022). <https://doi.org/10.1093/braincomms/fcac078>
- 1187 42 Jiang, J. & Nguyen, T. in *Linear and Generalized Linear Mixed Models and Their Applications*
1188 1-61 (Springer New York, 2021).
- 1189 43 Beasley, T. M., Erickson, S. & Allison, D. B. Rank-based inverse normal transformations are
1190 increasingly used, but are they merited? *Behav Genet* **39**, 580-595 (2009).
1191 <https://doi.org/10.1007/s10519-009-9281-0>
- 1192 44 Zaitlen, N. *et al.* Using extended genealogy to estimate components of heritability for 23
1193 quantitative and dichotomous traits. *PLoS Genet.* **9**, e1003520 (2013).
1194 <https://doi.org/10.1371/journal.pgen.1003520>
- 1195 45 Yang, C. *et al.* Genomic atlas of the proteome from brain, CSF and plasma prioritizes proteins
1196 implicated in neurological disorders. *Nat Neurosci* **24**, 1302-1312 (2021).
1197 <https://doi.org/10.1038/s41593-021-00886-6>
- 1198 46 Pietzner, M. *et al.* Mapping the proteo-genomic convergence of human diseases. *Science*
1199 **374**, eabj1541 (2021). <https://doi.org/10.1126/science.abj1541>
- 1200 47 Gudjonsson, A. *et al.* A genome-wide association study of serum proteins reveals shared loci
1201 with common diseases. *Nat Commun* **13**, 480 (2022). [https://doi.org/10.1038/s41467-021-](https://doi.org/10.1038/s41467-021-27850-z)
1202 [27850-z](https://doi.org/10.1038/s41467-021-27850-z)
- 1203 48 Zeng, J. *et al.* Signatures of negative selection in the genetic architecture of human complex
1204 traits. *Nat. Genet.* **50**, 746-753 (2018). <https://doi.org/10.1038/s41588-018-0101-4>
- 1205 49 Yang, J. *et al.* Common SNPs explain a large proportion of the heritability for human height.
1206 *Nat. Genet.* **42**, 565-569 (2010). <https://doi.org/ng.608> [pii]
- 1207 10.1038/ng.608
- 1208 50 Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide complex
1209 trait analysis. *Am. J. Hum. Genet.* **88**, 76-82 (2011).
1210 <https://doi.org/10.1016/j.ajhg.2010.11.011>
- 1211 51 Loh, P. R. *et al.* Contrasting genetic architectures of schizophrenia and other complex
1212 diseases using fast variance-components analysis. *Nat. Genet.* **47**, 1385-1392 (2015).
1213 <https://doi.org/10.1038/ng.3431>
- 1214 52 Elston, R. C., Buxbaum, S., Jacobs, K. B. & Olson, J. M. Haseman and Elston revisited. *Genet*
1215 *Epidemiol* **19**, 1-17 (2000). [https://doi.org/10.1002/1098-2272\(200007\)19:1<1::AID-](https://doi.org/10.1002/1098-2272(200007)19:1<1::AID-GEPI1>3.0.CO;2-E)
1216 [GEPI1>3.0.CO;2-E](https://doi.org/10.1002/1098-2272(200007)19:1<1::AID-GEPI1>3.0.CO;2-E)
- 1217 53 Jiang, L. *et al.* A resource-efficient tool for mixed model association analysis of large-scale
1218 data. *Nat. Genet.* **51**, 1749-1755 (2019). <https://doi.org/10.1038/s41588-019-0530-8>
- 1219 54 Yang, J., Zaitlen, N. A., Goddard, M. E., Visscher, P. M. & Price, A. L. Advantages and pitfalls in
1220 the application of mixed-model association methods. *Nat Genet* **46**, 100-106 (2014).
1221 <https://doi.org/10.1038/ng.2876>
- 1222 55 Galinsky, K. J. *et al.* Fast Principal-Component Analysis Reveals Convergent Evolution of
1223 ADH1B in Europe and East Asia. *Am. J. Hum. Genet.* **98**, 456-472 (2016).
1224 <https://doi.org/10.1016/j.ajhg.2015.12.022>

- 1225 56 Sidorenko, J. *et al.* The effect of X-linked dosage compensation on complex trait variation.
1226 *Nat Commun* **10**, 3009 (2019). <https://doi.org/10.1038/s41467-019-10598-y>
- 1227 57 Woo, J. J., Pouget, J. G., Zai, C. C. & Kennedy, J. L. The complement system in schizophrenia:
1228 where are we now and what's next? *Mol. Psychiatry* **25**, 114-130 (2020).
1229 <https://doi.org/10.1038/s41380-019-0479-0>
- 1230 58 Watanabe, K., Taskesen, E., van Bochoven, A. & Posthuma, D. Functional mapping and
1231 annotation of genetic associations with FUMA. *Nat Commun* **8**, 1826 (2017).
1232 <https://doi.org/10.1038/s41467-017-01261-5>
- 1233 59 Zhu, Z. *et al.* Integration of summary data from GWAS and eQTL studies predicts complex
1234 trait gene targets. *Nat. Genet.* **48**, 481-487 (2016). <https://doi.org/10.1038/ng.3538>
- 1235 60 GTEx Consortium. Genetic effects on gene expression across human tissues. *Nature* **550**, 204
1236 (2017). <https://doi.org/10.1038/nature24277>
- 1237 <https://www.nature.com/articles/nature24277#supplementary-information>
- 1238 61 Zhu, Z. *et al.* Causal associations between risk factors and common diseases inferred from
1239 GWAS summary data. *Nat Commun* **9**, 224 (2018). [https://doi.org/10.1038/s41467-017-](https://doi.org/10.1038/s41467-017-02317-2)
1240 [02317-2](https://doi.org/10.1038/s41467-017-02317-2)
- 1241 62 Trubetsky, V. *et al.* Mapping genomic loci implicates genes and synaptic biology in
1242 schizophrenia. *Nature* **604**, 502-508 (2022). <https://doi.org/10.1038/s41586-022-04434-5>
- 1243 63 Howard, D. M. *et al.* Genome-wide meta-analysis of depression identifies 102 independent
1244 variants and highlights the importance of the prefrontal brain regions. *Nat. Neurosci.* **22**,
1245 343-352 (2019). <https://doi.org/10.1038/s41593-018-0326-7>
- 1246 64 Mullins, N. *et al.* Genome-wide association study of more than 40,000 bipolar disorder cases
1247 provides new insights into the underlying biology. *Nat. Genet.* **53**, 817-829 (2021).
1248 <https://doi.org/10.1038/s41588-021-00857-4>
- 1249 65 Grove, J. *et al.* Identification of common genetic risk variants for autism spectrum disorder.
1250 *Nat. Genet.* **51**, 431-444 (2019). <https://doi.org/10.1038/s41588-019-0344-8>
- 1251 66 Demontis, D. *et al.* Discovery of the first genome-wide significant risk loci for attention
1252 deficit/hyperactivity disorder. *Nat. Genet.* **51**, 63-75 (2019). [https://doi.org/10.1038/s41588-](https://doi.org/10.1038/s41588-018-0269-7)
1253 [018-0269-7](https://doi.org/10.1038/s41588-018-0269-7)
- 1254 67 Watson, H. J. *et al.* Genome-wide association study identifies eight risk loci and implicates
1255 metabo-psychiatric origins for anorexia nervosa. *Nat Genet* **51**, 1207-1214 (2019).
1256 <https://doi.org/10.1038/s41588-019-0439-2>
- 1257 68 Marioni, R. E. *et al.* GWAS on family history of Alzheimer's disease. *Transl Psychiatry* **8**, 99
1258 (2018). <https://doi.org/10.1038/s41398-018-0150-6>
- 1259 69 van Rheenen, W. *et al.* Common and rare variant association analyses in amyotrophic lateral
1260 sclerosis identify 15 risk loci with distinct genetic architectures and neuron-specific biology.
1261 *Nat. Genet.* **53**, 1636-1648 (2021). <https://doi.org/10.1038/s41588-021-00973-1>
- 1262 70 International Multiple Sclerosis Genetics, C. Multiple sclerosis genomic map implicates
1263 peripheral immune cells and microglia in susceptibility. *Science* **365** (2019).
1264 <https://doi.org/10.1126/science.aav7188>
- 1265 71 Chiou, J. *et al.* Interpreting type 1 diabetes risk with genetics and single-cell epigenomics.
1266 *Nature* **594**, 398-402 (2021). <https://doi.org/10.1038/s41586-021-03552-w>
- 1267 72 de Lange, K. M. *et al.* Genome-wide association study implicates immune activation of
1268 multiple integrin genes in inflammatory bowel disease. *Nat. Genet.* **49**, 256-261 (2017).
1269 <https://doi.org/10.1038/ng.3760>
- 1270 73 Okada, Y. *et al.* Genetics of rheumatoid arthritis contributes to biology and drug discovery.
1271 *Nature* **506**, 376-381 (2014). <https://doi.org/10.1038/nature12873>
- 1272 74 Julia, A. *et al.* Genome-wide association study meta-analysis identifies five new loci for
1273 systemic lupus erythematosus. *Arthritis Res. Ther.* **20**, 100 (2018).
1274 <https://doi.org/10.1186/s13075-018-1604-1>

- 1275 75 Byrne, E. M. *et al.* Conditional GWAS analysis to identify disorder-specific SNPs for
1276 psychiatric disorders. *Mol Psychiatry* **26**, 2070-2081 (2021). [https://doi.org/10.1038/s41380-](https://doi.org/10.1038/s41380-020-0705-9)
1277 [020-0705-9](https://doi.org/10.1038/s41380-020-0705-9)
- 1278 76 Sudlow, C. *et al.* UK biobank: an open access resource for identifying the causes of a wide
1279 range of complex diseases of middle and old age. *PLoS Med.* **12**, e1001779 (2015).
1280 <https://doi.org/10.1371/journal.pmed.1001779>
- 1281 77 Consortium, U. K. *et al.* The UK10K project identifies rare variants in health and disease.
1282 *Nature* **526**, 82-90 (2015). <https://doi.org/10.1038/nature14962>
- 1283 78 Revez, J. A. *et al.* Genome-wide association study identifies 143 loci associated with 25
1284 hydroxyvitamin D concentration. *Nat Commun* **11**, 1647 (2020).
1285 <https://doi.org/10.1038/s41467-020-15421-7>
- 1286 79 Yang, J. *et al.* Conditional and joint multiple-SNP analysis of GWAS summary statistics
1287 identifies additional variants influencing complex traits. *Nat. Genet.* **44**, 369-375, S361-363
1288 (2012). <https://doi.org/10.1038/ng.2213>
- 1289 80 Ferkingstad, E. *et al.* Large-scale integration of the plasma proteome with genetics and
1290 disease. *Nat. Genet.* **53**, 1712-1721 (2021). <https://doi.org/10.1038/s41588-021-00978-w>
- 1291 81 Karlsson, M. *et al.* A single-cell type transcriptomics map of human tissues. *Sci Adv* **7** (2021).
1292 <https://doi.org/10.1126/sciadv.abh2169>
- 1293 82 Dunkelberger, J. R. & Song, W. C. Complement and its role in innate and adaptive immune
1294 responses. *Cell Res.* **20**, 34-50 (2010). <https://doi.org/10.1038/cr.2009.139>
- 1295 83 Schurz, H. *et al.* The X chromosome and sex-specific effects in infectious disease
1296 susceptibility. *Hum Genomics* **13**, 2 (2019). <https://doi.org/10.1186/s40246-018-0185-z>
- 1297 84 Poppelaars, F. *et al.* A Family Affair: Addressing the Challenges of Factor H and the Related
1298 Proteins. *Front. Immunol.* **12**, 660194 (2021). <https://doi.org/10.3389/fimmu.2021.660194>
- 1299 85 Kent, W. J. *et al.* The human genome browser at UCSC. *Genome Res* **12**, 996-1006 (2002).
1300 <https://doi.org/10.1101/gr.229102>
- 1301 86 Lee, B. T. *et al.* The UCSC Genome Browser database: 2022 update. *Nucleic Acids Res* **50**,
1302 D1115-D1122 (2022). <https://doi.org/10.1093/nar/gkab959>
- 1303 87 Jain, U., Otley, A. R., Van Limbergen, J. & Stadnyk, A. W. The complement system in
1304 inflammatory bowel disease. *Inflamm Bowel Dis* **20**, 1628-1637 (2014).
1305 <https://doi.org/10.1097/MIB.000000000000056>
- 1306 88 Cleyne, I. *et al.* Genome-Wide Copy Number Variation Scan Identifies Complement
1307 Component C4 as Novel Susceptibility Gene for Crohn's Disease. *Inflamm. Bowel Dis.* **22**,
1308 505-515 (2016). <https://doi.org/10.1097/MIB.0000000000000623>
- 1309 89 Cruz-Tapias, P., Rojas-Villarraga, A., Maier-Moore, S. & Anaya, J. M. HLA and Sjogren's
1310 syndrome susceptibility. A meta-analysis of worldwide studies. *Autoimmun Rev* **11**, 281-287
1311 (2012). <https://doi.org/10.1016/j.autrev.2011.10.002>
- 1312 90 Hanscombe, K. B. *et al.* Genetic fine mapping of systemic lupus erythematosus MHC
1313 associations in Europeans and African Americans. *Hum Mol Genet* **27**, 3813-3824 (2018).
1314 <https://doi.org/10.1093/hmg/ddy280>
- 1315 91 International, M. H. C. *et al.* Mapping of multiple susceptibility variants within the MHC
1316 region for 7 immune-mediated diseases. *Proc Natl Acad Sci U S A* **106**, 18680-18685 (2009).
1317 <https://doi.org/10.1073/pnas.0909307106>
- 1318 92 Langefeld, C. D. *et al.* Transancestral mapping and genetic load in systemic lupus
1319 erythematosus. *Nat Commun* **8**, 16021 (2017). <https://doi.org/10.1038/ncomms16021>
- 1320 93 Li, Y. R. *et al.* Meta-analysis of shared genetic architecture across ten pediatric autoimmune
1321 diseases. *Nat Med* **21**, 1018-1027 (2015). <https://doi.org/10.1038/nm.3933>
- 1322 94 Sakaue, S. *et al.* A cross-population atlas of genetic associations for 220 human phenotypes.
1323 *Nat. Genet.* **53**, 1415-1424 (2021). <https://doi.org/10.1038/s41588-021-00931-x>
- 1324 95 Wu, Y. L. *et al.* Phenotypes, genotypes and disease susceptibility associated with gene copy
1325 number variations: complement C4 CNVs in European American healthy subjects and those

- 1326 with systemic lupus erythematosus. *Cytogenet Genome Res* **123**, 131-141 (2008).
1327 [https://doi.org:10.1159/000184700](https://doi.org/10.1159/000184700)
1328 96 Tsang, A. S. M. W. P. *et al.* Comprehensive approach to study complement C4 in systemic
1329 lupus erythematosus: Gene polymorphisms, protein levels and functional activity. *Mol*
1330 *Immunol* **92**, 125-131 (2017). [https://doi.org:10.1016/j.molimm.2017.10.004](https://doi.org/10.1016/j.molimm.2017.10.004)
1331 97 Hoffman, G. E. *et al.* Comment on: What genes are differentially expressed in individuals
1332 with schizophrenia? A systematic review. *Mol Psychiatry* (2022).
1333 [https://doi.org:10.1038/s41380-022-01781-7](https://doi.org/10.1038/s41380-022-01781-7)
1334