



An expanded set of genome-wide association studies of brain imaging phenotypes in UK Biobank

Stephen M. Smith¹, Gwenaëlle Douaud¹, Winfield Chen², Taylor Hanayik¹, Fidel Alfaro-Almagro¹, Kevin Sharp³ and Lloyd T. Elliott²✉

UK Biobank is a major prospective epidemiological study, including multimodal brain imaging, genetics and ongoing health outcomes. Previously, we published genome-wide associations of 3,144 brain imaging-derived phenotypes, with a discovery sample of 8,428 individuals. Here we present a new open resource of genome-wide association study summary statistics, using the 2020 data release, almost tripling the discovery sample size. We now include the X chromosome and new classes of imaging-derived phenotypes (subcortical volumes and tissue contrast). Previously, we found 148 replicated clusters of associations between genetic variants and imaging phenotypes; in this study, we found 692, including 12 on the X chromosome. We describe some of the newly found associations, focusing on the X chromosome and autosomal associations involving the new classes of imaging-derived phenotypes. Our novel associations implicate, for example, pathways involved in the rare X-linked STAR (syndactyly, telecanthus and anogenital and renal malformations) syndrome, Alzheimer's disease and mitochondrial disorders.

UK Biobank (UKB) is now approximately halfway through imaging 100,000 volunteers; the early-2020 release of brain imaging data contained data from almost 40,000 participants. This spans six brain magnetic resonance imaging (MRI) modalities, allowing the study of many different aspects of brain structure, function and connectivity. In conjunction with other data being recorded by UKB, which include health outcomes, lifestyle, biophysical measures and genetics, UKB is a major resource for understanding the brain in human health and disease.

Previously, we presented genome-wide association studies (GWASs) of 3,144 brain imaging-derived phenotypes (IDPs), with a discovery sample of 8,428 individuals¹. At that point, we identified 148 replicated clusters of associations between genetic variants and the phenotypes. We found links between IDPs and genes involved in iron transport and storage, extracellular matrix and epidermal growth factor, development, pathway signaling and plasticity.

We have now expanded and enhanced this work, with an almost three-fold increase in sample size, an increase in the number of IDPs to almost 4,000 and with a focus on X chromosome associations carried out for the first time (including around 10 million variants with minor allele frequency (MAF) $\geq 1\%$). The new classes of IDPs, computed on behalf of UKB and released for general access, are as follows: subnuclei volumes in amygdala, brainstem, hippocampus and thalamus, Brodmann area FreeSurfer metrics and FreeSurfer-derived white-gray intensity contrasts. We have also greatly expanded our set of imaging confound variables², reducing the likelihood of finding artifactual associations. GWAS summary statistics and Manhattan plots for all 3,935 phenotypes are freely available for download from the Oxford Brain Imaging Genetics (BIG40) web server (available at <https://open.win.ox.ac.uk/ukbiobank/big40/>), which also includes detailed tables of all IDPs, all single-nucleotide polymorphisms (SNPs) tested, all

association clusters and an interactive viewer allowing for detailed interrogation of IDP associations with SNPs and nearby genes. We also provide a list of causal genetic variants for our top X chromosome clusters using a statistical fine mapping approach³.

We conducted sex-specific GWASs on autosomes (chromosomes 1–22) and the X chromosome, followed by a meta-analysis combining these, using Fisher's method. The X chromosome accounts for about 5% of the human genome and incorporates over 1,200 genes, including many that play a role in human cognition and development. However, testing for association with genetic variants on chromosome X requires special consideration⁴. Although genetic females inherit two copies of the X chromosome, genetic males inherit only a single copy, from the maternal line (here we refer to people with two X chromosomes as genetic females and people with one X and one Y chromosome as genetic males). The short pseudo-autosomal regions (PARs) on the ends of chromosome X are homologous with parts of chromosome Y and can be analyzed in the same way as autosomal chromosomes. For the non-PAR region, a mechanism has evolved to balance allele dosage differences between the genetic sexes (X chromosome inactivation (XCI)); during female development, one copy is randomly inactivated in each cell. This means that, maternally and paternally, inherited alleles would be expected to be expressed in different cell populations within the body approximately 50% of the time. However, this dosage compensation (DC) mechanism is imperfect; it is currently thought that only 60–75% of X-linked genes have one copy completely silenced in this way⁵.

To account for this in GWASs, it is common to assume full DC⁴: males are treated as homozygous females with genotypes coded as (0,2) according to whether they have 0 or 1 copy of the alternative allele. Joint analysis of genetic females and males with simulation studies⁶ suggest that type 1 error control under this approach is

¹Wellcome Centre for Integrative Neuroimaging (WIN FMRIB), University of Oxford, Oxford, United Kingdom. ²Department of Statistics and Actuarial Science, Simon Fraser University, Burnaby BC, Canada. ³Genomics PLC, Oxford, United Kingdom. ✉e-mail: lloyd_elliott@sfu.ca

reasonably robust to deviations from other assumptions (such as no sex-specific differences in allele frequencies), provided genetic sex is included as a covariate.

Recent studies have used the large sample sizes afforded by UKB to perform stratified analyses to estimate the degree of DC as a parameter across a broad variety of traits⁵. These studies suggest that only a small proportion of genes escape XCI, although the appropriate amount of DC shows considerable variation among traits. For educational attainment, Lee et al.⁵ estimated a DC factor of 1.45, but concluded that little power was lost in their joint analysis irrespective of the assumed model.

Although a joint analysis under a full DC model is a reasonable default, the power afforded by biobank-scale datasets also permits examination of possible sex-specific effects via stratified analyses. These stratified analyses can subsequently be meta-analyzed; if the meta-analysis is based on estimated effects (regression beta values), appropriately chosen weights can give results almost the same as those from a joint genetic male/female analysis corresponding to any assumed DC model⁵. Nevertheless, results will be biased if the assumed DC model differs from the truth. An unweighted meta-analysis based on *P* values using Fisher's method (explored in our research), although potentially less powerful, should avoid this possible bias (as it is not sensitive to any relative scaling in the regression model and, hence, the effect sizes in the two sex-separated GWASs) and still have value in confirming signals from a joint analysis.

The BIG40 web server includes a summary statistics resource including results for our discovery cohort and a full GWAS on all samples passing quality control (with discovery and replication cohorts combined). A browsable interface is also provided for both of these subsets. In addition, the results of the sex-separated GWAS on chromosome X and autosomes are provided. We have also provided the summary statistics to the European Bioinformatics Institute GWAS repository. More details about our resource are provided in the Online Methods.

Results

Overview of GWAS results. We conducted a GWAS using the 39,691 brain imaged samples in UKB. We divided these samples into a discovery cohort ($n=22,138$) and a replication cohort ($n=11,086$). The details for the imaging and genetics processing and the cohorts are given in the Online Methods. We applied automated methods for identifying local peak associations for each phenotype and also for aggregating peaks across phenotypes into clusters of association. A cluster is a set of phenotype/variant pairs such that all of the phenotype/variant pairs have a $-\log_{10}(P)$ value for association that exceeds a 7.5 genome-wide significance threshold, and such that all of the pairs have variants that are close with respect to genetic distance. We assigned each of the phenotype/variant pairs to one and only one cluster. We defined a cluster as replicating if at least one of the phenotype/variant pairs had nominal significance in our replication cohort ($P < 0.05$).

With these methods, we found 10,889 peak associations among all phenotypes and chromosomes (8,446 replicating at nominal significance) and found 1,282 clusters (692 replicating) after clustering the peak associations according to our automated methods (the number of replicating clusters reported in our previous work¹ was 148). The 692 replicating clusters are distributed across all chromosomes, with between 8 and 60 clusters per chromosome. We grouped the IDPs into 17 categories (Supplementary Table 1). Of the replicated associations among these 692 clusters, 16 of 17 categories are represented (the task functional MRI activation category is the only category without at least one association). The number of associations per category ranges between 12 for the category volume of white matter hyperintensities (lesions), which consists of just one IDP, and 1,954 for the regional and tissue volume

category. All of these associations are listed in Supplementary Table 4, and Manhattan plots along with quantile plots are provided on the BIG40 open web server.

Of all of our clusters of associations, 38 are on the X chromosome (12 replicating), and four of the X chromosome clusters have a phenotype/variant pair with association significance exceeding the more stringent Bonferroni-corrected level of $-\log_{10}(P) \geq 11.1$. (In this work, we adjust for computing multiple GWASs by applying the Bonferroni correction on top of the standard GWAS threshold of $-\log_{10}(P) \geq 7.5$, resulting in a 'Bonferroni' threshold of 11.1.) These four clusters are investigated below, and Manhattan plots for the IDPs most associated with these clusters are displayed in Fig. 1. We also investigate five novel clusters among the autosomal chromosomes.

We provide a fine mapping of the four X chromosome clusters using the CAVIAR software³ with results described in Supplementary Table 3, considering regions within 250 kbp of the lead associations for the clusters. For Cluster 1, this resulted in a region containing 1,211 genetic variants. We ran the CAVIAR software³ with default settings and recorded the genetic variant found to be most causal for each phenotype association within the cluster. For Cluster 1, CAVIAR reported rs2272737 (the same lead genetic variant found by our aggregation method) as the causal genetic variant for 81/97 of the phenotype associations. The situation was similar for the other three clusters; in each case, the lead genetic variant found by our aggregation method was among the genetic variants found to be causal by CAVIAR. (The number of genetic variants included in the 500-kb regions examined by CAVIAR in the remaining three clusters included 2,017, 731 and 935 genetic variants, respectively. The proportion of phenotype associations for which the lead genetic variant found by our aggregation method was also estimated to be the most probable causal variant by CAVIAR for these remaining three clusters was 10/21, 5/35 and 5/17, respectively.) The CAVIAR results are detailed in Supplementary Table 3, and further details provided in Methods.

BIG40 also provides summary statistics for a GWAS with the discovery and replication cohorts combined (a full scan). We also examined the heritability of each phenotype using linkage score regression⁷. The heritability for each of the phenotype categories is summarized in Fig. 2. With phenotypes for which the estimated h^2 value was more than one standard error greater than 0, the estimated value of h^2 ranged from 0.01 to 0.41. The highest estimated heritability was found in phenotypes involved in regional and tissue volumes (as in previous work¹), cortical gray-white contrast and the intra-cellular volume fraction diffusion imaging measure.

X chromosome results—overview and sex-specific tests. Full details for the lead associations for the X chromosome clusters (including clusters that do not replicate) are provided in Supplementary Table 2. A summary of all of the peak associations included in the X chromosome clusters is provided in Supplementary Table 3, and the full results for peak associations on all chromosomes are provided in Supplementary Table 4. In these tables, cluster numbers are given in the first column, and clusters are ordered based on the chromosome number (in ascending order, with the X chromosome first) and then by the $-\log_{10}(P)$ value of the lead association (in descending order). A summary of all replicating X chromosome clusters is provided in Table 1, and further details, including genes and expression quantitative trait loci (eQTL) for the four Bonferroni-significant clusters, are provided in Table 2. Figure 1 shows Manhattan plots for the lead associations in these four top X chromosome clusters, and these clusters are explored further below.

Genetic sex affects the brain in fundamental ways⁸. Our main GWAS analyses include sex as one of the confound variables (in part, as sex is a causal factor in some strong imaging confound effects, such as interaction of head size with image intensity and

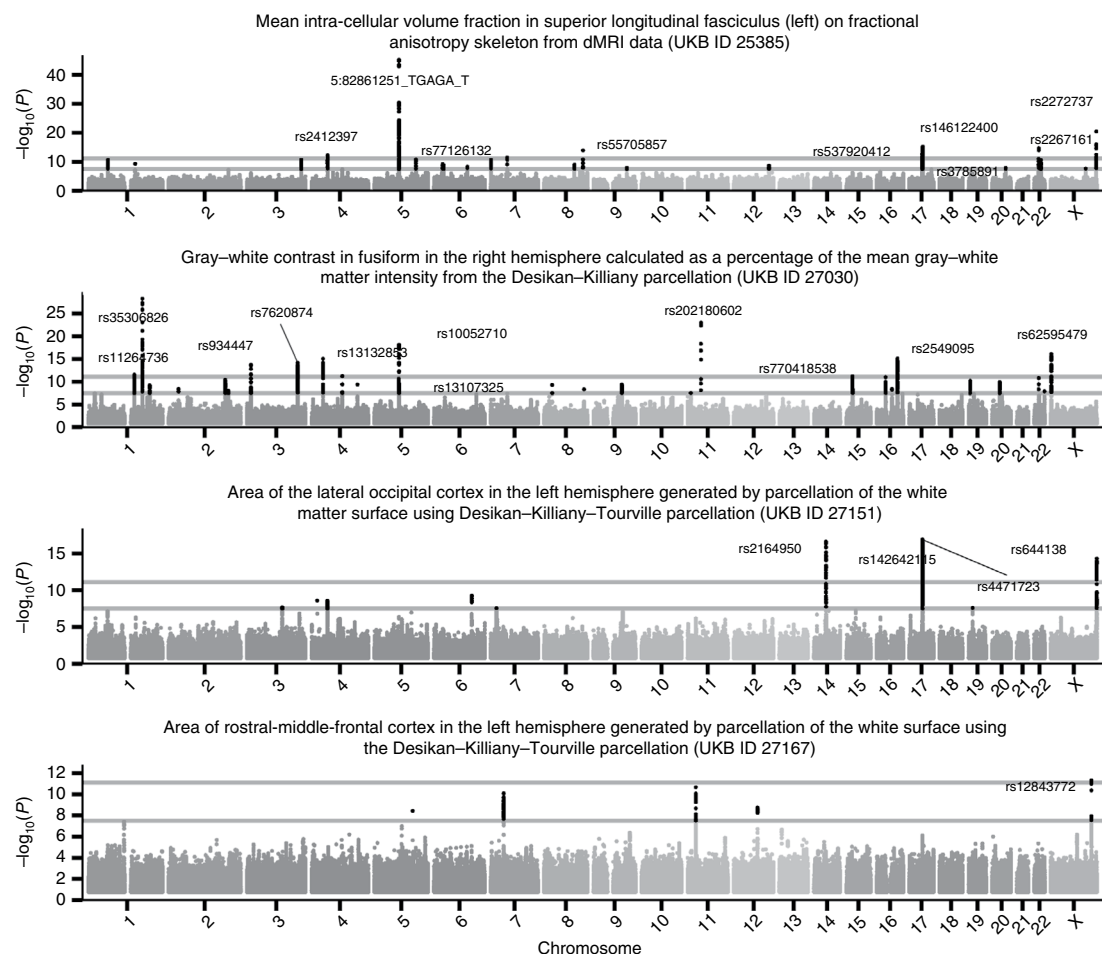


Fig. 1 | Manhattan plots for the four phenotypes achieving Bonferroni-corrected significance on the X chromosome. Genetic variants are labeled for peak associations achieving the Bonferroni level. Plot titles indicate phenotype definition (including the UKB ID field index from <http://biobank.ctsu.ox.ac.uk/crystal/field.cgi?id=25385>, <http://biobank.ctsu.ox.ac.uk/crystal/field.cgi?id=27030>, <http://biobank.ctsu.ox.ac.uk/crystal/field.cgi?id=27151> or <http://biobank.ctsu.ox.ac.uk/crystal/field.cgi?id=27167>). Black dots indicate associations that are significant at the genome-wide level, $-\log_{10}(P) \geq 7.5$. Gray lines show genome-wide + Bonferroni level (11.1) and genome-wide significance level (7.5). These associations involve diffusion MRI (dMRI) and the Desikan-Killiany and Desikan-Killiany-Tourville parcellations of white matter and gray matter.

head motion). To assess the quality of this deconfounding, and to explore associations on the X chromosome that are driven by genetic sex, we conducted two additional GWASs in which we restricted our discovery cohort to just genetic females and (separately) just genetic males. We then combined these two additional GWASs in a meta-analysis using Fisher's method. Clusters in our main analysis that are significant in the meta-analysis but not significant in one of the sex-specific scans might indicate sex-driven associations. Of the 12 replicating X chromosome clusters in the main analysis (with genetic male and female sexes combined in the discovery cohort), one cluster (Cluster 2) is significant at the level for one genetic sex but not significant for the other genetic sex; Cluster 2 might, therefore, be driven by genetic males. To provide more direct evidence for this, we performed two-tailed *z*-tests to determine if the beta coefficients differed significantly between the genetic sexes. For Clusters 1 and 2, we found that the beta coefficients are nominally different (Cluster 1: $P = 1.2 \times 10^{-2}$, beta coefficient for genetic females: -0.14 , beta coefficient for genetic males: -0.08 ; Cluster 2: $P = 1.0 \times 10^{-2}$, beta coefficient for genetic females: 0.07 , beta coefficient for genetic males: 0.13). The differences between the sex-specific beta coefficients for the lead associations of Clusters 3 and 4 were not significant. Among all of our associations (from

all chromosomes and all IDPs) with at least one of the sex-specific scans significant at the $-\log_{10}(P) \geq 11.1$ level, the signs of the effect sizes for genetic females and genetic males always matched. For significance at $-\log_{10}(P) \geq 7.5$, the signs matched 99.42% of the time (Extended Data Fig. 1).

Finally, we created an additional set of clusters (using the clustering method described in the Methods), based on the *P* values of the meta-analysis of the X chromosome (thresholding at the genome-wide significance level). The clustering of the meta-analysis X chromosome scan produced 23 clusters. Twenty of these had lead associations within 0.25 centimorgans (cM) of one of the discovery cohort clusters derived from the original GWAS that pooled all individuals of both genetic sexes (indicating strong concordance between the meta-analysis and the discovery cohort). Each of the four discovery cohort X chromosome clusters with lead association at the Bonferroni level (the first four rows of Table 1) overlap with a meta-analysis cluster, suggesting that these main clusters are not confounded by genetic sex. Of the remaining meta-analysis clusters, three do not overlap with any of the discovery cohort clusters. The lead-RSID/lead-phenotype pairs of these three non-overlapping clusters are rs5990961/V3742, rs142994659/V1233 and rs764953454/V3919 (the mapping between the

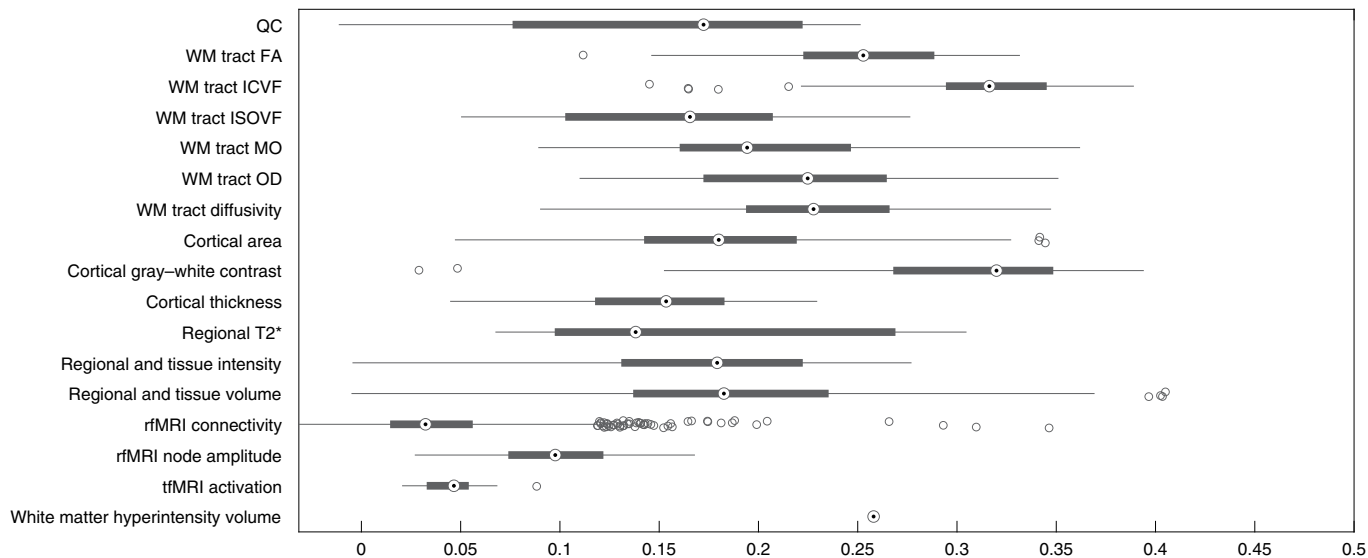


Fig. 2 | Heritability estimates (h^2) for phenotypes grouped according to IDP categories. Acronyms in y labels include quality control (QC), resting-state functional MRI (rfMRI), task functional MRI (tfMRI), and diffusion imaging phenotypes: white matter (WM), fractional anisotropy (FA), intra-cellular volume fraction (ICVF), isotropic or free water volume fraction (ISOVF), diffusion tensor mode (MO) and orientation dispersion index (OD). Box plots indicate medians, 25th and 75th quantiles and whiskers extending to maximal and minimal non-outliers (outliers are points exceeding 1.5 times the interquartile range from the median). More details for these 17 categories and heritability and standard errors for all phenotypes are provided in Supplementary Table 1.

#	RSID	Phenotype ID and brief description	Position	a1	a2	$-\log_{10}(P)$	n_{SNPs}
1	rs2272737	1943 White matter neurite density in left superior longitudinal fasciculus	152876080	C	T	20.459	97
2	rs62595479	1408 White-gray T1 intensity contrast in right fusiform gyrus	2163736	T	C	16.088	21
3	rs644138	0818 Area of left lateral-occipital cortex	154927581	G	A	14.319	35
4	rs12843772	0834 Area of left rostral-middle-frontal cortex	136648126	C	T	11.292	17
5	rs188847578	0076 Volume of gray matter in left supplementary motor cortex	127661277	G	T	10.247	3
10	rs193203210	1437 Volume of white matter hyperintensities	152634500	C	G	8.711	2
11	rs11539157	0202 Volume of left ventral diencephalon	68381264	C	A	8.659	1
15	rs149776026	3512 Functional connectivity, connection 1,084, dimensionality 100	124604066	A	G	8.193	1
16	rs5916169	2510 Functional connectivity, connection 82, dimensionality 100	5507316	C	T	8.12	1
23	rs6527976	1437 Volume of white matter hyperintensities	13808841	A	G	7.981	1
25	rs5930655	0136 Volume of gray matter in brain stem	133781599	A	C	7.925	1
30	rs6644158	0244 Volume of subiculum body in left hippocampus	113638468	G	A	7.723	1

References to the UKB phenotype definitions are provided in column 3. Summary statistics beta and s.e. are provided in Supplementary Table 1. The $-\log_{10}(P)$ values provided are for the main discovery cohort. Column n_{SNPs} indicates the number of peak phenotype/variant pairs included in the cluster. Bolded reference SNP ID (RSID) numbers indicate significance at the Bonferroni threshold in the main discovery cohort.

phenotype numbers and phenotype names is given in Supplementary Table 1). However, none of these three non-overlapping clusters achieved Bonferroni significance in the meta-analysis.

The differences in sensitivity among the main GWAS (pooling genetic male and female samples in the discovery cohort), the sex-specific GWAS and the Fisher meta-analysis are visualized concisely in Fig. 3. The histograms show the distributions of paired-difference $-\log_{10}(P)$ values. For the sex-specific comparisons, there are SNP/phenotype pairs having reduced sensitivity compared to the original all-subjects GWAS (likely due to reduced statistical power because of reduced subject numbers) and other pairs with increased sensitivity (likely because a given association is stronger for the sex in question than for the other sex). The meta-analysis

paired-difference distribution demonstrates that sex-separated GWASs followed by meta-analysis give increased sensitivity to finding genetic associations in the X chromosome.

Investigation of the four main X chromosome clusters. We now examine the four main X chromosome clusters in greater detail. The details for these additional investigations are summarized in Supplementary Table 5.

Cluster 1 comprises five SNPs, associated in total with 96 IDPs, all capturing differences in the properties of white matter tracts distributed throughout the cerebrum. This is described more in Table 2 and Supplementary Table 5. The top SNP (rs2272737, $P = 3.5 \times 10^{-21}$) is located about 10 kb away from, and is an eQTL

Table 2 | Context for the four X chromosome clusters with significance at the Bonferroni threshold

Phenotypes for cluster	Gene	Position	eQTLs	eQTL tissues	Sex
1. Diffusion MRI measures in distributed white matter tracts	<i>FAM58A</i> (<i>CCNQ</i>) intergenic	Xq28	<i>FAM58A</i>	Brain—cerebellum	—
2. Gray-white matter intensity contrast (temporal, limbic, default mode network)	<i>DHRX</i> intron	Xp22.33 / Yp11.2	<i>DHRX</i>	Blood, artery, heart, musculoskeletal	M
3. Occipital lobe gray matter area and volume, brainstem volume	<i>SPRY3</i> intron	Xq28	<i>F8A1</i> , <i>F8</i> , <i>BRCC3</i> , <i>TMLHE</i> , <i>RAB39B</i> , <i>CLIC2</i>	Brain—cortex, subcortex, cerebellum	—
4. Fronto-parietal gray matter area and volume	<i>ZIC3</i> intergenic	Xq26.3	<i>RP11-158M9.1</i>	Brain—cerebellum	—

Cluster number (bolded and matching cluster numbering in the Supplementary Tables) and a general description of the phenotypes involved in the cluster are given in the first column. Cytogenetic positions are provided: clusters with both X and Y positions are in a PAR. Note that the cytogenetic position Xq28 is on the edge of the non-PAR region and overlaps with the Y chromosome Yq12, although the lead association is in non-PAR. Information about eQTLs is provided. Column 'Sex' indicates whether the sex-specific scan is significant at the Bonferroni threshold of $-\log_{10}(P) \geq 11.1$ for one genetic sex but not significant for the other genetic sex (with M indicating the genetic sex for sex-driven effect in males and dash indicating no such sex-driven effect).

of, *FAM58A* (or *CCNQ*). The eQTLs reported throughout this work were assessed as significant according to the Genotype-Tissue Expression (GTEx) Project. Mutations in this gene lead to STAR syndrome, a rare X-linked developmental disorder recently identified⁹, for which notable brain variations have been observed, such as incomplete hippocampal inversion, thin corpus callosum, ventriculomegaly and cerebellar hypoplasia¹⁰. In addition, although *FAM58A* codes for an orphan cyclin with undescribed function, it has recently been shown to interact with *CDK10* (ref. ¹¹). Of particular relevance considering the many white matter IDPs associated with Cluster 1, mutation of the gene *CDK10* has been observed in a case study to lead to a rudimentary corpus callosum and paucity of white matter surrounding the lateral ventricles¹¹.

The SNPs of Cluster 1 are further associated with an array of non-imaging-derived phenotypes (nIDPs—phenotypes not derived from MRI) largely related to health (including diagnosed diseases and operative procedures), as well as some variables not necessarily health related (Supplementary Table 5). Interestingly, one SNP in Cluster 1 (rs1894299) was seen previously in a GWAS of type 2 diabetes¹². This SNP is located in an intron of *DUSP9* (*MKP4*), a gene that codes for a phosphatase whose overexpression specifically protects against stress-induced insulin resistance. This might be related to another consistent aspect of these nIDP associations with Cluster 1: the diet of UK Biobank participants with intake of sweet food and drinks (including desserts, puddings, beer and cider).

Cluster 2 comprises nine SNPs associated altogether with 17 IDPs, all of which are gray matter versus white matter intensity contrast, in limbic and temporal regions and in brain areas making up the default mode network¹³. The top SNP (rs62595479, $P = 8.2 \times 10^{-17}$) is located in a PAR of chromosome X (that is, a genetic region homologous between chromosomes X and Y) in an intron of *DHRX* and is an eQTL of the same gene. The genetic association with the gray-white intensity contrast IDP for this SNP was mainly driven by the male UK participants (as described above). The male-dominated aspect of the association between *DHRX* and the brain has also been observed in a study showing that four PAR genes, including *DHRX* (and *SPRY3*, see below), are upregulated in the blood of genetic male patients with ischemic stroke¹⁴.

Although most of the nIDP associations with the SNPs of Cluster 2 were related to diagnoses and operative procedures, half of these pointed to thyroid-related issues, in addition to the nIDP of workplace temperature, which might be related to thyroid function (Supplementary Table 5). Remarkably, the distribution of thyroid function modulation in the brain appears to consistently follow (in both positron emission tomography and functional MRI studies) that of the 17 regions associated with Cluster 2: mainly limbic and temporal areas, including the posterior cingulate cortex, orbitofrontal cortex, parahippocampal and fusiform gyrus¹⁵.

Cluster 3 includes nine genetic variants associated with 28 IDPs of local brain volume. All 28 IDPs are located in the occipital lobe, except for the volume of the brainstem and fourth ventricle. The top genetic variant (rs644138, $P = 4.8 \times 10^{-15}$) is located in a PAR in an intron of *SPRY3*. This genetic variant is also an eQTL in many brain regions of a variety of genes whose mutations are involved in developmental and neurodevelopmental disorders: *RAB39B*, which plays a role in normal neuronal development and dendritic process, is associated with cognitive impairment¹⁶, X-linked intellectual disability¹⁷ and Waisman syndrome in particular, an X-linked neurologic disorder characterized by delayed psychomotor development, impaired intellectual development and early-onset Parkinson's disease¹⁸; *TMLHE* is associated with X-linked autism¹⁹; *CLIC2* is associated with X-linked intellectual disability²⁰; and *BRCC3* is associated with an X-linked recessive syndromic form of moyamoya disease²¹. It is also an eQTL of *F8* and *F8A1* (*DXS522E/HAP40*), a likely candidate for the aberrant nuclear localization of mutant huntingtin in Huntington's disease. Considering the distribution of the brain IDPs in the occipital lobe, it, perhaps, lends additional credence to the consistent, but not yet understood, observation of volumetric and sulcal differences in these visual gray matter regions in Huntington's disease gene carriers²².

Another aspect of the genes for which the top genetic variant of Cluster 3 is an eQTL is cardiovascular issues—for instance, mutant *CLIC2* leads to atrial fibrillation, cardiomegaly and congestive heart failure²⁰. *F8* encodes a large plasma glycoprotein that functions as a blood coagulation factor, whereas mutations in *BRCC3* are linked to moyamoya syndrome, a rare blood vessel disorder in which certain arteries in the brain are blocked or constricted and that is accompanied by other symptoms, including hypertension, dilated cardiomyopathy and premature coronary heart disease²¹. In line with this, we found that Cluster 3 was associated in UKB participants with diagnosis of atrioventricular block and ventricular premature depolarization, as well as an operative procedure consisting of the replacement of two coronary arteries. In addition, Cluster 3 was consistently associated with many measures of physical growth, perhaps in line with the role of *SPRY3* as a fibroblast growth factor antagonist in vertebrate development, including height, lung function and capacity and body mass (Supplementary Table 5).

Finally, Cluster 4 comprised eight genetic variants associated with volume of gray matter regions in the dorsolateral prefrontal cortex and the lateral parietal cortex (supramarginal gyrus and opercular cortex). The top genetic variant in this cluster (rs12843772, $P = 5.1 \times 10^{-12}$) is located just <150 bp from *ZIC3*, which plays a key role in body pattern formation and left-right asymmetry. Mutations in this gene are thought to be involved in 1% of heterotaxy (situs ambiguous and inversus) in humans²³. This might explain why

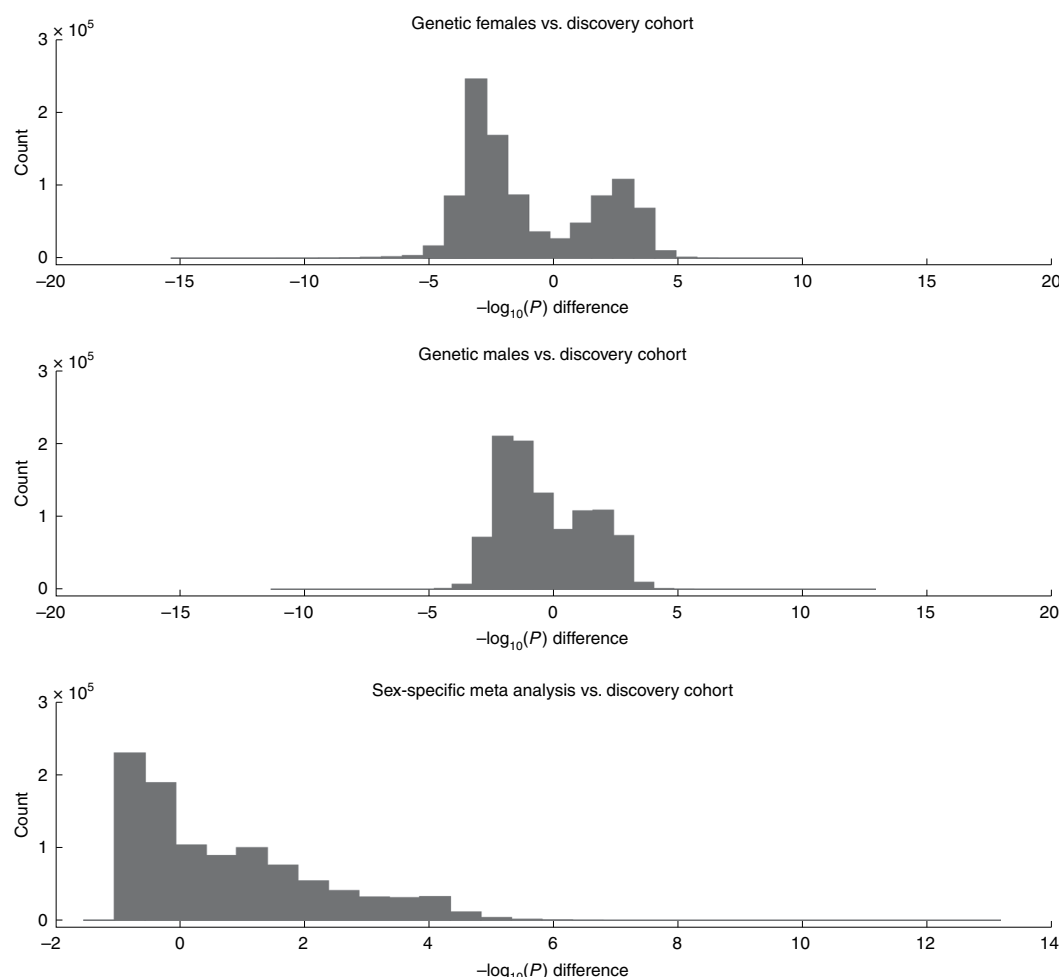


Fig. 3 | Paired-difference histograms for the sex-specific scans on the X chromosome. We plot histograms for the differences between the $-\log_{10}(P)$ values for genetic females (top), genetic males (middle) and the meta-analysis (bottom) versus the discovery scan (which includes genetic male and female samples together but did include a sex confound covariate). Differences are plotted for all associations for which the maximum $-\log_{10}(P)$ value over the four analyses is greater than 4.0, leading to the bimodal nature of the first two histograms. A total of 989,981 variant-IDP pairs pass this maximum filter. The bottom plot shows that there is greater statistical sensitivity when carrying out sex-specific GWASs on the X chromosome, and then combining the results with a meta-analysis, than by combining all individuals together in a simple standard GWAS.

the distribution of the higher-order gray matter regions associated with Cluster 4 follows the pattern of the fronto-parietal networks, which are notoriously left-versus-right segregated²⁴. In particular, the supramarginal gyrus is known to show the strongest asymmetries from an early developmental stage²⁵ and to be connected by white matter tracts that share a genetic influence with human handedness²⁶. *ZIC3* is also involved in neural tube development and closure, and mutations in this gene cause neural tube defects and cerebellum hypoplasia, consistent histological brain alterations with abnormal laterality and axial patterning, including a disorganized cerebral cortex²⁷.

Cluster 4 is also, in particular, strongly associated with insulin-like growth factor 1 (IGF-1) levels in the blood of UKB participants (Supplementary Table 5). IGF-1, in particular, controls brain development, plasticity and repair²⁸. More recently, it has emerged as a risk factor for dementia and, particularly, Alzheimer's disease²⁹, as it is a major regulator of $A\beta$ physiology and controls $A\beta$ clearance from the brain³⁰. Zoomed-in views of these associations are displayed in Extended Data Fig. 2.

With respect to XCI status of the genes discussed in this section, the gene *DHSRX* is located in PAR1 and is fully expressed in females³¹. Both *SPRY3* and *ZIC3* are subject to full inactivation³².

Investigation of five autosomal clusters. Of the 692 replicating clusters with GWAS-significant P value associations reported here, 48 involved a lead association among the autosomes with $-\log_{10}(P) \geq 11.1$ and replication (at nominal level) and were distinct from the loci reported in previous work^{1,33,34} (that is, more than 0.5 cM from a reported genetic variant and more than 1 Mbp from a gene reported as associated in refs. ^{33,34}). Of these 48 clusters, seven involved associations with phenotypes from among the three new classes of UKB brain imaging phenotypes analyzed here. Of these seven clusters, two (Cluster 271 with lead variant rs1368575 and Cluster 1070 with lead association rs3814883) were not related to previous brain imaging literature, were not intergenic and were without eQTLs. The remaining five clusters (Clusters 163, 270, 841, 1029 and 1067 from Supplementary Table 6) are investigated below. We provide a tabulation of all autosomal lead associations for our clusters with a cross-reference to refs. ^{1,33,34} in Supplementary Table 6.

The first of these most novel clusters related to previous brain imaging work (Cluster 163) included 12 distinct replicated variants associated with 31 IDPs, all but two related to the contrast between the intensity of the white matter and gray matter (tissue contrast (TC)) generated by FreeSurfer, widespread across

multiple cortical regions of interest³⁵. The top variant for this cluster (rs3832092, $P = 1.9 \times 10^{-13}$) is in an exon and an eQTL of *MARS2*, a mitochondrial methionyl-tRNA synthetase, and is associated in particular with parietal TC (precuneus and inferior and superior parietal cortex). Mutations of *MARS2* have been linked to spastic ataxia and neurodevelopmental delay as well as white matter abnormalities (leukoencephalopathy), cortical and cerebellar atrophy and corpus callosum thinning³⁶. This variant is further associated primarily with many nIDPs related to mass and body fat as well as nIDPs involving allergy, sleep and red blood cell count and shape (according to Open Targets Genetics).

The second cluster (Cluster 270) encompassed 13 replicated variants that were predominantly associated with the TC of 48 different cortical regions across the whole cerebrum. Its top variant (rs9290432, $P = 2.4 \times 10^{-17}$), associated mainly with temporal, parietal and prefrontal TC, is in an intron and eQTL of the gene *PLD1*. This gene codes for a phospholipid enzyme that has been shown to regulate the trafficking of the protein β APP³⁷, the precursor of amyloid beta, plaques of which are a critical factor in the pathogenesis of Alzheimer's disease. Increased expression of *PLD1* has been noted in the brains of patients with Alzheimer's disease, both in the hippocampus and temporal lobe³⁸, and *PLD1* is upregulated in the mitochondrial membrane in brains of patients with Alzheimer's disease³⁹. The top genetic variant of this cluster is, in addition, related to nIDPs of red blood cell count and shape (Open Targets Genetics).

The third cluster (Cluster 841) consisted of eight replicated loci all associated solely with measures of TC in 15 mainly higher-order cortical areas. We found that the lead variant in this cluster (rs549893, $P = 1.2 \times 10^{-12}$) is associated with TC in the superior temporal cortex, is located in an intron of *HMBS* and is an eQTL of both *HMBS* and *VPS11* in many brain tissues (cortical, subcortical and cerebellar). The gene *HMBS* codes for an enzyme of the biosynthetic pathway of heme production. As such, mutations in *HMBS* are known to be related to acute intermittent porphyria, an autosomal dominant defect in the biosynthesis of heme. Acute intermittent porphyria manifests itself mainly as a neurological disorder, involving the autonomous, central and peripheral nervous systems, with acute, life-threatening neurologic attacks. Deep cerebral white matter myelination abnormalities⁴⁰ and axonal neuropathy⁴¹ have been observed in patients and animal models, as well as deficiencies in mitochondrial complexes in *Hmbs* mutant brain, suggesting that mitochondrial energetic failure also plays an important role in the expression of the disease⁴². In addition, mutations in *VPS11* in five patients have been suggested as leading to infantile-onset leukoencephalopathy with brain white matter abnormalities, severe motor impairment, cortical blindness, intellectual disability and seizures, as well as in a significant reduction in myelination after extensive neuronal death in the hindbrain and midbrain in an animal model⁴³. The variant rs549893 is associated with nIDPs of body mass index, body fat mass and body fat percentage (Open Targets Genetics).

The fourth genetic cluster (Cluster 1029) comprised 17 replicated loci related, again, predominantly to TC in 45 widespread cortical areas. The lead variant for this cluster (rs140648465, $P = 4.1 \times 10^{-16}$) is associated with TC measures in the supramarginal, inferior and superior parietal cortex, and it is intergenic and is an eQTL in cortical and subcortical tissues of *RMDN3/FAM82A2*, coding for a regulator of microtubule dynamics. This tether protein is involved in facilitating the lipid transfer by increasing the contact sites between the endoplasmic reticulum and mitochondria, particularly in the brain⁴⁴. The protein FAM82A2 was found downregulated in the frontal and parietal lobes in patients with Parkinson's disease with dementia, with a dramatic reduction in the activity of the mitochondrial complexes⁴⁵. Recently, reduced levels of the RMDN3 protein (sometimes also named PTPIP51) have also been found in the temporal cortex of the brains of patients with Alzheimer's disease, suggesting that a disruption of endoplasmic reticulum-mitochondria

interactions mediated by *RMDN3* might be part of the neuropathological process in Alzheimer's disease⁴⁶. A proxy genetic variant for the top variant of this fifth cluster (rs8042729, in linkage disequilibrium (LD) $R^2 = 1.00$ with rs140648465) is further associated with nIDPs of blood count and hypertension (Open Targets Genetics).

The final interpretable cluster (Cluster 1067) comprised 31 replicated loci associated with a total of 51 IDPs of TC spread across higher-order cortical regions as well as 19 other IDPs of cortical thickness (superior prefrontal and parietal), regional volumes (hippocampus, thalamus and choroid plexus) and white matter diffusion. This cluster's most significant variant (rs2549095, $P = 5.3 \times 10^{-20}$) is associated with TC in an array of temporal, parietal and prefrontal cortical areas. This locus is in an exon of the gene *CLEC18A*, which codes for a C-type lectin that has only been recently characterized, and has been shown to be expressed in the brain, in microglia⁴⁷. This locus is also part of a 40-marker panel of genes that makes it possible to distinguish cardioembolic from large-vessel ischemic stroke with high accuracy⁴⁸. A proxy locus for the top genetic variant (rs72785089, LD $R^2 = 0.81$) is further associated with measures of body mass and percentage as well as blood count and alcohol intake frequency (Open Targets Genetics).

Discussion

A major component in the expansion of the UKB prospective epidemiological resource is the addition of tens of thousands of newly imaged participants and the increase in the richness of phenotypes that can be derived from the imaging data. Since we published our first large-scale GWAS of UKB brain imaging in 2018, the brain imaging has almost reached its halfway point, having now scanned nearly 50,000 volunteers, of which 40,000 samples have already passed quality control tests and been processed and released for use in research. As a result, the size of the available discovery sample for GWASs has nearly tripled, and the number of genetic variants passing reliability thresholds (such as $MAF \geq 1\%$) has increased by 30%, now reaching 10 million. It is now, therefore, a good time to update our large-scale GWASs, now with almost 4,000 imaging-derived phenotypes—thousands of distinct measures of brain structure, function, connectivity and microstructure. The number of replicated clusters of imaging-genetic associations identified has more than quadrupled since our previous work. We have made all of the GWAS summary statistics openly available via the new BIG40 brain imaging genetics server.

We also studied brain imaging associations in the X chromosome and autosomal associations with novel phenotypes. We identified four X chromosome clusters that replicate at the GWAS + Bonferroni level (increasing the standard 7.5 level according to Bonferroni correction for the number of IDPs tested). Among these four top X chromosome clusters, we found associations involving diffusion measures distributed in white matter tracts and gray-white matter intensity contrasts. We also found associations involving occipital lobe gray matter and fronto-parietal gray matter. These associations are relevant to a diverse set of variations in brain development and pathologies, including the recently identified STAR syndrome, Waisman syndrome, early-onset Parkinson's disease, X-linked autism, Alzheimer's disease and Huntington's disease. Cardiovascular conditions, such as ischemic stroke, moyamoya syndrome and premature coronary heart disease, are also implicated. The X chromosome genes *FAM58A* (*CCNQ*), *DHRXS*, *SPRY3*, *F8A1*, *F8*, *BRCC3*, *TMLHE*, *RAB39B*, *CLIC2* and *ZIC3* are involved in the associations we report.

The X chromosome is typically understudied in GWASs, and many of the X chromosome loci that we identify are now implicated in genetic associations with brain imaging phenotypes for the first time. Ploidy, the Barr body and potential confounding with genetic sex complicate X chromosome analysis. To address this, we performed sex-specific GWASs and a meta-analysis combining

the sex-separated analyses. We showed significant overlap between the sex-separated meta-analysis and the main GWAS (the discovery cohort pooling genetic females and males), providing evidence against confounding (and we note that the significant sex-specific effect signs never differed for genetic females and males for variants with $MAF \geq 1\%$). This also allowed us to investigate if a given association was sex driven (the association might be significant for only one genetic sex or have significantly different effect sizes between the two sex-separated GWASs). For example, we found evidence that Cluster 2 for our brain–gene associations is primarily driven by associations in genetic males.

Four of the autosomal clusters of association with gray–white tissue contrast were related to mitochondrial proteins and function. In the brain, mitochondrial disorders manifest themselves primarily as demyelinated lesions in the white matter or as stroke-like lesions in both gray and white matter. More rarely, they present with cortical and subcortical necrosis⁴⁹. Although TC shows strong changes with aging in prefrontal, lateral parietal, superior temporal and precuneus regions—those overall demonstrating the most significant associations with genetic variants—it is unclear what underlying tissue properties might be driving these effects in the MRI signal ratio³⁵. In this study, TC, being almost entirely the only type of IDP showing significant genome-wide associations with our mitochondrial-related genetic clusters, appears to be specific and sensitive to mitochondrial function. These genetic imaging findings are likely in line with the recent discovery that mitochondrial dysfunction, which emerges early during the aging process, might prompt the catabolism of the myelin lipids of the white matter as an adaptive response to address brain fuel and energy demand⁵⁰. This catabolic process perhaps explains the results previously observed with TC in aging³⁵.

Enhancing understanding of the mapping between genotype and phenotype might lead to advances in neuroscience and improvements in outcomes for brain pathologies. Deeply phenotyped resources, such as UKB, provide an opportunity to update the known genotype/phenotype mappings for a wide spectrum of brain imaging phenotypes as more samples and more phenotypes are released. It is crucial that such updates are provided on open platforms, accelerating the potential impact of brain imaging genetics. We have done that here with the open BIG40 resource (with a total of 1,282 clusters of associations) and with freely available summary statistics. We hope that this dissemination will be valuable for the next generation of neuroscience research.

Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41593-021-00826-4>.

Received: 27 July 2020; Accepted: 23 February 2021;
Published online: 19 April 2021

References

1. Elliott, L. T. et al. Genome-wide association studies of brain imaging phenotypes in UK Biobank. *Nature* **562**, 210–216 (2018).
2. F. Alfaro-Almagro, et al. Confound modelling in UK Biobank brain imaging. *Neuroimage* (in the press).
3. Hormozdiari, F., Kostem, E., Yong Kang, E., Pasaniuc, B. & Eskin, E. Identifying causal variants at loci with multiple signals of association. *Genetics* **198**, 497–508 (2014).
4. Clayton, D. Testing for association on the X chromosome. *Biostatistics* **9**, 593–600 (2008).
5. Lee, J. J. et al. Gene discovery and polygenic prediction from a genome-wide association study of educational attainment in 1.1 million individuals. *Nat. Genet.* **50**, 1112 (2018).
6. Özbek, U. et al. Statistics for X-chromosome associations. *Genet. Epidemiol.* **42**, 539–550 (2018).
7. Bulik-Sullivan, B. K. et al. LD score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.* **47**, 291–295 (2015).
8. Saleem F. & Rizvi S. W. Transgender associations and possible etiology: a literature review. *Cureus* **9**, e1984 (2017).
9. Unger, S. et al. Mutations in the cyclin family member *FAM58A* cause an X-linked dominant disorder characterized by syndactyly, telecanthus and anogenital and renal malformations. *Nat. Genet.* **40**, 287–289 (2008).
10. Bedeschi, M. F. et al. STAR syndrome plus: the first description of a female patient with the lethal form. *Am. J. Med. Genet.* **173**, 3226–3230 (2017).
11. Guen, V. J. et al. A homozygous deleterious *CDK10* mutation in a patient with agenesis of corpus callosum, retinopathy, and deafness. *Am. J. Med. Genet.* **176**, 92–98 (2018).
12. Suzuki, K. et al. Identification of 28 new susceptibility loci for type 2 diabetes in the Japanese population. *Nat. Genet.* **51**, 379–386 (2019).
13. Raichle, M. E. et al. A default mode of brain function. *Proc. Natl Acad. Sci. USA* **98**, 676–682 (2001).
14. Tian, Y. et al. Y chromosome gene expression in the blood of male patients with ischemic stroke compared with male controls. *Genet. Med.* **9**, 68–75 (2012).
15. Zhang, W. et al. Disrupted functional connectivity of the hippocampus in patients with hyperthyroidism: evidence from resting-state fMRI. *Eur. J. Radiol.* **83**, 1907–1913 (2014).
16. Vanmarsenille, L. et al. Increased dosage of *RAB39B* affects neuronal development and could explain the cognitive impairment in male patients with distal Xq28 copy number gains. *Hum. Mutat.* **35**, 377–383 (2014).
17. Giannandrea, M. et al. Mutations in the small GTPase gene *RAB39B* are responsible for X-linked mental retardation associated with autism, epilepsy, and macrocephaly. *Am. J. Hum. Genet.* **86**, 185–195 (2010).
18. Wilson, G. R. et al. Mutations in *RAB39B* cause X-linked intellectual disability and early-onset Parkinson disease with α -synuclein pathology. *Am. J. Hum. Genet.* **95**, 729–735 (2014).
19. Celestino-Soper, P. B. S. et al. Use of array CGH to detect exonic copy number variants throughout the genome in autism families detects a novel deletion in *TMLHE*. *Hum. Mol. Genet.* **20**, 4360–4370 (2011).
20. Takano, K. et al. An X-linked channelopathy with cardiomegaly due to a *CLIC2* mutation enhancing ryanodine receptor channel activity. *Hum. Mol. Genet.* **21**, 4497–4507 (2012).
21. Miskinyte, S. et al. Loss of BRCC3 deubiquitinating enzyme leads to abnormal angiogenesis and is associated with syndromic moyamoya. *Am. J. Hum. Genet.* **88**, 718–728 (2011).
22. Rosas, H. D. et al. Cerebral cortex and the clinical expression of Huntington's disease: complexity and heterogeneity. *Brain* **131**, 1057–1068 (2008).
23. Ware, S. M. et al. Identification and functional analysis of *ZIC3* mutations in heterotaxy and related congenital heart defects. *Am. J. Hum. Genet.* **74**, 93–105 (2004).
24. Witt, S. T., van Ettinger-Veenstra, H., Salo, T., Riedel, M. C. & Laird, A. R. What executive function network is that? An image-based meta-analysis of network labels. Preprint at *bioRxiv* <https://doi.org/10.1101/2020.07.14.201202> (2020).
25. Dubois, J. et al. Structural asymmetries of perisylvian regions in the preterm newborn. *Neuroimage* **52**, 32–42 (2010).
26. Wiberg, A. et al. Handedness, language areas and neuropsychiatric diseases: Insights from brain imaging and genetics. *Brain* **142**, 2938–2947 (2019).
27. Purandare, S. M. et al. A complex syndrome of left-right axis, central nervous system and axial skeleton defects in *Zic3* mutant mice. *Development* **129**, 2293–2302 (2002).
28. Dyer, A. H., Vahdatpour, C., Sanfeliu, A. & Tropea, D. The role of insulin-like growth factor 1 (IGF-1) in brain development, maturation and neuroplasticity. *Neuroscience* **325**, 89–99 (2016).
29. Westwood, A. J. et al. Insulin-like growth factor-1 and risk of Alzheimer dementia and brain atrophy. *Neurology* **82**, 1613–1619 (2014).
30. Bates, K. A. et al. Clearance mechanisms of Alzheimer's amyloid- β peptide: implications for therapeutic design and diagnostic tests. *Mol. Psychol.* **14**, 469–486 (2009).
31. Carrel, L. & Willard, H. F. X-inactivation profile reveals extensive variability in X-linked gene expression in females. *Nature* **434**, 400–404 (2005).
32. Balaton, B. P., Cotton, A. M. & Brown, C. J. Derivation of consensus inactivation status for X-linked genes from genome-wide studies. *Biol. Sex. Differ.* **6**, 1–11 (2015).
33. Zhao, B. et al. Genome-wide association analysis of 19,629 individuals identifies variants influencing regional brain volumes and refines their genetic co-architecture with cognitive and mental health traits. *Nat. Genet.* **51**, 1637–1644 (2019).
34. Zhao, B. et al. Large-scale GWAS reveals genetic architecture of brain white matter microstructure and genetic overlap with cognitive and mental health traits ($n = 17,706$). *Mol. Psychiatry* <https://doi.org/10.1038/s41380-019-0569-z> (2019).

35. Salat, D. H. et al. Age-associated alterations in cortical gray and white matter signal intensity and gray to white matter contrast. *Neuroimage* **48**, 21–28 (2009).
 36. Thiffault, I. et al. A new autosomal recessive spastic ataxia associated with frequent white matter changes maps to 2q33–34. *Brain* **129**, 2332–2340 (2006).
 37. Cai, D. et al. Phospholipase D1 corrects impaired β APP trafficking and neurite outgrowth in familial Alzheimer's disease-linked presenilin-1 mutant neurons. *Proc. Natl Acad. Sci. USA* **103**, 1936–1940 (2006).
 38. Krishnan, B., Kaye, R. & Taglialetta, G. Elevated phospholipase D isoform 1 in Alzheimer's disease patients' hippocampus: relevance to synaptic dysfunction and memory deficits. *Alzheimers Dement. (N Y)* **4**, 89–102 (2018).
 39. Jin, J. K. et al. Phospholipase D1 is up-regulated in the mitochondrial fraction from the brains of Alzheimer's disease patients. *Neurosci. Lett.* **407**, 263–267 (2006).
 40. Solis, C. et al. Acute intermittent porphyria: studies of the severe homozygous dominant disease provides insights into the neurologic attacks in acute porphyrias. *Arch. Neurol.* **61**, 1764–1770 (2004).
 41. Lindberg, R. L. et al. Porphobilinogen deaminase deficiency in mice causes a neuropathy resembling that of human hepatic porphyria. *Nat. Genet.* **12**, 195–199 (1996).
 42. Hamedan, C. et al. Mitochondrial energetic defects in muscle and brain of a *Hmbs*^{-/-} mouse model of acute intermittent porphyria. *Hum. Mol. Genet.* **24**, 5015–5023 (2015).
 43. Zhang, J. et al. A founder mutation in VPS11 causes an autosomal recessive leukoencephalopathy linked to autophagic defects. *PLoS Genet.* **12**, e1005848 (2016).
 44. Fecher, C. et al. Cell-type-specific profiling of brain mitochondria reveals functional and molecular diversity. *Nat. Neurosci.* **22**, 1731–1742 (2019).
 45. Garcia-Esparcia, P. et al. Mitochondrial activity in the frontal cortex area 8 and angular gyrus in Parkinson's disease and Parkinson's disease with dementia. *Brain Pathol.* **28**, 43–57 (2018).
 46. Lau, D. H. et al. Disruption of endoplasmic reticulum-mitochondria tethering proteins in post-mortem Alzheimer's disease brain. *Neurobiol. Dis.* **143**, 105020 (2020).
 47. Huang, Y. L. et al. Human CLEC18 gene cluster contains C-type lectins with differential glycan-binding specificity. *J. Biol. Chem.* **290**, 21252–21263 (2015).
 48. Jickling, G. C. & Sharp, F. R. Biomarker panels in ischemic stroke. *Stroke* **46**, 915–920 (2015).
 49. Finsterer, J. Central nervous system imaging in mitochondrial disorders. *Can. J. Neurol. Sci.* **36**, 143–153 (2009).
 50. Klosinski, L. P. et al. White matter lipids as a ketogenic fuel supply in aging female brain: implications for Alzheimer's disease. *EBioMedicine* **2**, 1888–1904 (2015).
- Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.
- © The Author(s), under exclusive licence to Springer Nature America, Inc. 2021

Methods

We describe the image processing and genetic pre-processing and the association studies, including the procedures that are X chromosome specific, and the details of our sex-specific analyses. Specific details about our protocols are provided in the Life Sciences Reporting Summary.

Image processing. We used brain IDPs from the “40k” (approximately 40,000 participants) UKB data release in early 2020, as processed by WIN FMRIB on behalf of UKB⁵¹. After removal of individuals as part of the genetic processing (see below), we used data from 33,224 individuals. These were then randomly split into a discovery sample of 22,138 individuals (11,624 genetic females) and a replication sample of 11,086 individuals (5,787 genetic females). The ages in the discovery sample were as follows: females: mean age = 63.6 ± 7.3 years, min = 45.1 years, max = 81.8 years; males: mean age = 65.0 ± 7.6 years, min = 46.1 years, max = 81.8 years. In the replication sample: females: mean age = 63.7 ± 7.4 years, min = 46.3 years, max = 81.6 years; males: mean age = 65.0 ± 7.6 years, min = 46.1 years, max = 81.0 years. The exact numbers of individuals vary across IDPs, according to patterns of missing data, with the maximum numbers given above (for IDPs with no missing data) and the minimum numbers being just 16% lower. The BIG40 online table listing the IDPs includes the exact number of individuals (in discovery and replication samples and in the sex-specific GWAS) for each IDP. The details for these IDPs (including long descriptions, category names and units) are summarized in Supplementary Table 1. The sample sizes we use represent the largest-to-date sizes released by UKB. No statistical methods were used to determine these sample sizes.

As described in detail in ref. ⁵², the UKB data includes six MRI modalities: T1-weighted and T2-weighted fluid-attenuated inversion recovery (T2-FLAIR) structural images, susceptibility-weighted MRI, diffusion MRI, task functional MRI and resting-state functional MRI. We (and colleagues) developed and applied an automated image processing pipeline on behalf of UKB⁵¹: <https://www.fmrrib.ox.ac.uk/ukbiobank/fbp>. This removes artifacts and renders images comparable across modalities and participants; it also generates thousands of IDPs—distinct measures of brain structure and function.

In this work, we used the 3,913 IDPs available from UKB, spanning a range of structural, diffusion and functional MRI summary measures (described in the central UKB brain imaging documentation (http://biobank.ctsu.ox.ac.uk/showcase/showcase/docs/brain_mri.pdf) and listed in full on the BIG40 server (<https://open.win.ox.ac.uk/ukbiobank/big40/>)).

We also used 16 quality control measures available from UKB as well as six compact summary functional connectivity features (derived from the hundreds of individual connectivity features¹). In this paper, we refer to all of the above 3,935 measures together as ‘the IDPs’. Each IDP’s *Nsubjects* × 1 data vector had outliers removed (set to missing, with outliers determined by being greater than six times the median absolute deviation from the median). We discarded individuals for whom 50 or more IDPs were missing (for any reason, which could be due to data acquisition incompleteness, data quality problems as described in ref. ⁵¹ or the above-described outlier removal).

The data were then split into discovery and replication samples, and the remaining steps below applied to each sample separately. Each IDP’s data vector was quantile normalized⁵², resulting in it being Gaussian distributed, with mean = 0 and s.d. = 1. Confounds were removed from the data, in a manner similar to that carried out in ref. ¹, including the 40 population genetic principal components supplied by UKB and with a greatly expanded new set of confounds². This includes confounds for age, head size, sex, head motion during functional MRI, scanner table position, imaging center and scan date-related slow drifts. To maximize GWAS interpretability, we regressed out all confounds listed and recommended in ref. ². For higher-order (non-linear and interaction) confounds, we used the same set of thresholds for automatic selection of these higher-order confounds. Given the slightly different set of individuals used here (for example, enforcing overlap with the genetics data) compared to those in ref. ², this resulted in the 602 ‘maximal’ set of confounds (reported in ref. ²) being reduced here to 597 confounds.

Genetic associations. We consider the 488,377 samples included in the Spring 2018 release of UKB and proceed with a pre-processing and discovery/reproduction paradigm similar to that described in ref. ¹. Of the samples, 39,944 are included among the 41,016 samples in UKB for which IDPs are available, after the genotyping quality control procedures for sample removal specified in ref. ⁵³. We removed samples without recent UK ancestry as determined by the *in.white.British.ancestry.subset* variable in the file *ukb_sqc_v2.txt* provided in the meta-data for UKB. This variable selects samples based on self-reported ancestry and genetic principal component thresholds. We also removed individuals based on relatedness, forming a maximally unrelated subset using the procedures recommended in ref. ⁵³. This resulted in a maximally unrelated subset of 34,298 samples with recent UK ancestry and accepted genotyping and imaging quality control. We divided this set into a discovery cohort and a reproduction cohort according to a random 2/3 and 1/3 (respective) proportion, resulting in a discovery cohort of 22,865 samples and a reproduction cohort of 11,433 samples.

For the X chromosome, fewer samples were available after exclusions due to variation in karyotype or additional aspects of genotype quality control described in ref. ⁵³. This resulted in discovery cohorts of 22,853 (instead of 22,865 as listed above) for the non-PAR and 22,844 for the PARs and replication cohorts of 11,430 for the non-PAR and 11,426 for the PARs (X chromosome aneuploidy was assessed using the putative.sex.chromosome.aneuploidy column of the genetic quality control file released by UKB⁵³).

Next, we considered quality control filters for the MAF, information score (INFO⁵⁴) and Hardy–Weinberg equilibrium (HWE). We applied filters for $MAF \geq 0.001$, $INFO \geq 0.3$ and $HWE -\log_{10}(P) \leq 7$. For the X chromosome, we applied these filters using the qctool software version 2.0.1 and the -infer-ploidy-from sex flag. This implies that genetic males contribute half as much as genetic females toward MAF and INFO for the non-PAR of the X chromosome. Furthermore, as has been recommended⁵⁵, for the X chromosome the HWE filters are applied after computing for genetic females only. After these filters, of the 93,095,645 genetic variants included in UKB in chromosomes 1–22, 16,445,196 remain. And of the 3,963,707 genetic variants in the X chromosome, 657,883 remain (639,835 in the non-PAR and 18,048 in the PARs).

GWASs. We performed linear association tests on the samples in the discovery cohort. We performed these tests between each of the 17,103,079 genetic variants and each of the 3,935 IDPs described above (10,119,893 of which have $MAF \geq 0.01$). The genotypes were provided by UKB, and details for imputation (including X chromosome imputation) and genetic principal component construction are provided in ref. ⁵³. We used bgenie software⁵³ to conduct the GWAS and record the effect sizes (beta), standard errors and $-\log_{10}(P)$ values for the associations. The effect sizes were recorded in the direction of the alternate allele. The phenotypes were scaled to have unit variance after deconfounding. The variants on chromosomes 1–22 were not scaled. Therefore, an effect size of 1.0 indicates that each copy of the alternate allele generally confers an increase in the phenotype by 1 s.d. For the non-PAR of X, the dosages for genetic males were scaled by a factor of 2.0 so that they lie in the range (0,2), respecting the Barr body.

We produced Manhattan plots for each of the 3,935 IDPs, plotting the $-\log_{10}(P)$ value for each variant (these Manhattan plots are provided on the BIG40 open web server, along with quantile plots). For the Manhattan plots, we applied an additional filter of $MAF \geq 0.01$ to all variants. A method for extracting hits (peak associations) from a GWAS was developed in ref. ¹. We applied that method to extract hits from the 3,935 scans with $-\log_{10}(P)$ values exceeding the GWAS threshold of 7.5 (again, with the $MAF \geq 0.01$ filter) and annotated the Manhattan plots with these hits (note that, in our tables, we also signified results passing a Bonferroni-corrected threshold for *P* values below $1.0 \times 10^{-7.5}$ divided by the number of IDPs: ~ 11.1). For each of the hits, we conducted a replication analysis by performing a linear association test on the samples in the replication cohort and recording the effect sizes, standard errors and $-\log_{10}(P)$ values for the replication. We extended the method from ref. ¹ to automatically generate clusters. A cluster is a set of phenotype/variant pairs such that each variant in the cluster is a peak association for its corresponding phenotype and such that all variants are within a 0.25-cM distance of the phenotype/variant pair with highest $-\log_{10}(P)$ value in the cluster. Each phenotype/variant pair identified as a hit appears in one and only one cluster. The details of this clustering method are provided in the Methods. We provide a software package called Peaks implementing these clustering methods and have released it under the open-source BSD 2-clause license.

Other methods for extracting lead associations include Bayesian methods such as CAVIAR³ for determination of causal genetic variants. We examined causal genetic variants in the top four X chromosome clusters by computing the LD matrix using plink2 (ref. ⁵⁶) for all variants within 250 kbp from the lead associations (our CAVIAR results are provided in Supplementary Table 3, and the results are summarized in the main text). Summary statistics for the associations of all genetic variants (with $MAF \geq 0.001$) for the discovery cohort (as well as the sex-separated and meta-analyses GWAS and a pooled discovery + replication GWAS; see below) are available for download on the BIG40 open web server.

Full scan. We also provide for download on the BIG40 open web server the summary statistics for a version of this GWAS conducted on the union of the discovery and replication cohorts considered in this study (that is, a maximal subset of unrelated samples with recent UK ancestry, among all samples in the UKB 2020 release of approximately 40,000 brain-imaged samples). The genetic variants considered in this scan are the same genetic variants passing our filters for the discovery cohort reported on in this paper (with $MAF \geq 0.001$). The sample size of each association test (after considering missing phenotypes and X chromosome exclusions due to aneuploidy) are also provided (the maximum number of included samples over all phenotypes in this full scan is 33,224).

Using the summary statistics from this full scan, we estimated the heritability of each phenotype. We used LD score regression⁷ to produce these estimates. LD scores were sourced from the European population of the 1000 Genomes Project⁵⁷. Results are listed in Supplementary Table 1.

Sex-specific scans. We also considered a scan for association on genetic male samples only and also a scan for association on genetic females only (in both cases, with samples from the discovery cohort). For these scans, in the non-pseudoautosomal X chromosome region, 11,885 genetic female samples were used and 10,968 genetic male samples were used. For the PARs, 11,882 genetic female samples were used and 10,962 genetic male samples were used. As in the main analysis, the number of samples per phenotype varies due to missingness (the sample sizes for each phenotype for genetic females and males are provided in Supplementary Table 1). We conducted these two scans for association between each deconfounded phenotype and each variant on the X chromosome and the autosome and recorded the effect sizes (beta), standard errors and $-\log_{10}(P)$ values. We then combined these two scans in a meta-analysis using Fisher's method⁵⁸, providing a $-\log_{10}(P)$ value that is more strongly controlled for sex-specific effects. The equation for the combined P value under this meta-analysis method is as follows:

$$1 - f_{\chi^2}(-2(\log P_m + \log P_f), 4) \quad (1)$$

Here, $f_{\chi^2}(\cdot, \nu)$ is the cumulative distribution function of a chi-squared random variable with ν degrees of freedom, and P_f and P_m are the P values of the genetic female and genetic male scans, respectively.

The Peaks algorithms. In GWASs with thousands of phenotypes, we must determine when two peak associations for two different phenotypes are related, in the sense that the peak variants are close together in genetic distance and/or LD. For studies in which the number of phenotypes (or strength and complexity of genetic associations) amounts to only a few supra-threshold results, such matching can be done by hand by examining recombination maps and LD diagrams. Such by-hand work is not feasible for studies with thousands of phenotypes. In this work, we provide a new automated method to co-register peak associations across many phenotypes. We provide a software package, Peaks, that implements these methods. The Peaks software provides an implementation of the algorithm in ref. ¹ (in the 'Identifying associated genetic loci' subsection of the Methods) for uncovering peak associations for each phenotype and an extension to co-localization of hits across phenotypes.

To combine peak associations into clusters that span phenotypes, we use a 'greedy' algorithm that delivers an optimally efficient clustering of genetic variant/phenotype pairs. The algorithm works by first identifying the peak associations for each phenotype (using the algorithm described in the previous subsection), then by iteratively extracting the genetic variant/phenotype pair with the top $-\log_{10}(P)$ and then by assigning all lead associations for all phenotypes within 0.25 cM of that genetic variant to the same cluster. The details are as follows:

1. For each chromosome, convert the chromosome's peak associations into an array.
2. For each chromosome, convert the chromosome's array into a binary max-heap keyed on the $-\log_{10}(P)$ of each genotype/phenotype pair, using the $O(n)$ running time heapify algorithm described in ref. ⁵⁹.
3. For each chromosome, while the chromosome's heap is not empty, extract the maximum genetic variant/phenotype pair from the heap to create a new cluster. This can be done in $O(\log(n))$ time. This extracted pair is the lead association of the cluster. Then, the 0.25-cM cover is removed from the chromosome's heap, which is also an $O(\log(n))$ operation. The cluster is outputted (including the removed aspects).
4. Because there are, at most, n extractions and deletions, the total running time is $O(n \log(n))$. This is optimal as it is tight to the sorting lower bound of $O(n \log(n))$ and because sorting (as a relaxation of clustering) has a running time no worse than the original problem.

The open-source Peaks software implementing these methods is available on GitHub at <https://github.com/wnldchen/peaks>.

Reporting Summary. Further information on research design is available in the Nature Research Reporting Summary linked to this article.

Data availability

Our resource includes openly released summary statistics and results for a variety of GWAS paradigms on the most recent release of 3,935 UKB brain imaging phenotypes. These results are released on BIG40 (<https://open.win.ox.ac.uk/ukbiobank/big40/>), the European Bioinformatics Institute and the Supplementary Material of this paper. An enumeration of the aspects of our resource is as follows:

- Summary statistics for our discovery cohort, available on BIG40 (Manhattan plots, full downloads and a browsable interface) and European Bioinformatics Institute under study accession numbers GCST90002426–GCST90006360 (ftp://ftp.ebi.ac.uk/pub/databases/gwas/summary_statistics/GCST90002426 to ftp://ftp.ebi.ac.uk/pub/databases/gwas/summary_statistics/GCST90006360)
- Details of the clusters of associations identified by Peaks, including summary statistics for replication, available on BIG40 and in Supplementary Tables 2–6
- Causal variants identified by CAVIAR for the four X chromosome clusters significant at the $-\log_{10}(P) \geq 11.1$ level (Supplementary Table 5)

- A full GWAS on all phenotypes and chromosomes, with the discovery and replication cohorts combined (available on the BIG40 website as a download and as a browsable interface)
- Sex-specific GWAS on the discovery cohort with genetic females and genetic males considered separately and combined through a Fisher meta-analysis (available as a download on BIG40)
- The heritability of each phenotype, assessed through LDSC on the full GWAS with discovery and replication cohorts combined (Supplementary Table 1)

Code availability

The following software packages and servers were used throughout this work:

- bgenie v1.3, software for efficient GWASs on high-dimensional phenotype data: <https://jmarchini.org/bgenie/>⁶³
- qctool v1.4 and v2.0.1, software for pre-processing genetic data: https://www.well.ox.ac.uk/~gav/qctool_v1/ and https://www.well.ox.ac.uk/~gav/qctool_v2/
- Peaks v1.0, novel software for extracting clusters from multi-phenotype GWASs: <https://github.com/wnldchen/peaks>
- PheWeb v1.1.19, a web server for browsing phenome-wide associations: <https://github.com/statgen/pheweb>⁶⁰
- BIG40 open web server for Brain Imaging Genetics: <https://open.win.ox.ac.uk/ukbiobank/big40>
- plink2 v2.0, alpha software for conducting GWASs and pre-processing of genetic data: <https://www.cog-genomics.org/plink/2.0/>⁵⁶
- CAVIAR v2.0, fine mapping software for extracting causal variants from summary statistics: <http://genetics.cs.ucla.edu/caviar/>³
- Open Targets Platform, an online web server for GWASs: <https://genetics-app.netlify.app/>⁶¹
- LDSC v1.0.1, software for heritability analysis from summary statistics (linkage score regression): <https://github.com/bulik/ldsc/>⁷
- The GTEx online resource: <https://gtexportal.org/home/>⁶²

References

1. Alfaro-Almagro, F. et al. Image processing and quality control for the first 10,000 brain imaging datasets from UK Biobank. *Neuroimage* **166**, 400–424 (2018).
2. Miller, K. L. et al. Multimodal population brain imaging in the UK Biobank prospective epidemiological study. *Nat. Neurosci.* **19**, 1523–1536 (2016).
3. Bycroft, C. et al. The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).
4. Marchini, J. & Howie, B. Genotype imputation for genome-wide association studies. *Nat. Rev. Genet.* **11**, 499–511 (2010).
5. König, I. R., Loley, C., Erdmann, J. & Ziegler, A. How to include chromosome X in your genome-wide association study. *Genet. Epidemiol.* **38**, 97–103 (2014).
6. Chang, C. C. et al. Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience* **4**, 7 (2015).
7. The 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
8. Fisher, R. A. Questions and answers #14. *Am. Statistician* **2**, 30–33 (1948).
9. Suchanek, M. A. Elementary yet precise worst-case analysis of Floyd's heap-construction program. *Fundam. Inform.* **120**, 75–92 (2012).
10. Gagliano, S. A. et al. Exploring and visualizing large-scale genetic associations by using PheWeb. *Nat. Genet.* **52**, 550–552 (2020).
11. Carvalho-Silva, D. et al. Open Targets Platform: new developments and updates two years on. *Nucleic Acids Res.* **47**, 1056–1065 (2019).
12. GTEx Consortium. Genetic effects on gene expression across human tissues. *Nature* **550**, 204–213 (2017).

Acknowledgements

This research was conducted, in part, using the UK Biobank Resource under application number 8107. We are grateful to UK Biobank for making the data available and to all UK Biobank study participants, who generously donated their time to make this resource possible. S.M.S. and G.D. were supported by Wellcome Trust Strategic Award 098369/Z/12/Z and Wellcome Trust Collaborative Award 215573/Z/19/Z. L.T.E. and W.C. were funded by NSERC grants RGPIN/05484-2019 and DGEGR/00118-2019 and an NSERC Undergraduate Student Research Award. G.D. was supported by an MRC Career Development Fellowship (MR/K006673/1). The BIG40 Open Data Server is provided by the Wellcome Centre for Integrative Neuroimaging, which is supported by center funding from the Wellcome Trust (203139/Z/16/Z). Compute resources were provided by the Oxford Biomedical Research Computing (BMRC) facility (a joint development between Oxford's Wellcome Centre for Human Genetics and the Big Data Institute, supported by Health Data Research UK and the NIHR Oxford Biomedical Research Centre). Some compute resources were also provided by ComputeCanada under the Resources for Research Groups program. We would like to thank the Resource Computing Managers at BMRC for their diligence in operation, with special thanks to J. Diprose, R. Esnouf and A. Huffman. We would also like to thank D. Mortimer, the Senior Informatics Officer at WIN FMRI, and M. Siegert, the Research Computing Director and Site Lead at WestGrid/ComputeCanada. We are grateful to S. Shi, A. Winkler, T. Nichols, P. McCarthy, D. Greve and B. Fischl for helpful input.

Author contributions

S.M.S. and L.T.E. co-directed this work. S.M.S., G.D., K.S. and L.T.E. interpreted the results and wrote the paper. All authors contributed to the analysis and the editing. W.C. wrote the Peaks software. T.H. created the BIG40 Pheweb resource. F.A. created novel brain imaging phenotypes for the UK Biobank.

Competing interests

The authors declare no competing financial interests.

Additional information

Extended data is available for this paper at <https://doi.org/10.1038/s41593-021-00826-4>.

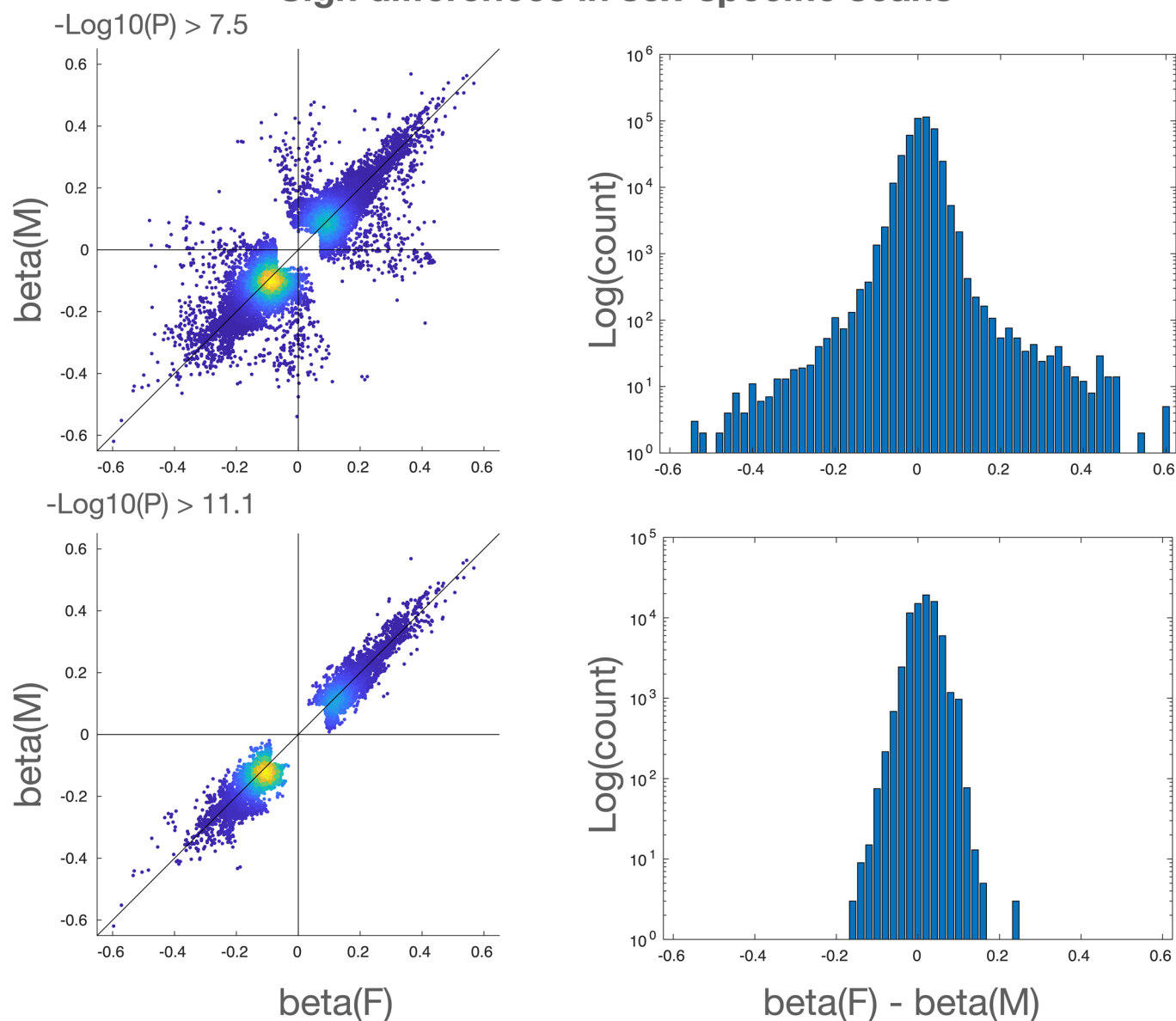
Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41593-021-00826-4>.

Correspondence and requests for materials should be addressed to L.T.E.

Peer review information *Nature Neuroscience* thanks Alex Fornito, Jason Stein, and the other, anonymous, reviewer(s) for their contribution to the peer review of this work.

Reprints and permissions information is available at www.nature.com/reprints.

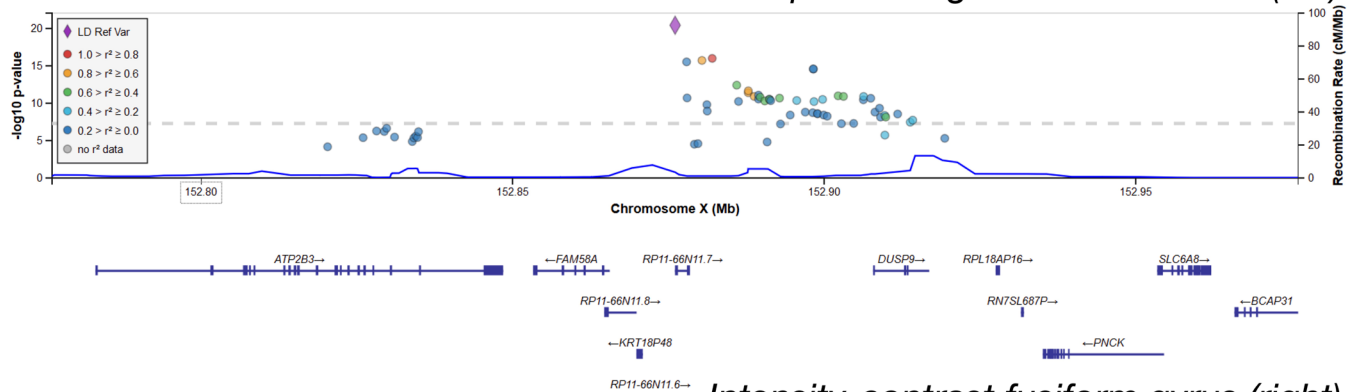
Sign differences in sex-specific scans



Extended Data Fig. 1 | Comparisons of effect sizes and signs for genetic females and males. *Top row:* Effect sizes for all associations with either genetic females or genetic males (or both) having $-\text{Log}_{10}(P) \geq 7.5$. *Bottom row:* effect sizes for all associations with either genetic females or genetic males (or both) having $-\text{Log}_{10}(P) \geq 11.1$. *Left column:* Scatter plots of effect sizes, indicating a small fraction (0.58%) of sign differences for $-\text{Log}_{10}(P) \geq 7.5$ and no sign differences (quadrants II and IV empty) for $-\text{Log}_{10}(P) \geq 11.1$ condition. *Right column:* Histograms of difference between effect sizes. Log y-scale indicates generally close matching of effect sizes.

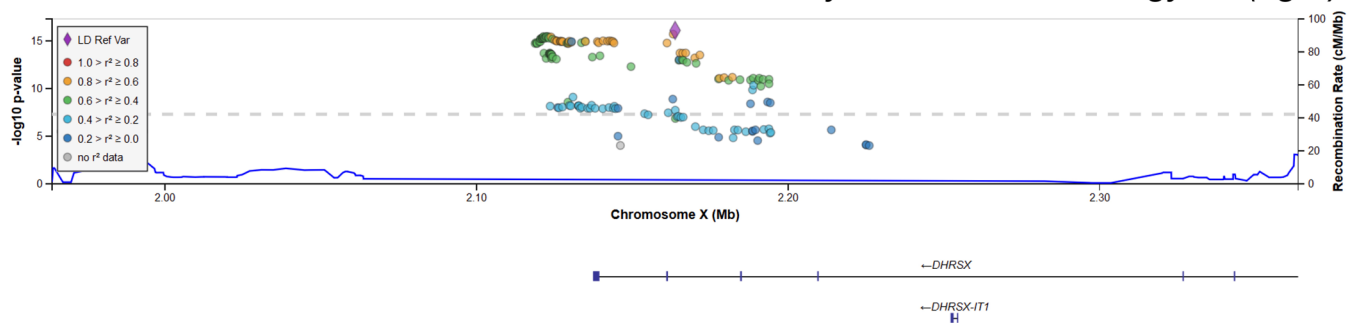
Category: WM tract ICVF

TBSS ICVF superior longitudinal fasciculus (left)



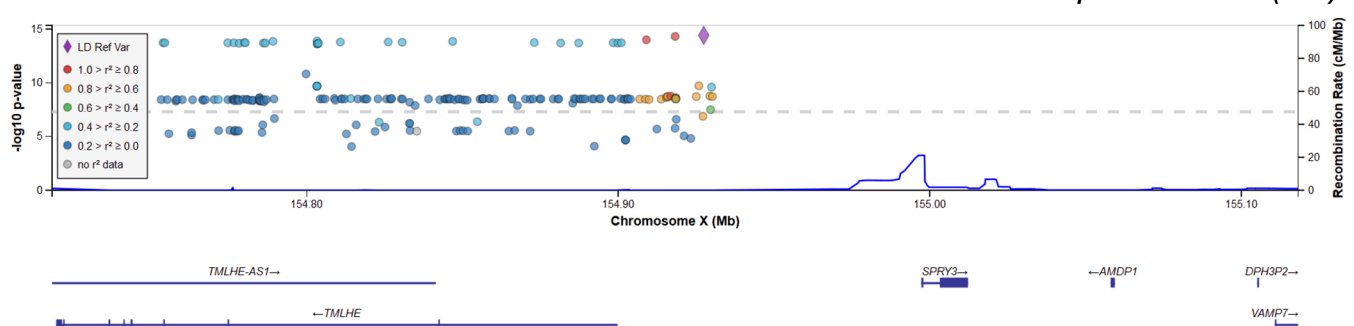
Category: cortical grey-white contrast

Intensity-contrast fusiform gyrus (right)



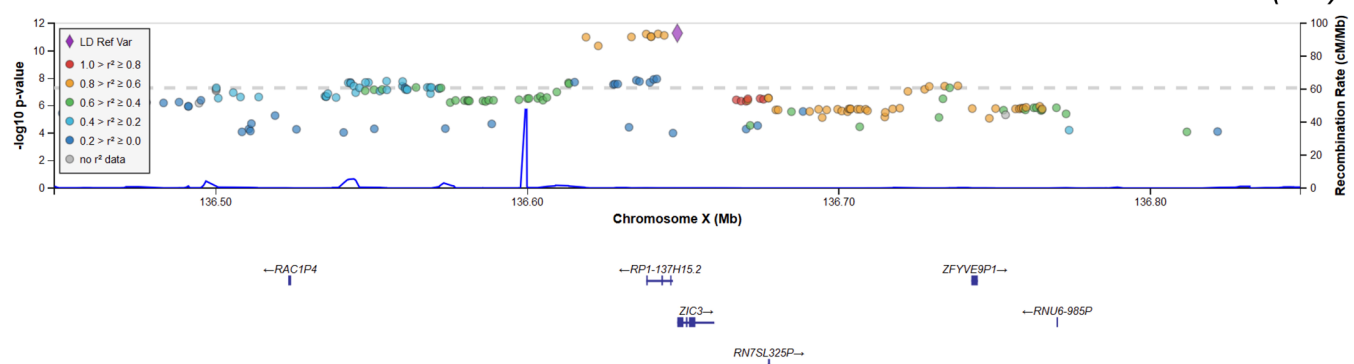
Category: cortical area

Area lateral occipital cortex (left)



Category: cortical area

Area rostral middle frontal cortex (left)



Extended Data Fig. 2 | Regional association plots of the significant variants in X. First row: Region around rs2272737 ($P = 3.5 \times 10^{-21}$). This variant is an eQTL of *FAM58A*. Second row: Region around rs62595479 ($P = 8.2 \times 10^{-17}$). This variant is located in a pseudoautosomal region (PAR1) of X, in an intron of *DHRX*. Third row: Region around rs644138 ($P = 4.8 \times 10^{-15}$). This variant is in an intron of *SPRY3* (and is an eQTL in brain tissue of various genes). Bottom row: Region around rs12843772 ($P = 5.1 \times 10^{-12}$) located ≤ 150 bp from *ZIC3*. The genomic positions of the loci and genes are based on Human Genome build hg19. Regions considered include all loci within 10 kbp of the hit.

Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

- | n/a | Confirmed |
|-------------------------------------|---|
| <input type="checkbox"/> | <input checked="" type="checkbox"/> The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> The statistical test(s) used AND whether they are one- or two-sided
<i>Only common tests should be described solely by name; describe more complex techniques in the Methods section.</i> |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> A description of all covariates tested |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals) |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> For null hypothesis testing, the test statistic (e.g. F , t , r) with confidence intervals, effect sizes, degrees of freedom and P value noted
<i>Give P values as exact values whenever suitable.</i> |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes |
| <input type="checkbox"/> | <input checked="" type="checkbox"/> Estimates of effect sizes (e.g. Cohen's d , Pearson's r), indicating how they were calculated |

Our web collection on [statistics for biologists](#) contains articles on many of the points above.

Software and code

Policy information about [availability of computer code](#)

Data collection No data collection software was used.

Data analysis

The following software packages and servers were used throughout this work:

- bgenie v1.3 software for efficient genome-wide association study on high dimensional phenotype data. <https://jmarchini.org/bgenie/> [Bycroft et al., 2018].
- qctool v1.4 and v2.0.1 software for preprocessing genetic data. https://www.well.ox.ac.uk/~gav/qctool_v1/ and https://www.well.ox.ac.uk/~gav/qctool_v2/
- Peaks v1.0 novel software for extracting clusters from multi-phenotype genome-wide association studies. <https://github.com/wnfldchen/peaks>
- PheWeb v1.1.19 a webserver for browsing phenome-wide associations. <https://github.com/statgen/pheweb> [Gagliano Taliun et al., 2020]
- Plink2 v2.0.0 general GWAS software. <https://www.cog-genomics.org/plink/2.0/>
- CAVIAR v2.2 causal variants software <http://genetics.cs.ucla.edu/caviar/>

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.

Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

The full set of GWAS results from this study is available on the Oxford BIG40 web browser (<https://open.win.ox.ac.uk/ukbiobank/big40/>), which allows users to browse associations by SNP, gene or phenotype. GWAS summary statistics can be downloaded for each phenotype from the same server.

Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

☒ Life sciences ☐ Behavioural & social sciences ☐ Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see nature.com/documents/nr-reporting-summary-flat.pdf

Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	We used brain IDPs from the "40k" (approximately 40,000 participants) data release in early 2020, as processed by WIN-FMRIB on behalf of UK Biobank [Alfaro-Almagro et al., 2018]. After removal of subjects as part of the genetic processing (see below), we used data from 33,224 subjects. These were then randomly split into a discovery sample of 22,138 subjects (11,624 genetic females) and a replication sample of 11,086 subjects (5,787 genetic females). We therefore used the maximum sample size available from this data source. This is by definition "sufficient" given that we are reporting associations that are statistically significant given this maximally available dataset.
Data exclusions	We consider the 488,377 samples included in the Spring 2018 release of the UK Biobank, and proceed with a preprocessing and discovery/reproduction paradigm similar to that described in Elliott et al. [2018]. Of the samples, 39,944 are included among the 41,016 samples in the UK Biobank for which IDPs are available, after the genotyping quality control procedures for sample removal specified in Bycroft et al. [2018]. We remove samples without recent UK ancestry as determined by the in.white.British.ancestry.subset variable in the file ukb_sqc_v2.txt provided in the meta-data for UK Biobank. This variable selects samples based on self reported ancestry and genetic principle component thresholds. We also remove subjects based on relatedness, forming a maximal unrelated subset using the procedures recommended in Bycroft et al. [2018]. This results in a maximally unrelated subset of 34,298 samples with recent UK ancestry and accepted genotyping and imaging quality control. We divide this set into a discovery and reproduction cohort according to a random 2/3 and 1/3 (respective) proportion, resulting in a discovery cohort of 22,865 samples and a reproduction cohort of 11,433 samples.
Replication	Full replication analysis carried out once, as described in full detail in Methods. Success of replication analysis fully described in main text.
Randomization	There is nothing in this study that pertains to randomization. We are using existing data released by UK Biobank.
Blinding	There is nothing in this study that pertains to blinding. We are using existing data released by UK Biobank.

Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

Materials & experimental systems

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> Antibodies
<input checked="" type="checkbox"/>	<input type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Palaeontology and archaeology
<input checked="" type="checkbox"/>	<input type="checkbox"/> Animals and other organisms
<input type="checkbox"/>	<input checked="" type="checkbox"/> Human research participants
<input checked="" type="checkbox"/>	<input type="checkbox"/> Clinical data
<input checked="" type="checkbox"/>	<input type="checkbox"/> Dual use research of concern

Methods

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input checked="" type="checkbox"/>	<input type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> MRI-based neuroimaging

Human research participants

Policy information about [studies involving human research participants](#)

Population characteristics

We used brain phenotypes from the ``40k'' (approximately 40,000 participants) UK Biobank data release in early 2020, as processed by WIN-FMRIB on behalf of UKB. After removal of subjects as part of the genetic processing (see Methods), we used data from 33,224 subjects. These were then randomly split into a discovery sample of 22,138 subjects (11,624 genetic females) and a replication sample of 11,086 subjects (5,787 genetic females). The ages in the discovery sample were: Females: mean age = 63.6 ± 7.3 years, min = 45.1, max = 81.8. Males: mean = 65.0 ± 7.6 , min = 46.1, max = 81.8. In the replication sample: Females: mean = 63.7 ± 7.4 , min = 46.3, max = 81.6. Males mean = 65.0 ± 7.6 , min = 46.1, max = 81.0. The exact numbers of subjects vary across IDPs, according to patterns of missing data, with the maximum numbers given above (for IDPs with no missing data), and the minimum numbers being just 16% lower. The BIG40 online table listing the IDPs includes the exact number of subjects (in discovery, replication samples and in the sex-specific GWAS) for each IDP. The details for these IDPs (including long descriptions, category names and units) are summarized in Supplementary Table S1.

Recruitment

We used existing Open data from UK Biobank and were not involved in recruitment. Recruitment was through the UK National Health Service, as described in:
Sudlow, C. et al. UK Biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. PLoS Med. 12, e1001779 (2015).

Ethics oversight

The UK Biobank has approval from the North West Multi-centre Research Ethics Committee (MREC) to obtain and disseminate data and samples from the participants (<http://www.ukbiobank.ac.uk/ethics/>), and these ethical regulations cover the work in this study. Written informed consent was obtained from all participants.

Note that full information on the approval of the study protocol must also be provided in the manuscript.