

# 1 **Genetics and Pathway Analysis of Normative Cognitive Variation** 2 **in the Philadelphia Neurodevelopmental Cohort**

3

4 Authors:

5 Shraddha Pai<sup>1</sup>, Shirley Hui<sup>1</sup>, Philipp Weber<sup>2</sup>, Owen Whitley<sup>1,3</sup>, Peipei Li<sup>4,5</sup>, Viviane Labrie<sup>4,5</sup>, Jan  
6 Baumbach<sup>2,6</sup>, Gary D Bader<sup>1,3,7,8</sup>

7

8 Affiliations:

- 9 1. The Donnelly Centre, University of Toronto, Toronto, Canada
- 10 2. Department of Mathematics and Computer Science, University of Southern Denmark, Odense,  
11 Denmark
- 12 3. Department of Molecular Genetics, University of Toronto, Toronto, Canada
- 13 4. Center for Neurodegenerative Science, Van Andel Research Institute, Grand Rapids, MI, USA
- 14 5. Division of Psychiatry and Behavioral Medicine, College of Human Medicine, Michigan State  
15 University, Grand Rapids, MI, USA
- 16 6. TUM School of Life Sciences Weihenstephan, Technical University of Munich, Munich, Germany.
- 17 7. Department of Computer Science, University of Toronto, Toronto, Canada
- 18 8. The Lunenfeld-Tanenbaum Research Institute, Mount Sinai Hospital, Toronto, Canada

19

20

21

22

\* [gary.bader@utoronto.ca](mailto:gary.bader@utoronto.ca)

23

24

## 25 Abstract

26

27 Identifying genes and cellular pathways associated with normative brain physiology and behavior  
28 could help discover molecular therapies that target specific psychiatric symptoms with minimal side  
29 effects. Linking genotype-phenotype associations from population-scale datasets to brain function is  
30 challenging because of the multi-level, heterogeneous nature of brain organization. To address this  
31 challenge, we developed a novel brain-focused gene and pathway prioritization workflow, which maps  
32 variants to genes based on knowledge of brain genome regulation, and subsequently to pathways,  
33 cells, diseases and drugs (21 resources). We applied this workflow to nine cognitive tasks from the  
34 Philadelphia Neurodevelopmental Cohort (subset of 3,319 individuals aged 8-21 years). We report  
35 genome-wide significance of variants associated with nonverbal reasoning within the 3' end of the  
36 *FBLN1* gene ( $p=4.6 \times 10^{-8}$ ), itself linked to fetal neurodevelopment and psychotic disorders. These  
37 findings suggest that nonverbal reasoning and *FBLN1* variation warrant further investigation in  
38 studies of psychosis. Multiple cognitive tasks demonstrated significant enrichment of variants in  
39 cellular pathways and brain-related gene sets, such as organ development, cell proliferation and  
40 nervous system dysfunction. Top-ranking genes in working memory associated pathways are  
41 genetically associated with multiple diseases with working memory deficits, including schizophrenia  
42 and Parkinson's disease, and with multiple drugs, suggesting that choice of therapy for memory  
43 deficits should consider disease context. Given the large amount of additional biological insight  
44 derived from our pathway analysis, versus a standard gene-based approach, we propose that “genes to  
45 behaviour” frameworks for modeling brain-related phenotypes, like RDoC, should include pathway  
46 information to create a “genes to pathways to behaviour” approach. Our workflow is broadly useful to  
47 put genotype-phenotype associations of brain-related phenotypes into the context of brain  
48 organization, function, disease and known molecular therapies.

## 49 Introduction

50 A major drive in the field of psychiatry is the reconceptualization of mental illnesses, diseases  
51 traditionally classified on the basis of clinical descriptions, as brain disorders treatable by  
52 neurobiologically-grounded therapies. The U.S. National Institute of Mental Health developed  
53 Research Domain Criteria (or RDoC), a framework to build a “genes to behaviour” model of the human  
54 brain that deconstructs behaviour into multiple domains mediated by different neuroanatomical  
55 regions, local cellular circuits and molecules ([https://www.nimh.nih.gov/research-  
56 priorities/rdoc/index.shtml](https://www.nimh.nih.gov/research-priorities/rdoc/index.shtml) ; <sup>1</sup>). The ambition of the RDoC framework is to develop neurobiologically-  
57 grounded taxonomy, biomarkers and treatments for mental illness to support use of precision  
58 medicine in psychiatry<sup>2</sup>. The neurobehavioural framework has been embraced by psychiatric  
59 researchers at multiple levels of brain research, including cross-disorder genome-wide association  
60 studies and genetic risk prediction<sup>3,4</sup>, and identification of genetic contributors to neuroimaging-based  
61 measures of brain activity and structure associated with disease<sup>5</sup>. Population-scale datasets that  
62 measure genotype and cognition-related phenotypes, such as the Philadelphia Neurodevelopmental  
63 Cohort<sup>6,7</sup>, the Adolescent Brain Cognitive Development (ABCD) dataset<sup>8</sup> and UK Biobank<sup>9</sup>, provide an  
64 attractive resource to build a molecules-to-behaviour model for brain-specific phenotypes. Moreover,  
65 maps of brain-specific genome regulation, such as those generated by the GTEx<sup>10</sup>, NIH Roadmap  
66 Epigenomics<sup>11</sup> and PsychENCODE<sup>12</sup> projects, now enable the effect of genetic variants to be  
67 interpreted in brain-relevant neuroanatomical and developmental contexts. However, integrating  
68 genotype-phenotype associations with these data resources to methodically infer variant impact on  
69 various levels of brain organization represents a major challenge, due to the large number of complex  
70 data sets that need to be integrated.

71  
72 In this work, we develop a novel brain-focused computational analysis workflow to identify genes,  
73 pathways and cellular functions, as well as gene-related brain functions, diseases and drugs. We apply  
74 this workflow to identify genes and functions associated with normative variation in nine cognitive  
75 phenotypes from the Philadelphia Neurodevelopmental Cohort (PNC). To our knowledge, little is  
76 known about the molecular basis of different cognitive phenotypes in humans, and the extent to which  
77 molecular and cellular players overlap across these. With extensive neurobehavioural and genotyping  
78 data available on 8,000 community youths aged 8-21 years, the PNC represents the largest publicly-  
79 available dataset of its kind for genotype-phenotype analysis of cognition<sup>6,7</sup>. All participants have  
80 computerized neurocognitive test battery (CNB) scores which measures speed and accuracy in  
81 multiple cognitive domains (e.g. emotion processing, executive function), and which has  
82 neurobehavioural validity (i.e. tasks known to activate specific brain regions), SNP-based  
83 heritability<sup>13</sup>, and disease relevance<sup>3,14</sup>. The CNB has also been characterized for demographic effects<sup>15</sup>  
84 and neuropsychological validation<sup>16</sup>, altogether providing a well-characterized set of phenotypes to  
85 study the genetic basis of specific cognitive abilities. While a number of CNB phenotypes demonstrate  
86 significant SNP-based heritability<sup>13</sup>, and reduced test scores have been genetically associated with  
87 psychiatric disease risk<sup>3</sup>, there has not been a methodical examination of the molecular players  
88 involved in individual phenotypes. We reasoned that identifying the genes, pathways, cellular and  
89 developmental context associated with these phenotypes could pinpoint genetic crosstalk between  
90 individual cognitive tasks and psychiatric and neurological diseases, and provide hypotheses for  
91 molecular therapy of corresponding cognitive impairments in disease.

## 92 Methods

### 93 Genetic imputation

94 The workflow for genomic imputation is shown in Supplementary Figure 1. Genotypes for four  
95 microarray genotyping platforms were downloaded from dbGaP (phs000607.v1). We performed  
96 genetic imputation for the Illumina Human610-Quad BeadChip, the Illumina HumanHap550  
97 Genotyping BeadChip v1.1, Illumina HumanHap550 Genotyping BeadChip v3, and the Affymetrix  
98 AxiomExpress platform (Supplementary Table 1, total of 6,502 samples before imputation), using the  
99 protocol recommended by the EMERGE consortium<sup>17</sup>. Imputation was performed as follows:

L00 **Step 1: Platform-specific plink Quality Control:** Quality control was first performed for each  
L01 microarray platform. Single nucleotide polymorphisms (SNPs) were limited to those on chr1-22. SNPs  
L02 in linkage disequilibrium (LD) were excluded (--indep-pairwise 50 5 0.2), and alleles were recoded  
L03 from numeric to letter (ACGT) coding. Samples were excluded if they demonstrated heterozygosity > 3  
L04 standard deviations (SD) from the mean, or if they were missing >=5% genotypes. Where samples had  
L05 pairwise Identity by Descent (IBD) > 0.185, one of the pair was excluded. Variants with minor allele  
L06 frequency (MAF) < 0.05 were excluded, as were those failing Hardy-Weinberg equilibrium with  $p < 1e-$   
L07  $6$  and those missing in >=5% samples.

L08 **Step 2: Convert coordinates to hg19.** LiftOver<sup>18</sup> was used to convert SNPs from hg18 to hg19;  
L09 Hap550K v1 data was in hg17 and was converted from this build to hg19.

L10 **Step 3: Strand-match check and prephasing:** ShapeIt v2.r790<sup>19</sup> was used to confirm that the allelic  
L11 strand in the input data matched that in the reference panel; where it did not, allele strands were  
L12 flipped (shapeit "--check" flag). ShapeIt was used to prephase the variants using the genetic\_b37  
L13 reference panel (downloaded from the Shapeit website,  
L14 [http://www.shapeit.fr/files/genetic\\_map\\_b37.tar.gz](http://www.shapeit.fr/files/genetic_map_b37.tar.gz))

L15 **Step 4: Imputation:** Genotypes were imputed using Impute2 v2.3.2<sup>20</sup> and a reference panel from the  
L16 1,000 Genomes (phase 1, prephased with Shapeit2, no singletons, 16 June 2014 release, downloaded  
L17 from

L18 [https://mathgen.stats.ox.ac.uk/impute/data\\_download\\_1000G\\_phase1\\_integrated\\_SHAPEIT2\\_16-06-](https://mathgen.stats.ox.ac.uk/impute/data_download_1000G_phase1_integrated_SHAPEIT2_16-06-14.html)  
L19 [14.html](https://mathgen.stats.ox.ac.uk/impute/data_download_1000G_phase1_integrated_SHAPEIT2_16-06-14.html)) was used for imputation, using the parameter settings "--use\_prephased\_g -Ne 20000 -seed  
L20 367946". Average concordance for all chromosomes was ~95%, indicating successful imputation  
L21 (Supplementary Figure 2). Imputed genotypes were merged across all platforms using software from  
L22 the Ritchie lab<sup>17</sup> (impute2-group-join.py, from <https://ritchielab.org/software/imputation-download>)  
L23 and converted to plink format. Following previous PNC genotype analysis<sup>13</sup>, only SNPs with info score  
L24 > 0.6 were retained, and deletions/insertions were excluded (plink "--snps-only just-acgt" flags). As  
L25 preliminary quality control, when merging across chromosomes, samples with missingness exceeding  
L26 99% were excluded, as were SNPs with MAF < 1% and with missingness exceeding 99%. This step  
L27 resulted in 10,845,339 SNPs and 6,327 individuals.

L28 **Step 5: Post-imputation Quality Control:** The HapMap3 panel was used to assign genetic ancestry  
L29 for samples, using steps from <sup>21</sup> (Supplementary Figure 3). Individuals within 5 SD of the centroid of  
L30 the HapMap3 CEU (Utah residents with Northern or Western European ancestry) or TSI (Tuscans in  
L31 Italy) cluster were assigned to belong to the respective groups, and were classified as being of  
L32 European descent; 3,441 individuals pass this filter. Individuals with >5% missing data were excluded,  
L33 as was one of each pair of individuals with IBS > 0.185 (47 individuals); 3,394 individuals passed this  
L34 filter. Variants that were symmetric or in regions of high LD (Supplementary Table 2) were excluded  
L35 (9,631,316 SNPs passed). Variants with >5% missingness were excluded (1,569,407 SNPs excluded).  
L36 Finally, SNPs with MAF < 0.01 (3,168,339 SNPs) and failing Hardy-Weinberg equilibrium (HWE) with  
L37  $p$  value <  $1e-6$  (373 SNPs) were excluded, resulting in 4,893,197 SNPs. Unlike Verma et al, quality  
L38 control steps were performed once, rather than repeated after samples were excluded. In sum, the

L39 imputation process resulted in 3,394 individuals and 4,893,197 SNPs available for downstream  
L40 analysis.

## L41 **Phenotype processing**

L42 Phenotype data was downloaded from dbGaP for 8,719 individuals. 637 individuals with severe  
L43 medical conditions (Medical rating=4) were excluded to avoid confounding the symptoms of their  
L44 conditions with performance on the cognitive tests<sup>13</sup>. Linear regression was used to regress out the  
L45 effect of age at test time (variable name: "age at cnb") and sex from sample-level phenotype scores,  
L46 and the residualized phenotype was used for downstream analysis.

L47 The nine phenotypes selected for pathway analysis were measures of overall performance accuracy in  
L48 Penn Computerized Neurocognitive Test Battery (CNB; Supplementary Table 3) and represented  
L49 major cognitive domains. Following regression, none of the phenotypes were significantly correlated  
L50 with age after Bonferroni correction, indicating that the age effect had been reduced (Supplementary  
L51 Table 4). Following guidelines from previous analyses on these data<sup>3</sup>, individuals with scores more  
L52 than four standard deviations from the mean for a particular test were excluded from the analysis of  
L53 the corresponding phenotype. For a given phenotype, only samples with a code indicating a valid test  
L54 score (codes "V" or "V2") were included; e.g. for pfmt\_tp (Penn Face Memory Test), only samples with  
L55 pfmt\_valid = "V" or "V2" were retained; the rest had scores set to NA. Finally, each phenotype was  
L56 dichotomized so that samples in the bottom 33<sup>rd</sup> percentile were relabeled as "poor" performers and  
L57 those in the top 33<sup>rd</sup> were set to be "good" performers; for a given phenotype, this process resulted in  
L58 ~1,000 samples in each group (Supplementary Table 3). Where an individual had good or poor  
L59 performance in multiple phenotypes, they were included in the corresponding group for each of those  
L60 phenotypes.

L61

## L62 **Genetic association analysis**

L63 For each of 9 CNB phenotypes, marginal SNP-level association was calculated using a mixed-effects  
L64 linear model (MLMA), using the leave-one-chromosome-out (LOCO) method of estimating polygenic  
L65 contribution (GCTA v1.97.7beta software<sup>22</sup>). In this strategy, a mixed-effect model is fit for each SNP:

$$L66 \quad y = a + bx + g + e$$

L67

L68 In this model,  $y$  is the binarized label (good/poor performers on a particular task),  $x$  measures the  
L69 effect of genotype (indicator variable encoded as 0, 1 or 2), and  $g$  represents the polygenic  
L70 contribution of all the SNPs in the genome (here, the ~4.89M imputed SNPs). In the LOCO variation,  $g$ -  
L71 is calculated using a chromosome-specific genetic relatedness matrix, one that excludes the  
L72 chromosome on which the candidate SNP is located<sup>22</sup>. SNPs and associated genes were annotated as  
L73 described in Supplementary Notes 1-4.

L74

## L75 **Hi-C Data Processing**

L76 We generated Hi-C data from the human prefrontal cortex<sup>23</sup> (Illumina HiSeq 2000 paired-end raw  
L77 sequence reads; n=1 sample; 746 Million reads; accession: GSM2322542  
L78 [<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSM2322542>]). Hi-C analysis involved Trim  
L79 Galore<sup>24</sup> (v0.4.3) for adapter trimming, HICUP<sup>25</sup> (v0.5.9) for mapping and performing quality control,  
L80 and GOTHIC<sup>26</sup> for identifying significant interactions (Bonferroni  $p < 0.05$ ), with a 40 kb resolution. Hi-  
L81 C gene annotation involved identifying interactions with gene promoters, defined as  $\pm 2$  kb of a gene  
L82 TSS. This analysis identified 303,464 interactions used for our study.

L83

## SNP to gene mapping for annotation and enrichment analyses

SNPs were mapped to genes using a combination of positional information, brain-specific expression Quantitative Trait Locus (eQTL) and higher-order chromatin interaction (hi-C) information. For eQTL-based mapping, we limited the search to significant eQTLs in brain tissue (GTEx v7 brain anterior cingulate cortex BA24, brain cortex, brain frontal cortex BA9, and hippocampus; downloaded from <https://www.gtexportal.org>; Supplementary Note 1<sup>10</sup>); of these, only SNPs located in open chromatin regions of brain-related samples were included (Roadmap Epigenomics 15-core chromatin state  $\leq 7$ )<sup>11</sup>. These included maps derived from neurospheres, angular gyrus, anterior caudate, germinal matrix, hippocampus, inferior temporal lobe, dorsolateral prefrontal cortex, substantia nigra, and fetal brain of both sexes (samples E053, E054, E067, E068, E069, E070, E071, E072, E073, E074, E081, E082, and E125), downloaded from <http://www.roadmapepigenomics.org/>. For 3D chromatin interaction mapping, SNPs were mapped to genes if these were located within a region where higher-order interaction was ascertained in the dorsolateral prefrontal cortex<sup>27</sup>; this region was constrained to be 250bp upstream and 500bp downstream of the gene's transcription start site; of these SNPs, only those overlapping brain enhancers were included<sup>11</sup>. These included enhancers in angular gyrus, hippocampus, inferior temporal lobe, and dorsolateral prefrontal cortex (samples E067, E071, E072, and E073; chromatin state "Enh" or "EnhG"). Finally, SNPs were positionally mapped to the nearest gene if the shortest distance to either transcription start site or end site was 60kb. This cutoff was selected because it maps the majority (90%) of SNPs to their nearest gene.

The order of mapping was as follows: SNPs that mapped to a gene via brain eQTL or hi-C interactions were not also positionally mapped to a gene. A SNP was allowed to map to genes using both eQTL and hi-C, and where SNPs mapped to multiple genes all associations were retained. SNPs without eQTL or hi-C mappings were positionally mapped to a gene. Where a SNP mapped to multiple genes, all associations were retained. These SNP-gene mappings were used for the gene set enrichment analysis described below, as well as to annotate SNPs from the GWAS analysis.

## Gene set enrichment analysis

For each of the nine CNB phenotypes, gene set enrichment analysis was performed using an implementation of GSEA for genetic variants<sup>28,29</sup>. GSEA was selected as it computes pathway enrichment scores using all available SNP information, which improves sensitivity, rather than using a hypergeometric model limited to SNPs passing a specific GWAS p-value cutoff. All SNPs were mapped to genes (as described in the "SNP-gene mapping for annotation and enrichment analyses" section) and each gene score is the best GWAS marginal p-value of all mapped SNPs. For each pathway, GSEA computes an enrichment score (ES) using the rank-sum of gene scores. The set of genes that appear in the ranked list before the rank-sum reaches its maximum deviation from zero, is called the "leading edge subset", and is interpreted as the core set of genes responsible for the pathway's enrichment signal. Following computation of the ES, a null distribution is created for each pathway by repeating genome-wide association tests with randomly label-permuted data and by computing ES from these permuted data; in this work, we use 100 permutations. Finally, the ES on the original data is normalized to the score computed for the same gene set for label-permuted data (Z-score of real ES relative to distribution of ES in label-permuted data), resulting in a Normalized Enrichment Score (NES) per pathway. The nominal p-value for the NES score is computed based on the null distribution and FDR correction is used to generate a q-value.

The first enrichment analysis used pathway information compiled from HumanCyc<sup>30</sup> (<http://humancyc.org>), NetPath (<http://www.netpath.org>)<sup>31</sup>, Reactome (<http://www.reactome.org>)<sup>32</sup>, NCI Curated Pathways<sup>33</sup>, mSigDB<sup>34</sup> (<http://software.broadinstitute.org/gsea/msigdb/>), and Panther<sup>35</sup> (<http://pantherdb.org/>) and Gene Ontology<sup>36</sup>

(Human\_GOBP\_AllPathways\_no\_GO\_iea\_May\_01\_2018\_symbol.gmt, downloaded from [http://download.baderlab.org/EM\\_Genesets/May\\_01\\_2018/Human/symbol/Human\\_GOBP\\_AllPathways\\_no\\_GO\\_iea\\_May\\_01\\_2018\\_symbol.gmt](http://download.baderlab.org/EM_Genesets/May_01_2018/Human/symbol/Human_GOBP_AllPathways_no_GO_iea_May_01_2018_symbol.gmt)); only pathways with 20-500 genes were used.

The second enrichment analysis used brain-related gene sets we compiled from various literature sources (see Supplementary Table 5 and Supplementary Note 5 for details). Gene sets included those identified through transcriptomic or proteomic assays in human brain tissue (i.e. direct measurement of expression), and genes associated with brain function by indirect inference (e.g. genetic association of nervous system disorders); both groups of gene sets were combined for this enrichment analysis. The transcriptomic/proteomic gene sets included: genes identified as markers for adult and fetal brain cell types through single-cell transcriptomic experiments<sup>37-39</sup>, genes enriched for brain-specific expression (Human Protein Atlas project (<https://www.proteinatlas.org>)); genes co-expressed with markers of various stages of human brain development (BrainSpan<sup>41</sup>); and genes encoding proteins altered in the schizophrenia synaptosomal proteome<sup>42</sup>. Other gene sets included: genes associated with schizophrenia, bipolar disorder, autism spectrum disorder and major depressive disorder through large-scale genetic association studies by the Psychiatric Genomics Consortium (Supplementary Note 5); genes associated with nervous system disorders by the Human Phenotype Ontology<sup>43</sup>. Genes in the second group were filtered to only include genes with detectable expression (including long non-coding RNA genes) in the fetal<sup>44</sup> or adult human brain<sup>40</sup>. A total of 1,343 gene sets were collected. Only gene sets with 20-500 genes were included in the analysis; 421 gene sets met these criteria and were included in the enrichment analysis.

## Enrichment map

An enrichment map was created to visualize the functional themes significant in enrichment analyses. We used the EnrichmentMap app v3.1.0<sup>45</sup> and Cytoscape v3.7.1<sup>46</sup> to create the map. Nodes in the map are pathways with FDR significance of  $FDR < 0.10$  and edges in the map connect nodes with at least a gene set similarity of 0.375 (using Jaccard + Overlap similarity).

## Leading edge gene interaction network

Genes contributing to pathway enrichment results (leading edge genes) were obtained as part of the implementation of GSEA for genetic variants<sup>28</sup>. The network was constructed from leading edge genes of pathways with  $q < 0.05$ . The online GeneMANIA service (v 3.6.0; <https://genemania.org>)<sup>47</sup> was used to obtain a gene-gene interaction network for leading edge genes (human database, default settings); the resulting network and edge attributes were downloaded. This network was imported into Cytoscape v3.7.1. Known drug associations were obtained from DGIdb<sup>48</sup> and GWAS associations with nervous system disorders were obtained from the NHGRI-EBI GWAS catalogue, via programmatic search using the TargetValidation.org API<sup>49,50</sup>. Cell type marker information was compiled from single cell RNA-seq datasets, including those for adult and fetal human brain<sup>37-39</sup>.

## Results

Figure 1a shows the workflow for the analysis performed in this work. Briefly, genotypes were imputed using a reference panel from the 1,000 Genomes Project<sup>51</sup>, and samples were limited to those of European genetic ancestry (Supplementary Figure 1-3, Supplementary Table 1). 3,394 individuals and ~4.9M SNPs passed the quality control and imputation process. Following quality control of phenotype data, 3,116 European samples passed both genotype and phenotype filters and were included in downstream analyses. We selected nine phenotypes from the Penn Computerized Neurocognitive Test Battery (CNB) representing overall accuracy in four cognitive domains: complex

277 cognition, executive function, declarative memory, and social processing (Supplementary Table 3).  
278 Measures included performance for verbal reasoning, nonverbal reasoning, spatial reasoning,  
279 attention allocation, working memory, recall tests for faces, words and objects, and emotion  
280 identification<sup>14</sup>. In all instances, age and sex was regressed out of the phenotype (Supplementary Table  
281 4) and samples were thereafter binarized into poor and good performers (bottom and top 33%  
282 percentile, respectively) resulting in ~1,000 samples per group for each phenotype (Supplementary  
283 Figure 4,5, Supplementary Table 3).

284  
285 For each of the nine phenotypes, we first performed SNP-level genome-wide association analysis using  
286 a mixed-effects linear model that included genome-wide genetic ancestry as a covariate (GCTA<sup>22</sup>).  
287 Among the nine phenotypes, 661 SNPs had suggestive levels of significance at the genome-wide level  
288 ( $p < 10^{-5}$ ; Figure 1b,c, Supplementary Figure 6,7, Supplementary Table 6). Over half of these SNPs are  
289 associated with tasks related to complex cognition (377 SNPs or 57%); 27% were associated with  
290 executive function (177 SNPs), 13% with declarative memory tasks (83 SNPs), and 4% with emotion  
291 identification (24 SNPs).

292  
293 We mapped SNPs to genes using brain eQTL information, brain-specific higher-order chromatin  
294 interactions<sup>10,27</sup> and positional information. We integrated our findings with functional annotation  
295 maps of the brain to identify the neurodevelopmental and psychiatric significance of these genes  
296 (Figure 1d, Supplementary Table 7). The 661 suggestive peaks map to 106 genes. ~14% (15 genes)  
297 have been genetically associated with diseases of the nervous system, including schizophrenia  
298 (*SNAP91*, *CORO7*), bipolar disorder (*FBLN1*), multiple sclerosis (*THEMIS*, *CLECL16A*), alcohol  
299 dependence (*MREG*, *KCNJ6*, *FSTL5*), and Alzheimer's disease (*NRXN1*) (11 or 13% genes;  
300 Supplementary Table 7). Nearly one-third of these genes are markers of various cell-types in the fetal  
301 and newborn brain, including neuronal progenitor cells, neurons, radial glia, astrocytes, and  
302 endothelial cells (31 genes, 29%;<sup>39</sup>), and one gene is a marker of adult brain cells (*THEMIS*)<sup>37</sup>. Seven  
303 genes are known to interact with drugs; a notable interaction is between *CACNA2D3*, a voltage gated  
304 Calcium channel with suggestive association with working memory (top SNP  $p = 3.9e-6$ ), and  
305 Gabapentin enacarbil, a drug used to treat epilepsy, neuralgia and restless legs syndrome<sup>52</sup>. One-sixth  
306 of suggestive peaks (112 SNPs or 17%) were predicted to have a functional consequence in brain  
307 tissue (Figure 1c, e), including nonsynonymous changes to protein sequence, presence in brain-  
308 specific promoters and enhancers, or association with changes in gene expression. In summary,  
309 genetic variants associated with typical variation in neurocognition map to genes implicated in human  
310 brain development, altered in psychiatric disease, and that are modulated by drugs used to treat  
311 neurological conditions.

312  
313 Nonverbal reasoning was the only phenotype with SNPs passing the cutoff for genome-wide  
314 significance (rs77601382 and rs5765534,  $p = 4.6 \times 10^{-8}$ ) (Figure 2). The peak is located in a ~33kb  
315 region (chr22:45,977,415-46,008,175) overlapping the 3' end of the Fibulin-1 (*FBLN1*) gene, including  
316 the last intron and exon (Figure 2b). To better understand the significance of this gene in brain  
317 function, we examined *FBLN1* expression in published fetal and adult transcriptomes, and single-cell  
318 data<sup>10,39,41</sup>. *FBLN1* transcription in the human brain is highest in the early stages of fetal brain  
319 development, with little to no expression in the adult (Figure 2c, Supplementary Figure 8); this is  
320 consistent with single-cell assays showing *FBLN1* to be a marker for dividing progenitor cells in the  
321 fetal brain (Figure 1d,<sup>39</sup>). *FBLN1* encodes a glycoprotein present in the extracellular matrix; this  
322 protein is a direct interactor of proteins involved in neuronal diseases, such as Amyloid Precursor  
323 Protein-1 (Supplementary Figure 9<sup>53</sup>). *FBLN1* expression is upregulated in the brain in schizophrenia  
324 and has been previously associated with genetic risk for bipolar disorder (Figure 1d,<sup>54,55</sup>). Therefore,  
325 we conclude that *FBLN1*, associated with nonverbal reasoning test performance, shows characteristics



of a gene involved in neurodevelopment and the dysregulation of which could increase risk for psychotic disorders of neurodevelopmental origin.

We then performed pathway analysis for all nine selected CNB phenotypes using a rank-based pathway analysis strategy that includes all SNPs used in the association analysis (GSEA<sup>28,34</sup>, 100 permutations; 4,102 pathways tested). SNPs were mapped to genes using brain-specific eQTL, chromatin interaction and positional information, using the same method as described above. Four out of nine phenotypes demonstrated significant enrichment of top-ranking genetic variants in pathways ( $q < 0.1$ ; Figure 3a, Supplementary Tables 8-10). These included tasks in complex cognition (spatial reasoning), declarative memory (object and face memory), and executive function (working memory). The working memory phenotype showed significant enrichment of variants in pathways related to development, including neural development ( $q < 0.05$ ; Figure 3a, Supplementary Tables 8-10). To understand how genes contributing to pathway enrichment could be related to brain function, we annotated the corresponding leading edge genes with prior knowledge about associations with nervous system disorders, drug interactions and transcription in brain cell types<sup>37-39,48,49</sup>. Out of 355 leading edge genes, over half are known brain cell markers (228 genes or 64%), roughly one-third have known drug interactions (129 genes or 36%), and ~14% are associated with nervous system disease (51 genes) (pathway  $q < 0.10$ , Figure 3b, Supplementary Table 10). Among disease-associated genes were those associated with autism (*CSDE1*), multiple sclerosis (*CYP27B1*, *EOMES*), depression (*ROBO1*), glaucoma and wet macular degeneration (*LHCGR*). None of the SNPs associated with leading-edge genes (416 SNPs) overlapped suggestive or significant GWAS SNPs (661 SNPs).

To identify enrichment specific to brain-related processes and mental illness, we performed a second enrichment analysis using gene sets curated from the literature (Supplementary Note 5). These included gene sets derived from transcriptomic and proteomic profiles of the developing and adult healthy brain and brains affected by mental illness, genome-wide association studies and terms from phenotype ontology (421 gene sets tested, Supplementary Note 5, Supplementary Table 5, Supplementary Data 1). Six gene sets were significantly enriched ( $q < 0.10$ ), with five associated with working memory and the sixth with verbal reasoning (Figure 3c, Supplementary Table 11). A cluster of related gene sets related to autonomic nervous system dysfunction and a gene set related to locomotor dysfunction achieved significance at  $q < 0.05$ . Only one out of 157 SNPs associated with leading-edge genes overlaps with suggestive SNPs from GWAS analysis. Roughly 13% of the 134 leading edge genes are associated with nervous system disorders (18 genes), one-fifth have known drug targets (27 genes, 20%), and over half (81 genes or 60%) are markers of brain cell-types (Figure 3c,d; Supplementary Table 12, 13). Five genes have all three attributes: *SNCA*, *CAV1*, *LRRK2*, *ERBB4* and *MAPT* (Figure 3d, Supplementary Table 13). One example is Alpha-synuclein (*SNCA*, top SNP  $p = 2.6e-4$ ), which has been genetically associated with risk for developing Parkinson's disease<sup>56</sup>, is a marker of excitatory neurons in the fetal brain<sup>39</sup>, and is a drug target of BIIB504<sup>48</sup>. Another example is ERB-B2 receptor tyrosine kinase 4 (*ERBB4*), which has been genetically associated with mood disorders and unipolar depression<sup>57</sup>, is a target of 24 drugs and is a marker of inhibitory neurons in the fetal brain. Other leading edge genes have been associated with schizophrenia, autism spectrum disorder, Parkinson's disease, Alzheimer's disease, depression and mood disorders (Figure 3d, Supplementary Table 13). In summary, genetic variants associated with normative variation in a range of neurocognitive phenotypes are enriched in pathways and gene sets related to cell proliferation, brain development, nervous system dysfunction and mental disorders.

## Discussion

This study identifies molecular variants and cellular processes that contribute to normal human variation in specific cognitive domains. Consistent with heritability estimates, we find that the number of variant-level associations and enriched pathways varies considerably by phenotype (Figure 4). In particular, we find an enrichment of genetic variants associated with complex cognitive phenotypes (75-219 suggestive peaks), consistent with heritability estimates of up to 0.30-0.41 for these phenotypes<sup>13</sup>. A variety of cognitive phenotypes are enriched for variants in pathways. Moreover, the set of variants driving pathway enrichment has almost no overlap with suggestive variants from the GWAS analysis (no overlap for brain-related gene sets; a single SNP, rs9367669, overlaps for pathway sets). These results suggest that a molecules-to-behaviour research framework that includes genes and molecules, should also include pathways as a way to uncover new biological insights into existing genotype databases. Previous research in other polygenic psychiatric disorders, such as schizophrenia and major depression<sup>58</sup>, has also shown an enrichment of disease-associated molecules in pathways. We suggest that the Research Domain Criteria (RDoC) matrix be updated to add a level for pathways, above that of genes and molecules and below cells. This modification will help associate additional genetic signal with brain related phenotypes, which otherwise would be missed if just considering SNPs and genes.

Variants, genes and pathways associated with typical variation in neurocognitive phenotypes, demonstrate evidence for a role in neurodevelopment, modulating gene expression in the fetal and adult brain and increasing risk for psychiatric disease (Figure 1, Supplementary Table 6, 7, 10, 13). Multiple lines of evidence suggest that *FBLN1*, the gene associated with genome-wide significant SNPs for nonverbal reasoning, is dysregulated in disease. In addition to the evidence provided in our results (Figure 1d, Figure 2c, Supplementary Figure 8,9), *FBLN1* has been associated with other rare genetic syndromes and protein levels of *FBLN1* have been associated with altered risk for ischaemic stroke<sup>59,60</sup>. However, the mechanism by which *FBLN1* contributes to normal brain function is not known. We also do not exclude the possibility that suggestive peaks we identified within *FBLN1* may affect the function of neighbouring genes. One such gene is Ataxin-10 (*ATXN10*), in which a pentanucleotide repeat expansion causes spinocerebellar atrophy and ataxia<sup>61</sup>.

An advantage of using a rank-based gene set enrichment analysis method, as compared to hypergeometric tests, is that the method ranks and prioritizes a subset of genes (leading edge genes) within a potentially large gene set (>100 genes), which are responsible for driving the enrichment statistic. In this work, we found five neurocognitive phenotypes with significant enrichment of high-ranking variants in pathways. We annotated leading edge genes to identify those that are jointly related to working memory, which demonstrated significant enrichment in both gene set analyses, and psychiatric disease (Figure 3). For instance, among the leading edge genes contributing to working memory were genes previously associated with Parkinson's disease, Alzheimer's disease, schizophrenia, autism, and depression, all of which have been associated with working memory impairments<sup>62-67</sup>. We note, however, that the individual genes connecting any given disease to working memory are different. For instance, among leading edge genes for working memory, *ERBB4* is associated with depression, whereas *SNCA* is associated with Parkinson's disease (Figure 3c, Supplementary Table 13). One implication of this partially overlapping gene network is that the therapeutic targets that may be relevant for working memory deficits may depend on what disease the patient has, as a different subset of the "working memory gene network" is affected by each condition.

This work contributes towards an understanding of the molecular underpinnings of human brain-related behaviour and could help to identify genetic contributors towards the heterogeneity in phenotypes associated with multiple brain-related disorders<sup>68,69</sup>. Our analysis is limited to univariate

genetic effects, and future work will explore the contribution of interactions between individual SNPs, possibly explaining lack of SNP-level or pathway-level signal in some of the phenotypes studied here<sup>70</sup>. Our findings also suggest that different cognitive phenotypes may be vulnerable to genetic alterations in different cellular pathways. Such exploration could identify disease-specific molecular targets that impinge on the same neurocognitive phenotype. Finally, we propose that research frameworks for linking genotype to phenotype for brain-related traits include cellular pathways as an organizational layer to support uncovering additional genetic signal from available genetic data.

## Conflict of Interest

The authors declare no conflict of interest.

## Acknowledgements

We thank Shefali Verma for advice on genome imputation, and Sarah Gagliano for guidance on genetic analysis.

## References

- 1 Insel, T. *et al.* Research domain criteria (RDoC): toward a new classification framework for research on mental disorders. (2010).
- 2 Stein, M. B. & Smoller, J. W. Precision Psychiatry-Will Genomic Medicine Lead the Way? *JAMA psychiatry* **75**, 663-664, doi:10.1001/jamapsychiatry.2018.0375 (2018).
- 3 Germine, L. *et al.* Association between polygenic risk for schizophrenia, neurocognition and social cognition across development. *Translational psychiatry* **6**, e924, doi:10.1038/tp.2016.147 (2016).
- 4 Cross-Disorder Group of the Psychiatric Genomics Consortium. Identification of risk loci with shared effects on five major psychiatric disorders: a genome-wide analysis. *Lancet (London, England)* **381**, 1371-1379, doi:10.1016/s0140-6736(12)62129-1 (2013).
- 5 Thompson, P. M. *et al.* The ENIGMA Consortium: large-scale collaborative analyses of neuroimaging and genetic data. *Brain imaging and behavior* **8**, 153-182, doi:10.1007/s11682-013-9269-5 (2014).
- 6 Calkins, M. E. *et al.* The Philadelphia Neurodevelopmental Cohort: constructing a deep phenotyping collaborative. *Journal of child psychology and psychiatry, and allied disciplines* **56**, 1356-1369, doi:10.1111/jcpp.12416 (2015).
- 7 Satterthwaite, T. D. *et al.* Neuroimaging of the Philadelphia neurodevelopmental cohort. *NeuroImage* **86**, 544-553, doi:10.1016/j.neuroimage.2013.07.064 (2014).
- 8 Jernigan, T. L. & Brown, S. A. Introduction. *Developmental cognitive neuroscience* **32**, 1-3, doi:10.1016/j.dcn.2018.02.002 (2018).
- 9 Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203-209, doi:10.1038/s41586-018-0579-z (2018).
- 10 Battle, A., Brown, C. D., Engelhardt, B. E. & Montgomery, S. B. Genetic effects on gene expression across human tissues. *Nature* **550**, 204-213, doi:10.1038/nature24277 (2017).
- 11 Kundaje, A. *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317-330, doi:10.1038/nature14248 (2015).
- 12 Akbarian, S. *et al.* The PsychENCODE project. *Nature neuroscience* **18**, 1707-1712, doi:10.1038/nn.4156 (2015).

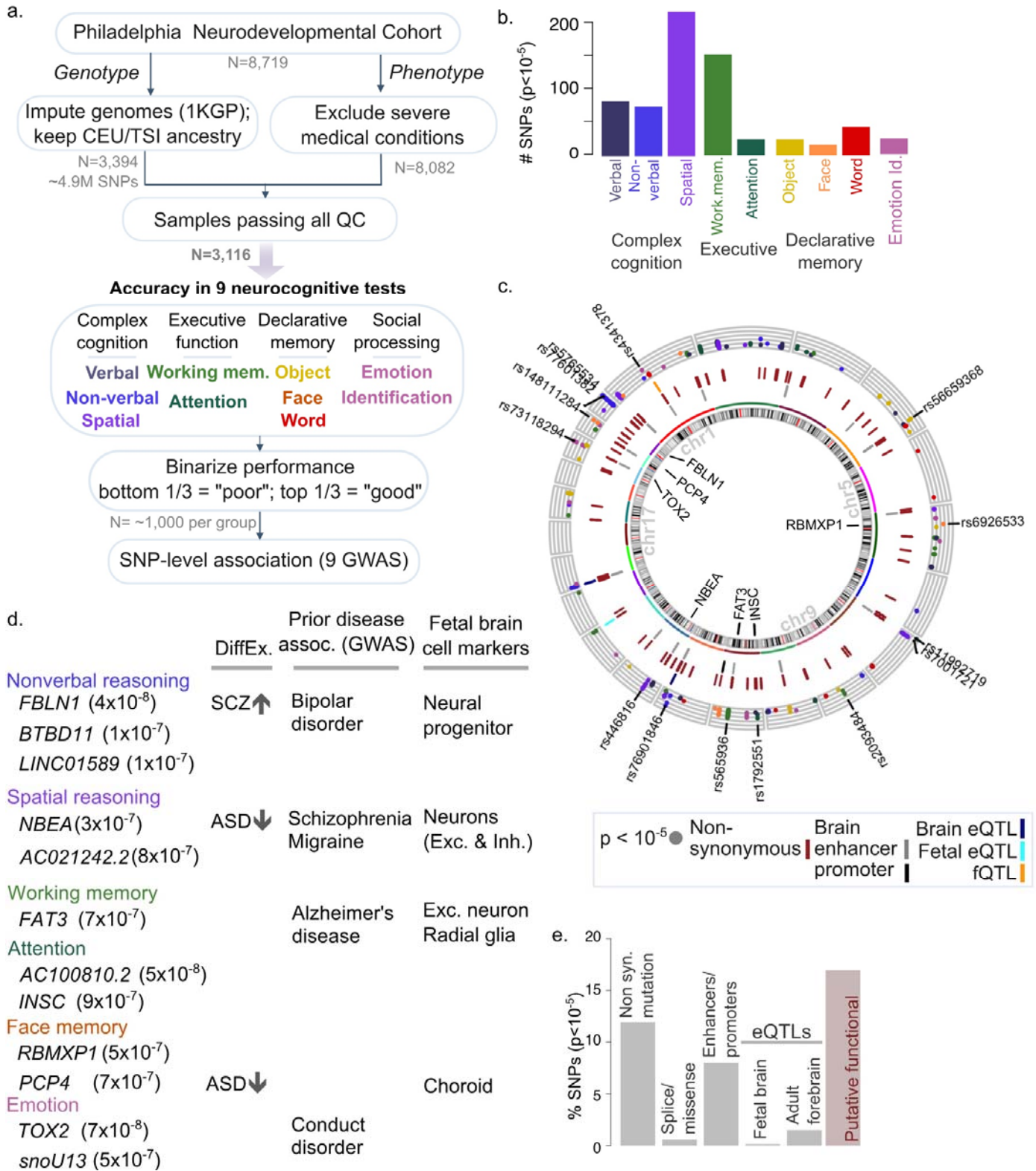
- 13 Robinson, E. B. *et al.* The genetic architecture of pediatric cognitive abilities in the Philadelphia  
14 Neurodevelopmental Cohort. *Molecular psychiatry* **20**, 454-458, doi:10.1038/mp.2014.65  
15 (2015).
- 16 Gur, R. C. *et al.* A cognitive neuroscience-based computerized battery for efficient measurement  
17 of individual differences: standardization and initial construct validation. *Journal of*  
18 *neuroscience methods* **187**, 254-262, doi:10.1016/j.jneumeth.2009.11.017 (2010).
- 19 Gur, R. C. *et al.* Age group and sex differences in performance on a computerized neurocognitive  
20 battery in children age 8-21. *Neuropsychology* **26**, 251-265, doi:10.1037/a0026712 (2012).
- 21 Moore, T. M., Reise, S. P., Gur, R. E., Hakonarson, H. & Gur, R. C. Psychometric properties of the  
22 Penn Computerized Neurocognitive Battery. *Neuropsychology* **29**, 235-246,  
23 doi:10.1037/neu0000093 (2015).
- 24 Verma, S. S. *et al.* Imputation and quality control steps for combining multiple genome-wide  
25 datasets. *Frontiers in Genetics* **5**, 370 (2014).
- 26 Hinrichs, A. S. *et al.* The UCSC Genome Browser Database: update 2006. *Nucleic acids research*  
27 **34**, D590-598, doi:10.1093/nar/gkj144 (2006).
- 28 Delaneau, O., Zagury, J. F. & Marchini, J. Improved whole-chromosome phasing for disease and  
29 population genetic studies. *Nature methods* **10**, 5-6, doi:10.1038/nmeth.2307 (2013).
- 30 Howie, B. N., Donnelly, P. & Marchini, J. A flexible and accurate genotype imputation method for  
31 the next generation of genome-wide association studies. *PLoS genetics* **5**, e1000529,  
32 doi:10.1371/journal.pgen.1000529 (2009).
- 33 Anderson, C. A. *et al.* Data quality control in genetic case-control association studies. *Nature*  
34 *protocols* **5**, 1564-1573, doi:10.1038/nprot.2010.116 (2010).
- 35 Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide complex trait  
36 analysis. *American journal of human genetics* **88**, 76-82, doi:10.1016/j.ajhg.2010.11.011 (2011).
- 37 Schmidt, A. *et al.* Acute effects of heroin on negative emotional processing: relation of amygdala  
38 activity and stress-related responses. *Biological psychiatry* **76**, 289-296,  
39 doi:10.1016/j.biopsych.2013.10.019 (2014).
- 40 Krueger F. *Trim Galore!*, <[http://www.bioinformatics.babraham.ac.uk/projects/trim\\_galore/](http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/)> (  
41 Wingett, S. *et al.* HiCUP: pipeline for mapping and processing Hi-C data. *F1000Research* **4**, 1310,  
42 doi:10.12688/f1000research.7334.1 (2015).
- 43 Mifsud, B. *et al.* GOTHIC, a probabilistic model to resolve complex biases and to identify real  
44 interactions in Hi-C data. *PloS one* **12**, e0174744, doi:10.1371/journal.pone.0174744 (2017).
- 45 Schmitt, A. D. *et al.* A Compendium of Chromatin Contact Maps Reveals Spatially Active Regions  
46 in the Human Genome. *Cell reports* **17**, 2042-2059, doi:10.1016/j.celrep.2016.10.061 (2016).
- 47 Wang, K., Li, M. & Bucan, M. Pathway-based approaches for analysis of genomewide association  
48 studies. *American journal of human genetics* **81**, 1278-1283, doi:10.1086/522374 (2007).
- 49 Wang, K., Li, M. & Hakonarson, H. Analysing biological pathways in genome-wide association  
50 studies. *Nature reviews. Genetics* **11**, 843-854, doi:10.1038/nrg2884 (2010).
- 51 Romero, P. *et al.* Computational prediction of human metabolic pathways from the complete  
52 human genome. *Genome biology* **6**, R2, doi:10.1186/gb-2004-6-1-r2 (2005).
- 53 Kandasamy, K. *et al.* NetPath: a public resource of curated signal transduction pathways.  
54 *Genome biology* **11**, R3, doi:10.1186/gb-2010-11-1-r3 (2010).
- 55 Fabregat, A. *et al.* The Reactome pathway Knowledgebase. *Nuc Acids Res* **44**, D481-487,  
56 doi:10.1093/nar/gkv1351 (2016).
- 57 Schaefer, C. F. *et al.* PID: the Pathway Interaction Database. *Nucleic acids research* **37**, D674-679,  
58 doi:10.1093/nar/gkn653 (2009).
- 59 Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for  
60 interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences*  
61 *of the United States of America* **102**, 15545-15550, doi:10.1073/pnas.0506580102 (2005).

- 35 Mi, H. *et al.* The PANTHER database of protein families, subfamilies, functions and pathways.  
36 *Nucleic acids research* **33**, D284-288, doi:10.1093/nar/gki078 (2005).
- 37 Consortium, T. G. O. The Gene Ontology Resource: 20 years and still GOing strong. *Nucleic acids*  
38 *research* **47**, D330-d338, doi:10.1093/nar/gky1055 (2019).
- 39 Darmanis, S. *et al.* A survey of human brain transcriptome diversity at the single cell level.  
40 *Proceedings of the National Academy of Sciences of the United States of America* **112**, 7285-7290,  
41 doi:10.1073/pnas.1507125112 (2015).
- 42 Lake, B. B. *et al.* Neuronal subtypes and diversity revealed by single-nucleus RNA sequencing of  
43 the human brain. *Science (New York, N.Y.)* **352**, 1586-1590, doi:10.1126/science.aaf1204  
44 (2016).
- 45 Nowakowski, T. J. *et al.* Spatiotemporal gene expression trajectories reveal developmental  
46 hierarchies of the human cortex. *Science (New York, N.Y.)* **358**, 1318-1323,  
47 doi:10.1126/science.aap8809 (2017).
- 48 Yu, N. Y. *et al.* Complementing tissue characterization by integrating transcriptome profiling  
49 from the Human Protein Atlas and from the FANTOM5 consortium. *Nucleic acids research* **43**,  
50 6787-6798, doi:10.1093/nar/gkv608 (2015).
- 51 Kang, H. J. *et al.* Spatio-temporal transcriptome of the human brain. *Nature* **478**, 483-489,  
52 doi:10.1038/nature10523 (2011).
- 53 Velasquez, E. *et al.* Synaptosomal Proteome of the Orbitofrontal Cortex from Schizophrenia  
54 Patients Using Quantitative Label-Free and iTRAQ-Based Shotgun Proteomics. *Journal of*  
55 *proteome research* **16**, 4481-4494, doi:10.1021/acs.jproteome.7b00422 (2017).
- 56 Kohler, S. *et al.* Expansion of the Human Phenotype Ontology (HPO) knowledge base and  
57 resources. *Nucleic acids research* **47**, D1018-d1027, doi:10.1093/nar/gky1105 (2019).
- 58 Zhong, S. *et al.* A single-cell RNA-seq survey of the developmental landscape of the human  
59 prefrontal cortex. *Nature* **555**, 524-528, doi:10.1038/nature25980 (2018).
- 60 Merico, D., Isserlin, R., Stueker, O., Emili, A. & Bader, G. D. Enrichment map: a network-based  
61 method for gene-set enrichment visualization and interpretation. *PloS one* **5**, e13984,  
62 doi:10.1371/journal.pone.0013984 (2010).
- 63 Shannon, P. *et al.* Cytoscape: a software environment for integrated models of biomolecular  
64 interaction networks. *Genome research* **13**, 2498-2504, doi:10.1101/gr.1239303 (2003).
- 65 Franz, M. *et al.* GeneMANIA update 2018. *Nucleic acids research* **46**, W60-w64,  
66 doi:10.1093/nar/gky311 (2018).
- 67 Cotto, K. C. *et al.* DGIdb 3.0: a redesign and expansion of the drug-gene interaction database.  
68 *Nucleic acids research* **46**, D1068-d1073, doi:10.1093/nar/gkx1143 (2018).
- 69 Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide association studies,  
70 targeted arrays and summary statistics 2019. *Nucleic acids research* **47**, D1005-d1012,  
71 doi:10.1093/nar/gky1120 (2019).
- 72 Carvalho-Silva, D. *et al.* Open Targets Platform: new developments and updates two years on.  
73 *Nucleic acids research* **47**, D1056-d1065, doi:10.1093/nar/gky1133 (2019).
- 74 1000 Genomes Project Consortium *et al.* A global reference for human genetic variation. *Nature*  
75 **526**, 68-74, doi:10.1038/nature15393 (2015).
- 76 Yaltho, T. C. & Ondo, W. G. The use of gabapentin enacarbil in the treatment of restless legs  
77 syndrome. *Therapeutic advances in neurological disorders* **3**, 269-275,  
78 doi:10.1177/1756285610378059 (2010).
- 79 Stark, C. *et al.* BioGRID: a general repository for interaction datasets. *Nucleic acids research* **34**,  
80 D535-539, doi:10.1093/nar/gkj109 (2006).
- 81 Wang, D. *et al.* Comprehensive functional genomic resource and integrative model for the  
82 human brain. *Science (New York, N.Y.)* **362**, doi:10.1126/science.aat8464 (2018).

- 55 Greenwood, T. A., Akiskal, H. S., Akiskal, K. K. & Kelsoe, J. R. Genome-wide association study of  
56 temperament in bipolar disorder reveals significant associations with three novel Loci.  
57 *Biological psychiatry* **72**, 303-310, doi:10.1016/j.biopsych.2012.01.018 (2012).
- 58 Simon-Sanchez, J. *et al.* Genome-wide association study reveals genetic risk underlying  
59 Parkinson's disease. *Nature genetics* **41**, 1308-1312, doi:10.1038/ng.487 (2009).
- 60 Nagel, M. *et al.* Meta-analysis of genome-wide association studies for neuroticism in 449,484  
61 individuals identifies novel genetic loci and pathways. *Nature genetics* **50**, 920-927,  
62 doi:10.1038/s41588-018-0151-7 (2018).
- 63 Consortium, N. a. P. A. S. o. P. G. Psychiatric genome-wide association study analyses implicate  
64 neuronal, immune and histone pathways. *Nature neuroscience* **18**, 199-209,  
65 doi:10.1038/nn.3922 (2015).
- 66 Palumbo, P. *et al.* Clinical and molecular characterization of an emerging chromosome  
67 22q13.31 microdeletion syndrome. *American journal of medical genetics. Part A* **176**, 391-398,  
68 doi:10.1002/ajmg.a.38559 (2018).
- 69 Vadgama, N., Lamont, D., Hardy, J., Nasir, J. & Lovering, R. C. Distinct proteomic profiles in  
70 monozygotic twins discordant for ischaemic stroke. *Molecular and cellular biochemistry*,  
71 doi:10.1007/s11010-019-03501-2 (2019).
- 72 Matsuura, T. *et al.* Large expansion of the ATTCT pentanucleotide repeat in spinocerebellar  
73 ataxia type 10. *Nature genetics* **26**, 191-194, doi:10.1038/79911 (2000).
- 74 Forbes, N. F., Carrick, L. A., McIntosh, A. M. & Lawrie, S. M. Working memory in schizophrenia: a  
75 meta-analysis. *Psychological medicine* **39**, 889-905, doi:10.1017/s0033291708004558 (2009).
- 76 Huntley, J. D. & Howard, R. J. Working memory in early Alzheimer's disease: a  
77 neuropsychological review. *International journal of geriatric psychiatry* **25**, 121-132,  
78 doi:10.1002/gps.2314 (2010).
- 79 Kehagia, A. A., Barker, R. A. & Robbins, T. W. Neuropsychological and clinical heterogeneity of  
80 cognitive impairment and dementia in patients with Parkinson's disease. *The Lancet. Neurology*  
81 **9**, 1200-1213, doi:10.1016/s1474-4422(10)70212-x (2010).
- 82 Rose, E. J. & Ebmeier, K. P. Pattern of impaired working memory during major depression.  
83 *Journal of affective disorders* **90**, 149-161, doi:10.1016/j.jad.2005.11.003 (2006).
- 84 Stopford, C. L., Thompson, J. C., Neary, D., Richardson, A. M. & Snowden, J. S. Working memory,  
85 attention, and executive function in Alzheimer's disease and frontotemporal dementia. *Cortex; a  
86 journal devoted to the study of the nervous system and behavior* **48**, 429-446,  
87 doi:10.1016/j.cortex.2010.12.002 (2012).
- 88 Wang, Y. *et al.* A Meta-Analysis of Working Memory Impairments in Autism Spectrum  
89 Disorders. *Neuropsychology review* **27**, 46-61, doi:10.1007/s11065-016-9336-y (2017).
- 90 Clementz, B. A. *et al.* Identification of Distinct Psychosis Biotypes Using Brain-Based  
91 Biomarkers. *The American journal of psychiatry* **173**, 373-384,  
92 doi:10.1176/appi.ajp.2015.14091200 (2016).
- 93 Jeste, S. S. & Geschwind, D. H. Disentangling the heterogeneity of autism spectrum disorder  
94 through genetic findings. *Nature reviews. Neurology* **10**, 74-81, doi:10.1038/nrneuro.2013.278  
95 (2014).
- 96 Wang, W. *et al.* Pathway-based discovery of genetic interactions in breast cancer. *PLoS genetics*  
97 **13**, e1006973, doi:10.1371/journal.pgen.1006973 (2017).
- 98 Smedley, D. *et al.* The BioMart community portal: an innovative alternative to large, centralized  
99 data repositories. *Nucleic acids research* **43**, W589-598, doi:10.1093/nar/gkv350 (2015).
- 100 Watanabe, K., Taskesen, E., van Bochoven, A. & Posthuma, D. Functional mapping and  
101 annotation of genetic associations with FUMA. *Nature communications* **8**, 1826,  
102 doi:10.1038/s41467-017-01261-5 (2017).
- 103 Robinson, J. T. *et al.* Integrative genomics viewer. *Nature biotechnology* **29**, 24-26,  
104 doi:10.1038/nbt.1754 (2011).

- 510 74 Thorvaldsdottir, H., Robinson, J. T. & Mesirov, J. P. Integrative Genomics Viewer (IGV): high-  
511 performance genomics data visualization and exploration. *Briefings in bioinformatics* **14**, 178-  
512 192, doi:10.1093/bib/bbs017 (2013).
- 513 75 Gold, J. M. Cognitive deficits as treatment targets in schizophrenia. *Schizophrenia research* **72**,  
514 21-28, doi:10.1016/j.schres.2004.09.008 (2004).
- 515
- 516

## Figures

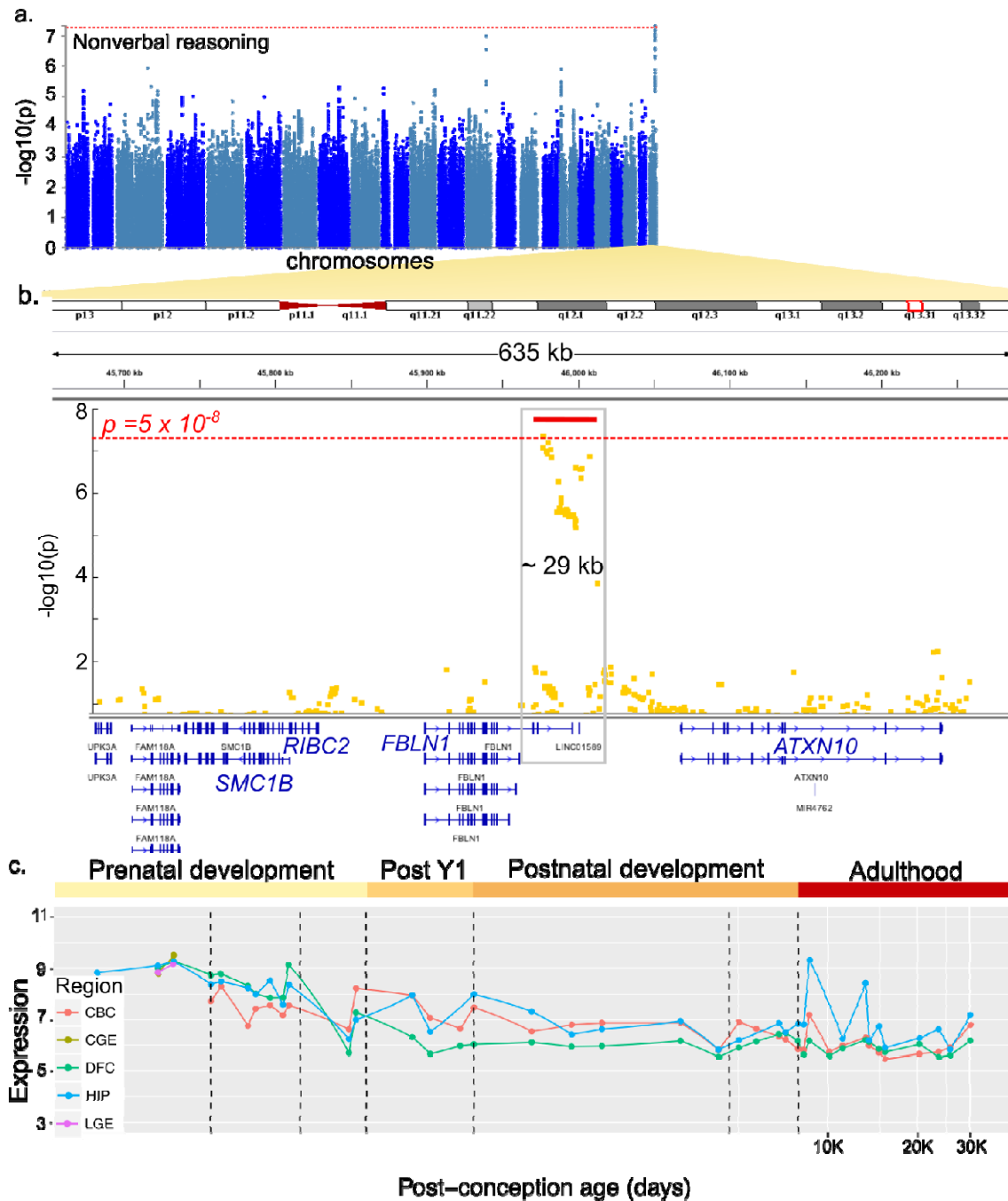


517  
518

519  
520



521 **Figure 1.** Genome-wide association analysis for neurocognitive phenotypes from the Philadelphia  
522 Neurodevelopmental Cohort.  
523 a. Workflow. Genotypes were imputed (1KGP reference), and limited to European samples. Samples  
524 with severe medical conditions were removed and invalid test scores excluded. Nine neurocognitive  
525 test scores were binarized after age and sex had been regressed out. GWAS was performed for  
526 accuracy for each of these nine phenotypes.  
527 b. Breakdown of SNPs achieving suggestive significance, by phenotype (top).  
528 c. Suggestive and significant SNPs and associated genes. The outermost ring shows the location of  
529 suggestive peaks ( $p < 10^{-5}$ ), coloured by phenotype (see b); y-axis shows  $-\log_{10}(\text{SNP } p)$ , so that SNPs  
530 with stronger significance are higher. SNPs with  $p < 10^{-7}$  are labeled. The tracks with ticks indicate  
531 functional consequences of associated SNPs. The track closest to the middle indicate SNPs overlapping  
532 brain enhancers (light gray) or promoters (black). The dark red middle track indicates SNPs with  
533 nonsynonymous variation, including NMD transcript, missense or splice variants<sup>71</sup>. The outermost  
534 track indicates QTL associations, including eQTL in adult prefrontal cortex (dark blue), fetal brain  
535 (cyan), or neuronal cell proportions in the adult brain (fQTL; orange). Genes associated with top SNPs  
536 are indicated within the circle.  
537 d. Genes associated with top SNPs ( $p < 3 \times 10^{-7}$ ) with prior knowledge about relevance to brain  
538 development and psychiatric disorders. Columns indicate differential expression in  
539 neurodevelopmental disorders<sup>54</sup> (SCZ = schizophrenia; ASD= autism), significant association with a  
540 nervous system disorder<sup>49</sup>, or status as marker gene for specific cell types in fetal brain<sup>39</sup>.  
541 e. Breakdown of functional consequence of top SNPs and by functional consequence (bottom).  
542 Consequence shown is limited to effect on protein sequence<sup>71</sup>, presence in enhancers or promoters in  
543 adult cortical regions<sup>11</sup>, eQTL in fetal brain, or adult forebrain. Final bar shows cumulative proportion  
544 of putatively functional SNPs.



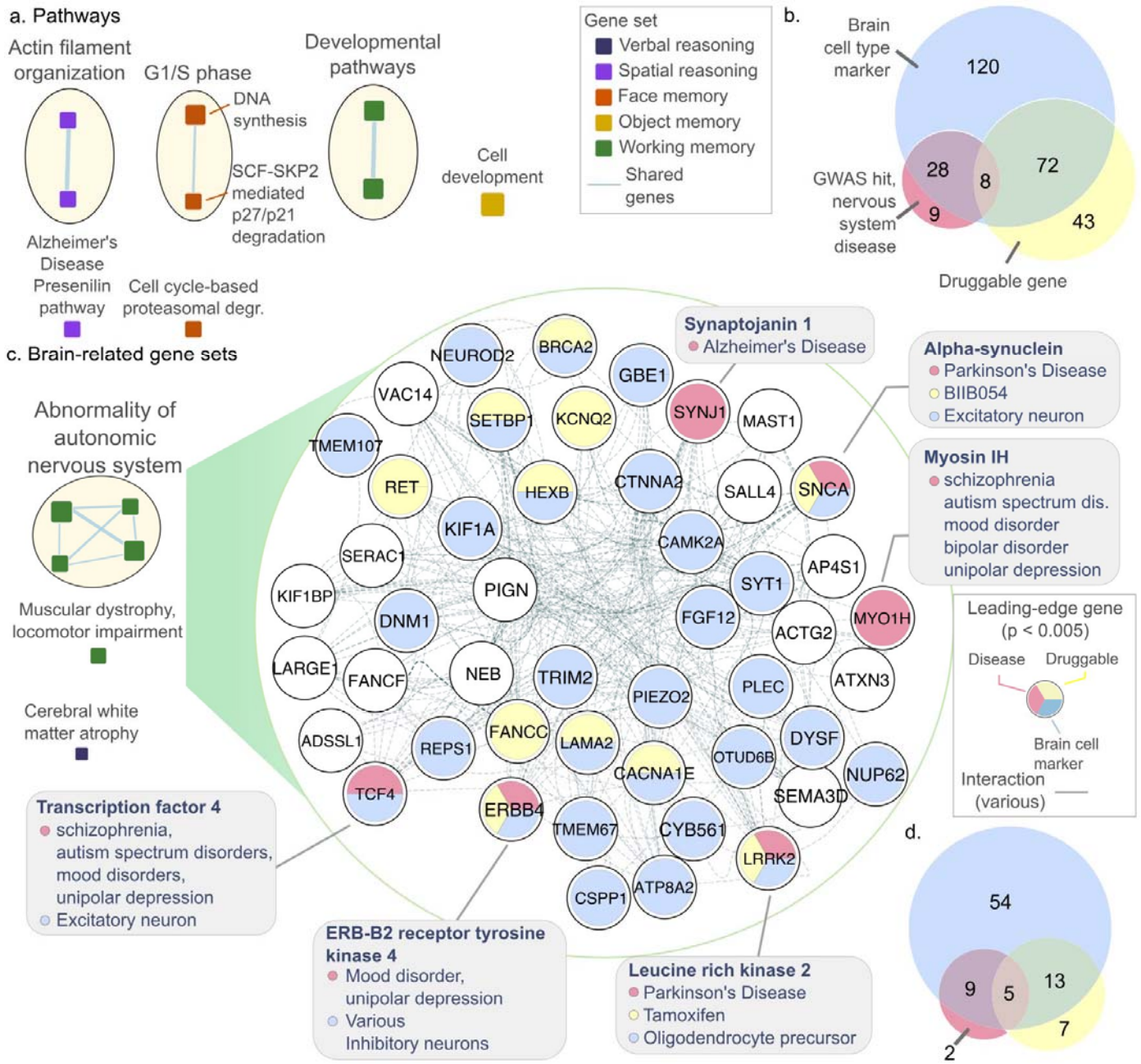
**Figure 2.** Genome-wide significance of *FBLN1* region for binarized performance in nonverbal reasoning

a. Manhattan plot of univariate SNP association with binarized performance in nonverbal reasoning (N=1,024 poor vs. 1,023 good performers; 4,893,197 SNPs). Plot generated using FUMA<sup>72</sup>.

b. Detailed view of hit region at chr22q13. Two SNPs pass genome-wide significance threshold, rs77601382 and rs74825248 ( $p=4.64e-8$ ). View using Integrated Genome Viewer (v2.3.93<sup>73,74</sup>). The red bar indicates the region with increased SNP-level association.

c. *FBLN1* transcription in the human brain through the lifespan. Data from BrainSpan<sup>41</sup>. Log-transformed normalized expression is shown for cerebellar cortex (CBC), central ganglionic eminence (CGE) and lateral ganglionic eminence (LGE), dorsal frontal cortex (DFC), and hippocampus (HIP).

557



**Figure 3.** Pathway and gene set enrichment analysis for neurocognitive task performance

a. Pathways significantly enriched for genetic variation in neurocognitive task performance (GSEA, 100 permutations,  $q < 0.1$ , Supplementary Tables 8, 9, 10). Nodes indicate pathways, with fill indicating phenotype and yellow bubbles denoting clusters of related gene sets; edges indicate shared genes.

b. Number of leading edge genes associated with transcription in specific brain cell types (blue), drug targets (yellow) or genetic associations with specific nervous system disorders (pink) (pathways with  $q < 0.10$ ,  $N=355$  genes).

c. Brain-related gene sets enriched for genetic variation in task performance. Left: Significant gene sets; legend same as panel a (Supplementary Tables 11,12, 13). Right: Top leading edge genes in enriched brain-related gene sets ( $N=48$  genes,  $p < 5e-3$ , pathways with  $q < 0.05$ ). Nodes show genes and fill indicates genes associated with brain cell types, drugs or genetic associations with nervous system disorders (white indicates absence of association). Edges indicate known interactions (GeneMANIA<sup>47</sup>). Genes with disease associations have been highlighted in grey pullout bubbles.

558

559

560

561

562

563

564

565

566

567

568

569

570

571

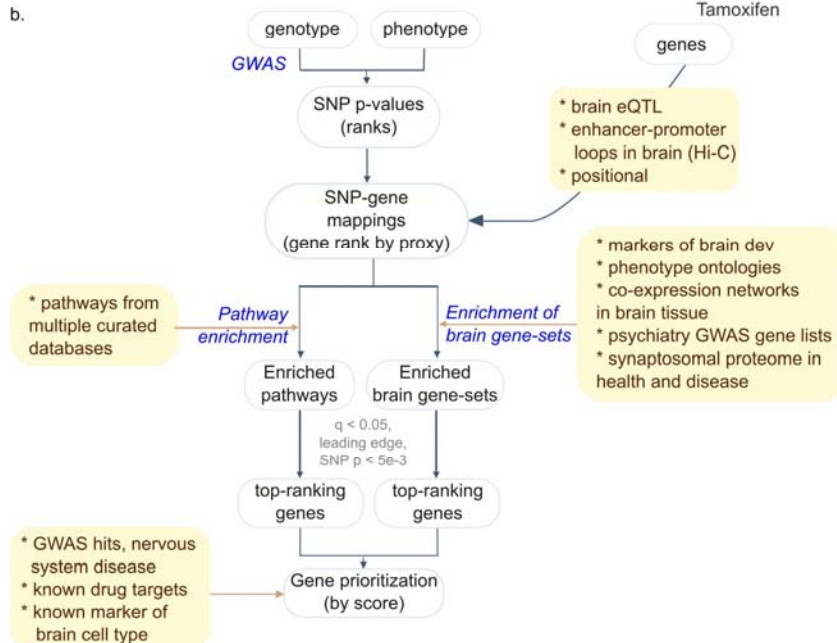
572

573 d. Leading edge genes in brain-related gene sets associated with disease, drugs or brain cell types  
574 (N=134 genes); legend as in b.

a.

Cognitive domain	Complex cognition (reasoning)			Executive function			Declarative memory		Social processing
	Verbal	Nonverbal	Spatial	Attention	Working mem.	Object	Face	Word	Emotion Ident.
<b>Task</b>									
<b>Variants #, <math>p &lt; 1 \times 10^{-5}</math></b>	75	75	219	24	153	24	16	43	24
<b>Genes</b> $p < 1 \times 10^{-5}$ $* p < 5 \times 10^{-8}$	<i>CEP162</i> <i>EEF1AKMT1</i> <i>MST1R</i> <i>NAV2</i> <i>NOTO</i> <i>NRXN1</i> <i>PGBD5</i> <i>PRKN</i> <i>PTPRQ</i> <i>RBM6</i> <i>XPO4</i>	<i>FBLN1*</i> <i>BTBD11</i> <i>C16orf96</i> <i>CDIP1</i> <i>CLSTN2</i> <i>CORO7</i> <i>FNBP1L</i> <i>MGRN1</i> <i>MREG</i> <i>NMRAL1</i> <i>PTPRD</i> <i>SLC26A3</i>	<i>ACADS</i> <i>BLOC1S5</i> <i>BLOC1S5-TXNDC5</i> <i>CDC45</i> <i>EEF1E1-BLOC1S5</i> <i>NBEA</i> <i>PRDM2</i> <i>SPPL3</i>	<i>INSC</i> <i>SRGAP3</i>	<i>CACNA2D3</i> <i>DNAH14</i> <i>FAT3</i> <i>KMT5A</i> <i>NCAM2</i> <i>SNAP91</i> <i>TECPR2</i> <i>THEMIS</i> <i>WISP3</i>	<i>BMP5</i> <i>CDH23</i> <i>FERMT1</i> <i>IGFBP7</i> <i>MBNL2</i> <i>PTPN13</i>	<i>ADAMTS14</i> <i>ETS1</i> <i>PCP4</i>	<i>FSTL5</i> <i>KCNC4</i> <i>TENM2</i>	<i>C6orf10</i> <i>CLEC16A</i> <i>EPG5</i> <i>KCNJ6</i> <i>NUP210L</i> <i>TOX2</i>
<b>Pathway themes</b> $q < 0.1$ $* q < 0.05$		Actin filament org.; Presenilin pathway			Developmental* pathways	Reg. of cell devel.		G1-S phase of cell cycle; NOD-like receptor sig.	
<b>Brain-related gene set</b> $q < 0.1$ $* q < 0.05$			White matter atrophy		Autonomic* nervous system dysfunc. Inability to walk*				
<b>Cells (fetal brain), gene assoc.</b>		Neural progenitors	Excitatory & inhibitory neurons	Excitatory neurons Radial glia					
<b>Disease, gene assoc.</b>		Schizoph. Bipolar disorder	Schizoph. Migraine	Alzheimer's	Parkinson's Autism Depression Alzheimer's Schizoph.				Conduct disorder
<b>Disease-Task association</b>				Attention <sup>1</sup> deficit; Psychosis risk	Schizoph. <sup>1</sup> Bipolar dis. TBI Schizoph. Bipolar dis. TBI		Schizoph. <sup>1,2</sup> Epilepsy		Various neuropsych. <sup>1</sup> Depression Schizoph.
<b>Drugs, gene assoc.</b>					AEE788 BIIB054 Chlorambucil Ergocalciferol Everolimus Imagabalin Nerispiridine Ocriplasmin Omeprazole Regorafenib Tamoxifen				

b.



576 **Figure 4.** a. Association of top genes, gene sets, and pathways with different levels of brain  
577 organization. Each column shows data for an individual phenotype, grouped by domain; rows show  
578 associations at increasingly higher levels (from top to bottom), and finally with drug targets. All results  
579 are from this work unless otherwise cited. Circles indicate relative number of suggestive variant peaks  
580 ( $p < 10^{-5}$ ) from GWAS (median=43; mean=73.4), with numbers indicated below (asterisk:  $p < 5 \times 10^{-8}$ ),  
581 and genes are those mapped to top-ranking SNPs ( $p \leq 1 \times 10^{-5}$ ) (only protein-coding genes; noncoding  
582 genes listed in Supplementary Table 14). Pathways and brain-related gene sets shown are those  
583 passing  $q < 0.1$  in enrichment analysis (red asterisk:  $q \leq 0.05$ ). Fetal brain cell associations are as  
584 shown in Figure 1d. Gene-disease associations combine those for top GWAS SNPs (Figure 1d) and from  
585 gene set enrichment analysis; drug associations are from the latter (Supplementary Tables 10 and 14).  
586 Prior associations of alterations in phenotype or task-based brain activation as described in <sup>14</sup>(1) or  
587 <sup>75</sup>(2).

588 b. Proposed workflow for gene prioritization, as used in this work. When provided with genotype-  
589 phenotype data, SNPs are first prioritized by assigning an association statistic (e.g. by GWAS). Gene set  
590 enrichment analysis is performed to identify groups of genes with subthreshold phenotype  
591 association. SNP-gene mappings use brain-specific maps of genome regulation, prioritizing evidence-  
592 based association over positional mapping. Enrichment of pathways and brain-related gene sets are  
593 simultaneously performed using a rank-based method such as GSEA, which provides a leading edge  
594 subset for subsequent prioritization. Leading edge genes are annotated with clinical attributes of  
595 interest, such as druggability, prior disease association and evidence for expression in particular brain  
596 cell types, and the combination of attributes can be turned into a prioritization score.

## Tables

Phenotype	N	# lead SNPs (p < 1e-5)	Indiv. Sig. SNPs (p < 1e-6)	SNP p	Gene
<b>Complex Cognition</b>					
Verbal	2,068	83	-		
Non-verbal	2,047	75	rs77601382	4.6X10 <sup>-8</sup>	<i>FBLN1</i>
			rs76901846	1.0 X10 <sup>-7</sup>	<i>BTBD11</i>
			rs5765534	1.4 X10 <sup>-7</sup>	
Spatial	2,024	219	rs446816	2.6 X10 <sup>-7</sup>	<i>NBEA</i>
			rs7001721	8.5 X10 <sup>-7</sup>	
<b>Executive Function</b>					
Working memory	2,047	153	rs565936	6.6 X10 <sup>-7</sup>	<i>FAT3</i>
			rs2093484	9.3 X10 <sup>-7</sup>	
Attention	2,041	24	rs11992719	5.1 X10 <sup>-8</sup>	
			rs1792551	9.3 X10 <sup>-7</sup>	<i>INSC</i>
<b>Social processing</b>					
Emotion Identification <sup>o</sup>	2,068	24	rs73118294	7.1 X10 <sup>-8</sup>	<i>TOX2</i>
			rs4341378	4.9 X10 <sup>-7</sup>	
<b>Declarative memory</b>					
Face memory	2,066	16	rs6926533; rs148111284	5.4 X10 <sup>-7</sup> ; 6.9 X10 <sup>-7</sup>	<i>RBMXP1</i> <i>PCP4</i>
Word memory	2,073	43	-		
Object memory	2,070	24	rs56659368	3.2 X10 <sup>-7</sup>	

**Table 1.** Genetic variants significantly associated with neurocognitive phenotypes in the Philadelphia Neurodevelopmental Cohort (PNC) dataset. For each test in the PNC neurocognitive test battery, GCTA was run to obtain SNP-level (marginal) p-values associated with binarized (good or poor) performance. Top SNPs (p < 1.0x10<sup>-6</sup>) are shown above (full list of suggestive SNPs in Supplementary Table 5). SNPs were mapped to genes based on expression modulation, chromatin interaction of positional information. Only protein-coding genes shown here; additional non-coding RNA associations shown in Supplementary Table 7).