

Quantifying concordant genetic effects of *de novo* mutations on multiple disorders

Hanmin Guo^{1,2}, Lin Hou^{1,2,3}, Yu Shi⁴, Sheng Chih Jin⁵, Xue Zeng^{6,7}, Boyang Li⁸, Richard P. Lifton^{6,7}, Martina Brueckner^{6,9}, Hongyu Zhao^{6,8,10,#}, Qiongshi Lu^{11,#,*}

¹ Center for Statistical Science, Tsinghua University, Beijing, China

² Department of Industrial Engineering, Tsinghua University, Beijing, China

³ MOE Key Laboratory of Bioinformatics, School of Life Sciences, Tsinghua University, Beijing, China

⁴ Yale School of Management, Yale University, New Haven, CT, USA

⁵ Department of Genetics, Washington University School of Medicine, St. Louis, MO, USA

⁶ Department of Genetics, Yale University, New Haven, CT, USA

⁷ Laboratory of Human Genetics and Genomics, Rockefeller University, New York, USA

⁸ Department of Biostatistics, Yale School of Public Health, New Haven, CT, USA

⁹ Department of Pediatrics, Yale University School of Medicine, New Haven, CT, USA

¹⁰ Program of Computational Biology and Bioinformatics, Yale University, New Haven, CT, USA

¹¹ Department of Biostatistics and Medical Informatics, University of Wisconsin-Madison, Madison, WI, USA

These authors should be considered shared last author

* Correspondence:

Dr. Qiongshi Lu

Department of Biostatistics and Medical Informatics

University of Wisconsin–Madison

425 Henry Mall

Madison, WI, USA 53706

qlu@biostat.wisc.edu

31 **Abstract**

32

33 Exome sequencing on tens of thousands of parent-proband trios has identified numerous
34 deleterious *de novo* mutations (DNMs) and implicated risk genes for many disorders. Recent
35 studies have suggested shared genes and pathways are enriched for DNMs across multiple
36 disorders. However, existing analytic strategies only focus on genes that reach statistical
37 significance for multiple disorders and require large trio samples in each study. As a result, these
38 methods are not able to characterize the full landscape of genetic sharing due to polygenicity and
39 incomplete penetrance. In this work, we introduce EncoreDNM, a novel statistical framework to
40 quantify shared genetic effects between two disorders characterized by concordant enrichment
41 of DNMs in the exome. EncoreDNM makes use of exome-wide, summary-level DNM data,
42 including genes that do not reach statistical significance in single-disorder analysis, to evaluate
43 the overall and annotation-partitioned genetic sharing between two disorders. Applying
44 EncoreDNM to DNM data of nine disorders, we identified abundant pairwise enrichment
45 correlations, especially in genes intolerant to pathogenic mutations and genes highly expressed
46 in fetal tissues. These results suggest that EncoreDNM improves current analytic approaches and
47 may have broad applications in DNM studies.

48

49

50 Introduction

51
52 *De novo* mutations (DNMs) can be highly deleterious and provide important insights into the
53 genetic cause for disease¹. As the cost of sequencing continues to drop, whole-exome
54 sequencing (WES) studies conducted on tens of thousands of family trios have pinpointed
55 numerous risk genes for a variety of disorders²⁻⁴. In addition, accumulating evidence suggests
56 that risk genes enriched for pathogenic DNMs may be shared by multiple disorders⁵⁻⁹. These
57 shared genes could reveal biological pathways that play prominent roles in disease etiology and
58 shed light on clinically heterogeneous yet genetically related diseases⁷⁻⁹.

59
60 Most efforts to identify shared risk genes directly compare genes that are significantly associated
61 with each disorder^{10,11}. There have been some successes with this approach in identifying shared
62 genes and pathways (e.g., chromatin modifiers) underlying developmental disorder (DD), autism
63 spectrum disorder (ASD), and congenital heart disease (CHD), thanks to the large trio samples
64 in these studies^{3,4,12}, whereas findings in smaller studies remain suggestive^{13,14}. Even in the
65 largest studies to date, statistical power remains moderate for risk genes with weaker effects^{3,15}.
66 It is estimated that more than 1,000 genes associated with DD remain undetected³. Therefore,
67 analytic approaches that only account for top significant genes cannot capture the full landscape
68 of genetic sharing in multiple disorders. Recently, a Bayesian framework was proposed to jointly
69 analyze DNM data of two diseases and improve risk gene mapping⁹. Although some parameters
70 in this framework can quantify shared genetics between diseases, the statistical property of these
71 parameter estimates have not been studied. There is a pressing need for powerful, robust, and
72 interpretable methods that quantify concordant DNM association patterns for multiple disorders
73 using exome-wide DNM counts.

74
75 Recent advances in estimating genetic correlations using summary data from genome-wide
76 association studies (GWAS) may provide a blueprint for approaching this problem in DNM
77 research¹⁶. Modeling “omnigenic” associations as independent random effects, linear mixed-
78 effects models leverage genome-wide association profiles to quantify the correlation between
79 additive genetic components of multiple complex traits¹⁷⁻²⁰. These methods have identified
80 ubiquitous genetic correlations across many human traits and revealed significant and robust
81 genetic correlations that could not be inferred from significant GWAS associations alone²¹⁻²⁴.

82
83 Here, we introduce EncoreDNM (**E**nrichment **c**orrelation **e**stimator for **D**e **N**ovo **M**utations), a
84 novel statistical framework that leverages exome-wide DNM counts, including genes that do not
85 reach exome-wide statistical significance in single-disorder analysis, to estimate concordant DNM
86 associations between disorders. EncoreDNM uses a generalized linear mixed-effects model to
87 quantify the occurrence of DNMs while accounting for *de novo* mutability of each gene and
88 technical inconsistencies between studies. We demonstrate the performance of EncoreDNM
89 through extensive simulations and analyses of DNM data of nine disorders.

90 Results

91 92 Method overview

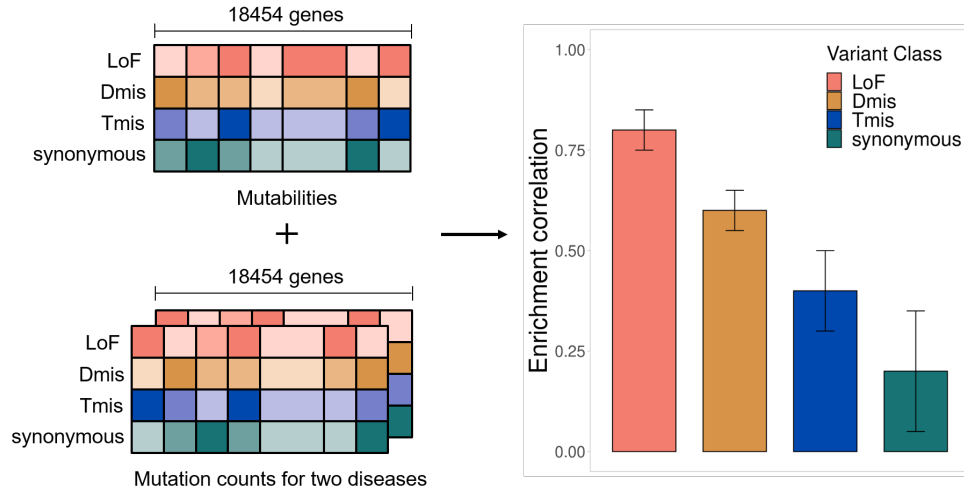
93
94 DNM counts in the exome deviate from the null (i.e., expected counts based on *de novo* mutability)
95 when mutations play a role in disease etiology. Disease risk genes will show enrichment for
96 deleterious DNMs in probands and non-risk genes may be slightly depleted for DNM counts. Our
97 goal is to estimate the correlation of such deviation between two disorders, which we refer to as
98 the DNM enrichment correlation. More specifically, we use a pair of mixed-effects Poisson
99 regression models to quantify the occurrence of DNMs in two studies.

$$\begin{aligned} 100 \quad & \begin{bmatrix} Y_{i1} \\ Y_{i2} \end{bmatrix} \sim \text{Poisson} \left(\begin{bmatrix} \lambda_{i1} \\ \lambda_{i2} \end{bmatrix} \right), \\ 101 \quad & \log \left(\begin{bmatrix} \lambda_{i1} \\ \lambda_{i2} \end{bmatrix} \right) = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} + \log \left(\begin{bmatrix} 2N_1 m_i \\ 2N_2 m_i \end{bmatrix} \right) + \begin{bmatrix} \phi_{i1} \\ \phi_{i2} \end{bmatrix}, \\ 102 \quad & \begin{bmatrix} \phi_{i1} \\ \phi_{i2} \end{bmatrix} \sim \text{MVN} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix} \right). \end{aligned}$$

103 Here, Y_{i1}, Y_{i2} are the DNM counts for the i -th gene and N_1, N_2 are the number of parent-proband
104 trios in two studies, respectively. The log Poisson rates of DNM occurrence are decomposed into
105 three components: the elevation component, the background component, and the deviation
106 component. The elevation component β_k ($k = 1, 2$) is a fixed effect term adjusting for systematic,
107 exome-wide bias in DNM counts. One example of such bias is the batch effect caused by different
108 sequencing and variant calling pipelines in two studies. The elevation parameter β_k tends to be
109 larger when DNMs are over-called with higher sensitivity and smaller when DNMs are under-
110 called with higher specificity²⁵. The background component $\log(2N_k m_i)$ is a gene-specific fixed
111 effect that reflects the expected mutation counts determined by the genomic sequence context
112 under the null²⁶. m_i is the *de novo* mutability for the i -th gene, and $2N_1 m_i$ and $2N_2 m_i$ are the
113 expected DNM counts in the i -th gene under the null in two studies. The deviation component
114 ϕ_{ik} is a gene-specific random effect that quantifies the deviation of DNM profile from what is
115 expected under the null. ϕ_{i1} and ϕ_{i2} follow a multivariate normal distribution with dispersion
116 parameters σ_1 and σ_2 and a correlation ρ . DNM enrichment correlation is denoted by ρ and is
117 our main parameter of interest. It quantifies the concordance of DNM burden in two disorders.

118
119 Parameters in this model can be estimated using a Monte Carlo maximum likelihood estimation
120 (MLE) procedure. Standard errors of the estimates are obtained through inversion of the observed
121 Fisher information matrix. In practice, we use annotated DNM data as input and fit mixed-effects
122 Poisson models for each variant class separately: loss of function (LoF), deleterious missense
123 (Dmis, defined as MetaSVM-deleterious), tolerable missense (Tmis, defined as MetaSVM-
124 tolerable), and synonymous (**Figure 1**). More details about model setup and parameter estimation
125 are discussed in **Methods**.

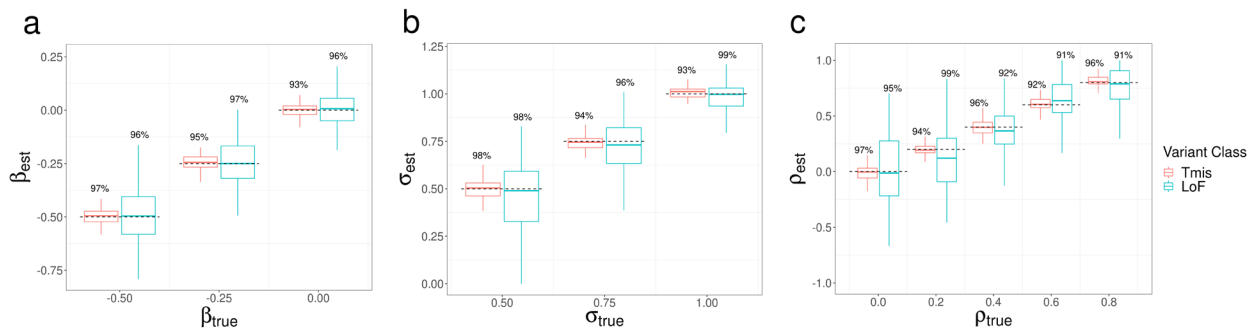
126



127
128 **Figure 1. EncoreDNM workflow.** The inputs of EncoreDNM are *de novo* mutability of each gene and exome-wide,
129 annotated DNM counts from two studies. We fit a mixed-effects Poisson model to estimate the DNM enrichment
130 correlation between two disorders for each variant class.
131

132 **Simulation results**

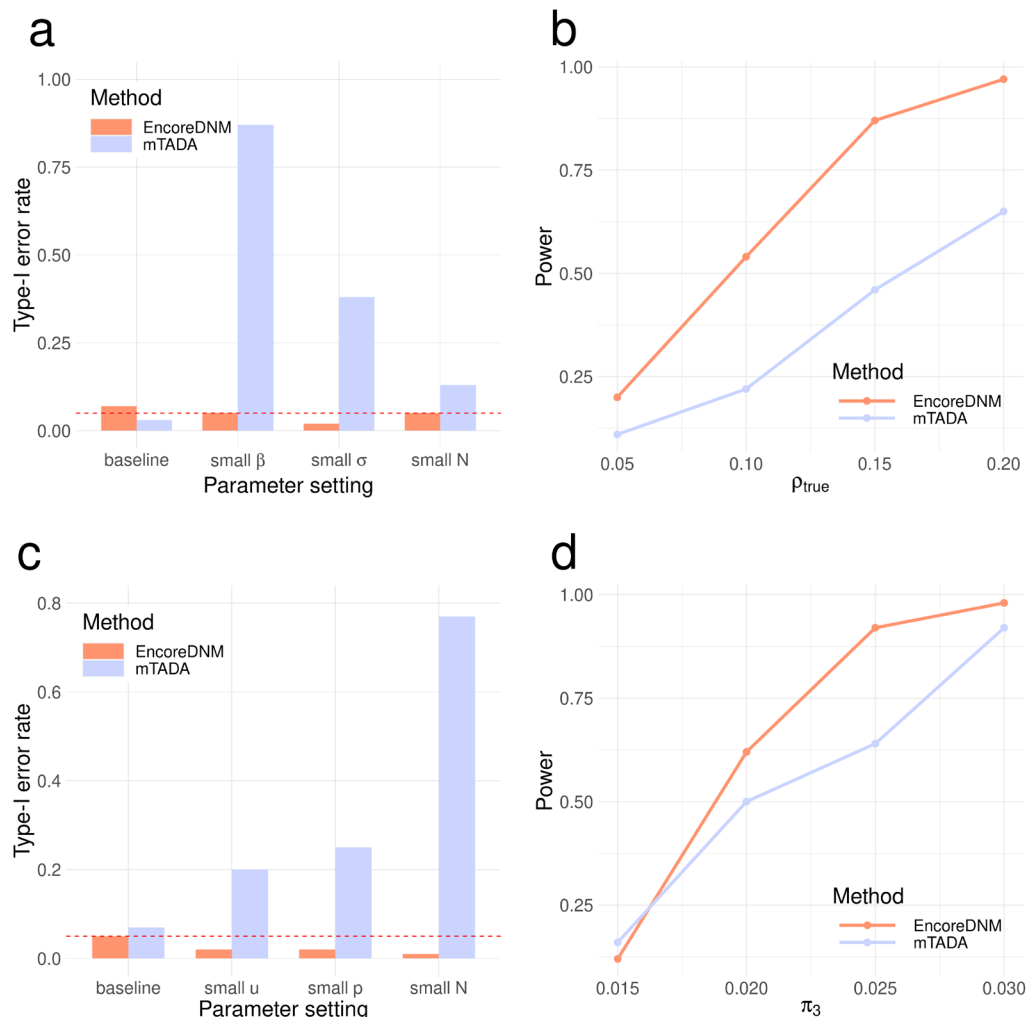
133
134 We conducted simulations to assess the parameter estimation performance of EncoreDNM in
135 various settings. We focused on two variant classes, i.e., Tmis and LoF variants, since they have
136 the highest and lowest median mutabilities in the exome. We used EncoreDNM to estimate the
137 elevation parameter β , dispersion parameter σ , and enrichment correlation ρ (**Methods**).
138 Under various parameter settings, EncoreDNM always provided unbiased estimation of the
139 parameters (**Figure 2** and **Supplementary Figures 1-2**). Furthermore, the 95% Wald confidence
140 intervals achieved coverage rates close to 95% under all simulation settings, demonstrating the
141 effectiveness of EncoreDNM to provide accurate statistical inference.
142



143
144 **Figure 2. Parameter estimation results of EncoreDNM.** (a) Boxplot of β estimates in single-trait analysis with σ
145 fixed at 0.75. (b) Boxplot of σ estimates in single-trait analysis with β fixed at -0.25. (c) Boxplot of ρ estimates in
146 cross-trait analysis with β and σ fixed at -0.25 and 0.75. True parameter values are marked by dashed lines. The
147 number above each box represents the coverage rate of 95% Wald confidence intervals. Each simulation setting was
148 repeated 100 times.
149

150 Next, we compared the performance of EncoreDNM with mTADA⁹, a Bayesian framework that
151 could estimate the proportion of shared risk genes for two disorders. First, we simulated DNM
152 data under the mixed-effects Poisson model. We evaluated two methods across a range of

153 combinations of elevation parameter, dispersion parameter, and sample size for two disorders.
 154 The type-I error rates for our method were well-calibrated in all parameter settings, but mTADA
 155 produced false positive findings when the observed DNM counts were relatively small (e.g., due
 156 to reduced elevation or dispersion parameters or a lower sample size; **Figure 3a**). We also
 157 assessed the statistical power of two approaches under a baseline setting where type-I errors for
 158 both methods were controlled. As enrichment correlation increased, EncoreDNM achieved
 159 universally greater statistical power compared to mTADA (**Figure 3b**).
 160



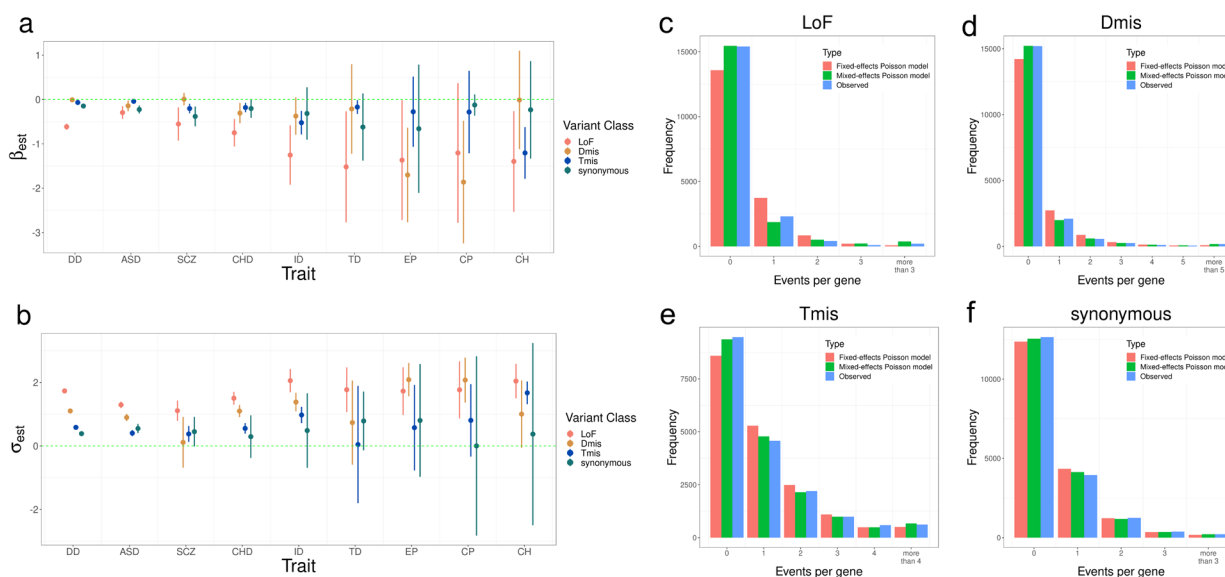
161 **Figure 3. Comparison of EncoreDNM and mTADA.** (a) Type-I error rates under a mixed-effects Poisson regression
 162 model. (β, σ, N) were fixed at (-0.25, 0.75, 5000) under the baseline setting, (-1, 0.75, 5000) under a setting with small
 163 β , (-0.25, 0.5, 5000) under a setting with small σ , and (-0.25, 0.75, 1000) under a setting with small N for two disorders.
 164 (b) Statistical power of two methods under a mixed-effects Poisson regression model as the enrichment correlation
 165 increases. Parameters (β, σ, N) were fixed at (-0.25, 0.75, 5000) for both disorders. (c) Type-I error rates under a
 166 multinomial model. (u, p, N, π^S) were fixed at (0.95, 0.25, 5000, 0.1) under the baseline setting, (0.75, 0.25, 5000, 0.1)
 167 under a setting with small u (i.e., reduced total DNM counts), (0.95, 0.15, 5000, 0.1) under a setting with small p (i.e.,
 168 fewer probands explained by DNMs), and (0.95, 0.25, 1000, 0.1) under a setting with lower sample size. (d) Statistical
 169 power under a multinomial model with varying proportion of shared causal genes. Parameters (u, p, N, π^S) were fixed
 170 at (0.95, 0.25, 5000, 0.1) for both disorders. Each simulation setting was repeated 100 times.
 171
 172

173 To ensure a fair comparison, we also considered a mis-specified model setting where we
 174 randomly distributed the total DNM counts for each disorder into all genes with an enrichment in
 175 causal genes (**Methods**). EncoreDNM showed well-controlled type-I error across all simulation
 176 settings, whereas severe type-I error inflation arose for mTADA when the total mutation count, the
 177 proportion of probands that can be explained by DNMs, or the sample size were small (**Figure**
 178 **3c**). Furthermore, we compared the statistical power of two methods under this model in a
 179 baseline setting where type-I error was controlled. EncoreDNM showed higher statistical power
 180 compared to mTADA as the fraction of shared causal genes increased (**Figure 3d**).

181 182 Pervasive enrichment correlation of damaging DNMs among developmental disorders

183
 184 We applied EncoreDNM to DNM data of nine disorders (**Supplementary Table 1; Methods**): DD
 185 ($n=31,058$; number of trios)³, ASD ($n=6,430$)⁴, schizophrenia (SCZ; $n=2,772$)¹⁵, CHD ($n=2,645$)¹²,
 186 intellectual disability (ID; $n=820$)², Tourette disorder (TD; $n=484$)²⁷, epileptic encephalopathies (EP;
 187 $n=264$)¹³, cerebral palsy (CP; $n=250$)¹⁴, and congenital hydrocephalus (CH; $n=232$)²⁸. In addition,
 188 we also included 1,789 trios comprising healthy parents and unaffected siblings of ASD probands
 189 as controls²⁹.

190



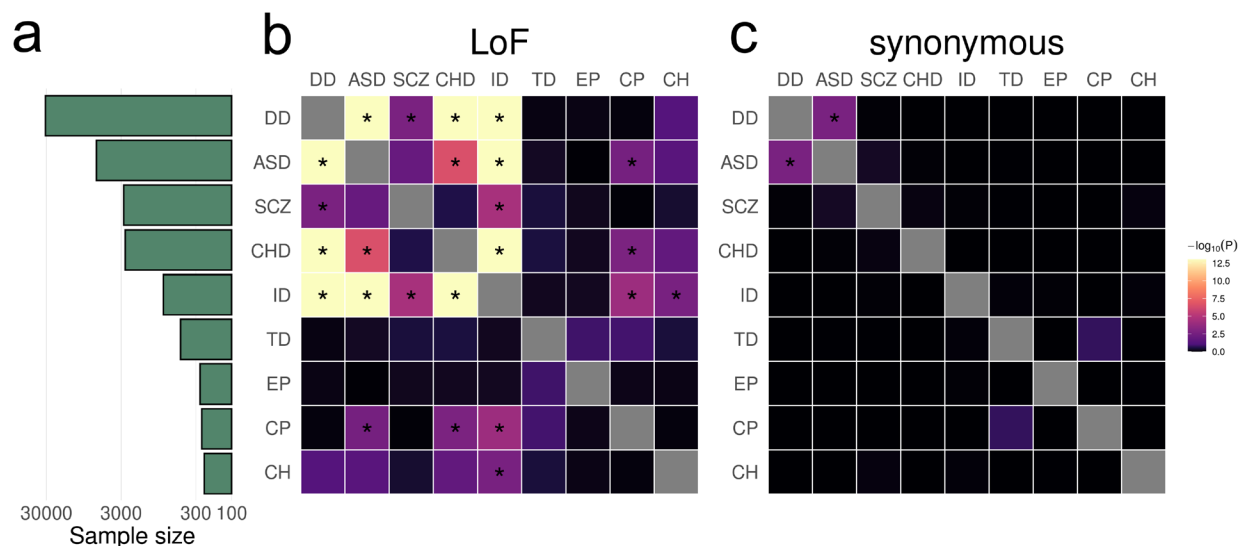
191
 192 **Figure 4. Model fitting results for nine disorders.** (a, b) Estimation results of β and σ for nine disorders and four
 193 variant classes. Error bars represent 1.96*standard errors. (c-f) Distribution of DNM events per gene in four variant
 194 classes for DD. Red and green bars represent the expected frequency of genes under the fixed-effects and mixed-
 195 effects Poisson regression models, respectively. Blue bars represent the observed frequency of genes.
 196

197 We first performed single-trait analysis for each disorder. The estimated elevation parameters (i.e.,
 198 β) were negative for almost all disorders and variant classes (**Figure 4a**), with LoF variants
 199 showing particularly lower parameter estimates. This may be explained by more stringent quality
 200 control in LoF variant calling¹² and potential survival bias³⁰. It is also consistent with a depletion
 201 of LoF DNMs in healthy control trios⁷. The dispersion parameter estimates (i.e., σ) were higher

202 for LoF variants than other variant classes (**Figure 4b**), which is consistent with our expectation
 203 that LoF variants have stronger effects on disease risk and should show a larger deviation from
 204 the null mutation rate in disease probands. We compared the goodness of fit of our proposed
 205 mixed-effects Poisson model to a simpler fixed-effects model without the deviation component
 206 (**Methods**). The expected distribution of recurrent DNM counts showed substantial and
 207 statistically significant improvement under the mixed-effects Poisson model (**Figures 4c-f** and
 208 **Supplementary Figure 3**).

209
 210 Next, we estimated pairwise DNM enrichment correlations for 9 disorders. In total, we identified
 211 25 pairs of disorders with significant correlations at a false discovery rate (FDR) cutoff of 0.05
 212 (**Figure 5** and **Supplementary Figure 4**), including 12 significant correlations for LoF variants, 7
 213 for Dmis variants, 5 for Tmis variants, and only 1 significant correlation for synonymous variants.
 214 Notably, all significant correlations are positive (**Supplementary Table 2**). No significant
 215 correlation was identified between any disorder and healthy controls (**Supplementary Figure 5**).
 216 The number of identified significant correlations for each disorder was proportional to the sample
 217 size in each study (Spearman correlation = 0.70) with controls being a notable outlier
 218 (**Supplementary Figure 6**).

219



220
 221 **Figure 5. EncoreDNM identifies pervasive enrichment correlations of damaging DNMs among nine disorders.**
 222 (a) shows sample size (i.e., number of trios) for each disease. X-axis denotes sample size on the log scale. (b, c)
 223 Heatmap of enrichment correlations for LoF and synonymous DNMs among nine disorders. Significant correlations
 224 (FDR<0.05) are marked by asterisks. Results with $-\log_{10} P > 13$ are truncated to 13 for visualization purpose.

225
 226 We identified highly concordant and significant LoF DNM enrichment among DD, ASD, ID, and
 227 CHD, which is consistent with previous reports^{8-10,31}. SCZ shows highly significant LoF
 228 correlations with DD and ID ($p=2.0e-3$ and $3.7e-5$), hints at a correlation with ASD ($p=0.012$), but
 229 does not correlate strongly with CHD. The positive enrichment correlation between ASD and CP
 230 in LoF variants ($\rho=0.81$, $p=3.3e-3$) is consistent with their co-occurrence³². The high enrichment
 231 correlation between ID and CP in LoF variants ($\rho=0.68$, $p=1.0e-4$) is consistent with the

232 associations between ID and motor or non-motor abnormalities caused by CP³³. A previous study
233 also suggested significant genetic sharing of ID and CP by overlapping genes harboring rare
234 damaging variants¹⁴. Here, we obtained consistent results after accounting for *de novo*
235 mutabilities and potential confounding bias.

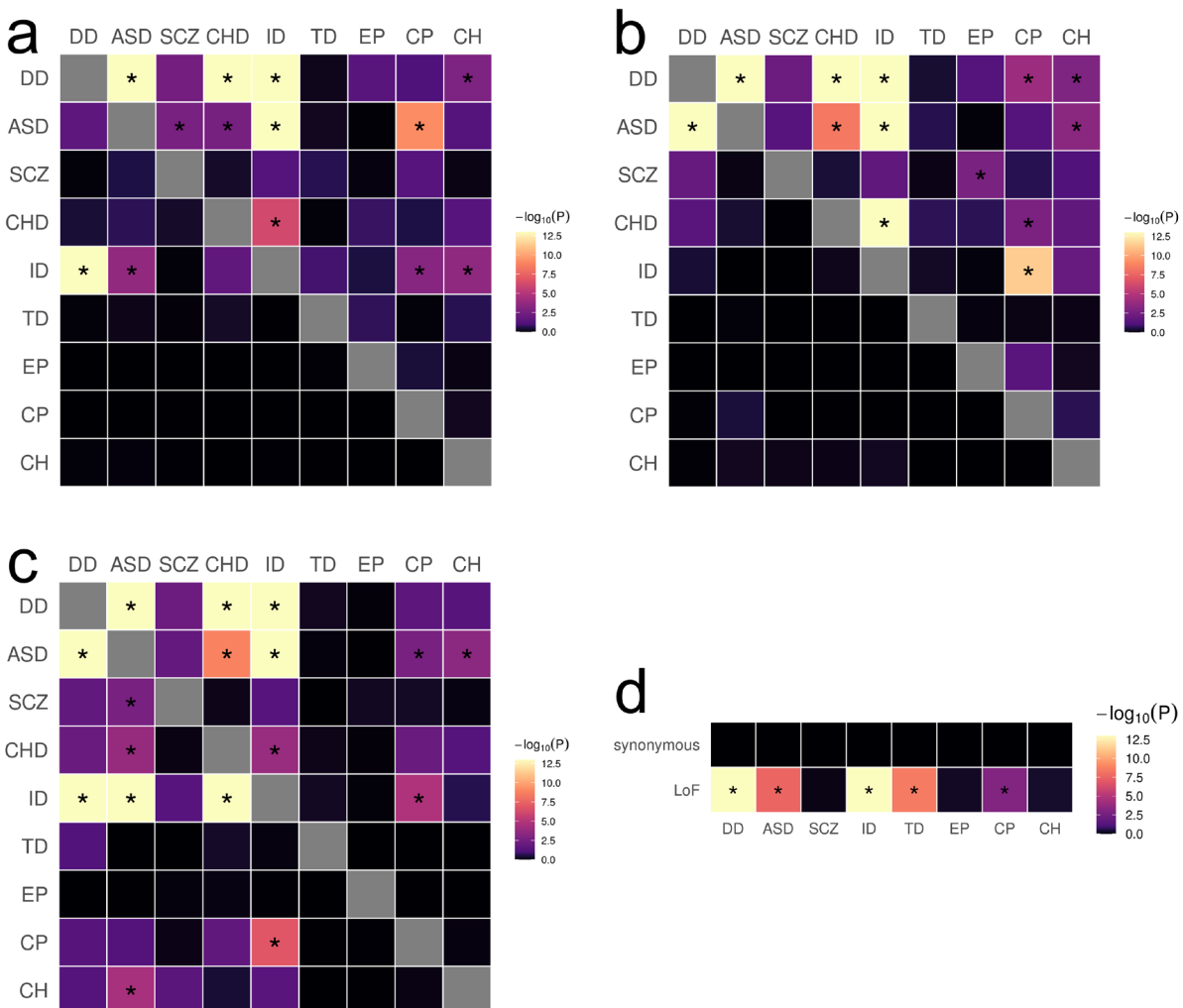
236
237 Some significant correlations identified in our analysis are consistent with phenotypic associations
238 in epidemiological studies, but have not been reported using genetic data to the extent of our
239 knowledge. For example, the LoF enrichment correlation between CHD and CP ($\rho=0.88$, $p=1.7e-$
240 3) is consistent with findings that reduced supply of oxygenated blood in fetal brain due to cardiac
241 malformations may be a risk factor for CP³⁴. The enrichment correlation between ID and CH in
242 LoF variants ($\rho=0.63$, $p=2.4e-3$) is consistent with lower intellectual performance in a proportion
243 of children with CH³⁵.

244
245 Genes showing pathogenic DNMs in multiple disorders may shed light on the mechanisms
246 underlying enrichment correlations (**Supplementary Table 3**). We identified five genes (i.e.
247 *CTNNB1*, *NBEA*, *POGZ*, *SPRED2*, and *KMT2C*) with LoF DNMs in five different disorders and
248 21 genes had LoF DNMs in four disorders (**Supplementary Table 4**). These 26 genes with LoF
249 variants in at least four disorders were significantly enriched for 63 gene ontology (GO) terms with
250 FDR<0.05 (**Supplementary Table 5**). Chromatin organization ($p=7.8e-11$), nucleoplasm ($p=2.8e-$
251 10), chromosome organization ($p=6.8e-10$), histone methyltransferase complex ($p=1.4e-9$), and
252 positive regulation of gene expression ($p=2.2e-9$) were the most significantly enriched GO terms.
253 One notable example consistently included in these gene sets is *CTNNB1* (**Supplementary**
254 **Figure 7**). It encodes β -catenin, is one of the only two genes reaching genome-wide significance
255 in a recent WES study for CP¹⁴, and also harbors multiple LoF variants in DD, ID, ASD, and CHD.
256 It is a fundamental component of the canonical Wnt signaling pathway which is known to confer
257 genetic risk for ASD³⁶. We also identified 157 recurrent LoF mutations in 45 genes
258 (**Supplementary Table 6**). Most of these recurrent mutations were identified in DD due to its large
259 sample size, but one mutation was identified in joint comparison of other disorders. *FBXO11*,
260 encoding the F-box only protein 31, shows two recurrent p.Ser831fs LoF variants in ASD and CH
261 (**Supplementary Figure 8**; $p=1.9e-3$; **Methods**). The F-box protein constitutes a substrate-
262 recognition component of the SCF (SKP1-cullin-F-box) complex, an E3-ubiquitin ligase complex
263 responsible for ubiquitination and proteasomal degradation³⁷. DNMs in *FBXO11* have been
264 previously implicated in severe ID individuals with autistic behavior problem³⁸ and
265 neurodevelopmental disorder³⁹.

266
267 For comparison, we also applied mTADA to the same nine disorders and control trios. In total,
268 mTADA identified 117 disorder pairs with significant genetic sharings at an FDR cutoff of 0.05
269 (**Supplementary Table 7** and **Supplementary Figure 9**). Notably, we identified significant
270 synonymous DNM correlations for all 36 disorder pairs and between all disorders and healthy
271 controls (**Supplementary Figure 10**). These results are consistent with the simulation results and
272 suggest a substantially inflated false positive rate in mTADA.

273
274
275
276
277
278
279
280
281

Further, we applied cross-trait linkage disequilibrium (LD) score regression¹⁸ to five of the nine disorders with publicly available GWAS summary statistics (**Supplementary Table 8**): ASD (n=46,350)⁴⁰, SCZ (n=161,405)⁴¹, cognitive performance (used as a proxy for ID; n=257,841)⁴², TD (n=14,307)⁴³, and epilepsy (n=44,889)⁴⁴. In total, we identified 6 trait pairs with significant genetic correlations at an FDR cutoff of 0.05 (**Supplementary Table 9**), suggesting consistent findings made from GWAS and DNM data (Spearman correlation = 0.70; **Supplementary Figure 11**).



282
283
284
285
286
287
288
289
290
291

Figure 6. DNM enrichment correlations in disease-relevant gene sets. (a) Enrichment correlations in High-pLI genes (upper triangle) and Low-pLI genes (lower triangle) for LoF variants. (b) Enrichment correlations in HBE genes (upper triangle) and LBE genes (lower triangle) for LoF variants. (c) Enrichment correlations in HHE genes (upper triangle) and LHE genes (lower triangle) for LoF variants. (d) Enrichment correlations in CHD-related pathways for LoF and synonymous variants. Significant correlations (FDR<0.05) are marked by asterisks. Results with $-\log_{10} P > 13$ are truncated to 13 for visualization purpose.

292 **Partitioning DNM enrichment correlation by gene set**

293
294 To gain biological insights into the shared genetic architecture of nine disorders, we repeated
295 EncoreDNM correlation analysis in several gene sets. First, we defined genes with high/low
296 probability of intolerance to LoF variants using pLI scores⁴⁵, and identified genes with high/low
297 brain expression (HBE/LBE)⁴⁶ (**Methods; Supplementary Table 10**). We identified 11 and 12
298 disorder pairs showing significant enrichment correlations for LoF DNMs in high-pLI genes and
299 HBE genes, respectively (**Figure 6a-b**). We observed fewer significant correlations for Dmis and
300 Tmis variants in these gene sets (**Supplementary Figures 12-13**). All identified significant
301 correlations were positive (**Supplementary Tables 11-12**). No significant correlations were
302 identified for synonymous variants (**Supplementary Figures 12-13**) or between disorders and
303 controls (**Supplementary Figures 14-15**).

304
305 We observed a clear enrichment of significant correlations in disease-relevant gene sets. Overall,
306 high-pLI genes showed substantially stronger correlations across disorders than genes with low
307 pLI (one-sided Kolmogorov-Smirnov test; $p=2.3e-6$). Similarly, enrichment correlations were
308 stronger in HBE genes than in LBE genes ($p=8.8e-7$). Among the 11 disorder pairs showing
309 significant enrichment correlations in high-pLI genes, two pairs, i.e., ASD-SCZ ($\rho=0.68$, $p=2.4e-$
310 3) and DD-CH ($\rho=0.43$, $p=1.5e-3$), were not identified in the exome-wide analysis. We also
311 identified four novel disorder pairs with significant correlations in HBE genes, including DD-CP
312 ($\rho=0.80$, $p=9.5e-5$), DD-CH ($\rho=0.67$, $p=1.4e-3$), ASD-CH ($\rho=0.82$, $p=4.7e-4$), and SCZ-EP
313 ($\rho=0.66$, $p=2.0e-3$). These novel enrichment correlations are consistent with known comorbidities
314 between these disorders^{47,48} and findings based on significant risk genes^{8,28,49,50}.

315
316 Furthermore, we estimated DNM enrichment correlations in genes with high/low expression in
317 mouse developing heart (HHE/LHE)⁷ (**Methods; Supplementary Table 10**). We identified 9
318 significant enrichment correlations for LoF variants in HHE genes (**Figure 6c**). Strength of
319 enrichment correlations did not show a significant difference between HHE and LHE genes
320 ($p=0.846$), possibly due to a lack of cardiac disorders in our analysis. Finally, we estimated
321 enrichment correlations between CHD and other disorders in known pathways for CHD⁵¹
322 (**Methods; Supplementary Table 10**). We identified 5 significant correlations for LoF variants
323 (**Figure 6d**), including a novel correlation between CHD and TD ($\rho=0.93$, $p=3.3e-9$). Of note,
324 arrhythmia caused by CHD is a known risk factor for TD⁵². In these analyses, all significant
325 enrichment correlations were positive (**Supplementary Tables 13-14**) and other variant classes
326 showed generally weaker correlations than LoF variants (**Supplementary Figures 16-17**). We
327 did not observe significant correlations in these gene sets between disorders and controls
328 (**Supplementary Figures 18-19**).

329
330
331

332 Discussion

333

334 In this paper, we introduced EncoreDNM, a novel statistical framework to quantify correlated DNM
335 enrichment between two disorders. Through extensive simulations and analyses of DNM data for
336 nine disorders, we demonstrated that our proposed mixed-effects Poisson regression model
337 provides unbiased parameter estimates, shows well-controlled type-I error, and is robust to
338 exome-wide technical biases. Leveraging exome-wide DNM counts and genomic context-based
339 mutability data, EncoreDNM achieves superior fit for real DNM datasets compared to simpler
340 models and provides statistically powerful and computationally efficient estimation of DNM
341 enrichment correlation. Further, EncoreDNM can quantify concordant genetic effects for user-
342 defined variant classes within pre-specified gene sets, thus is suitable for exploring diverse types
343 of hypotheses and can provide crucial biological insights into the shared genetic etiology in
344 multiple disorders.

345

346 Multi-trait analyses of GWAS data have revealed shared genetic architecture among many
347 neuropsychiatric traits^{22,53,54}. These findings have led to the identification of pleiotropic variants,
348 genes, and hub genomic regions underlying many traits and have revealed multiple
349 psychopathological factors jointly affecting human neurological phenotypes^{55,56}. Although
350 emerging evidence suggests that causal DNMs underlying several disorders with well-powered
351 studies (e.g., CHD and neurodevelopmental disorders⁷) may be shared, our understanding of the
352 extent and the mechanism underlying such sharing remains incomplete. Applied to DNM data for
353 nine disorders, EncoreDNM identified pervasive enrichment correlations of DNMs. We observed
354 particularly strong correlations in pathogenic variant classes (e.g., LoF and Dmis variants) and
355 disease-relevant genes (e.g., genes with high pLI and genes highly expressed in relevant tissues).
356 Genes underlying these correlations were significantly enriched in pathways involved in chromatin
357 organization and modification and gene expression regulation. The DNM correlations were
358 substantially attenuated in genes with lower expression and genes with frequent occurrences of
359 LoF variants in the population. A similar attenuation was observed in less pathogenic variant
360 classes (e.g., synonymous variants). Further, no significant correlations were identified between
361 any disorder and healthy controls. These results lay the groundwork for future investigations of
362 pleiotropic mechanisms of DNMs.

363

364 Our study has some limitations. First, a main goal in DNM research is to identify disease risk
365 genes. EncoreDNM leverages exome-wide DNM counts to quantify shared genetic basis in
366 multiple disorders but does not improve the analysis of gene-disease associations. Second,
367 EncoreDNM assumes probands from different input studies to be independent. In rare cases
368 when two studies have overlapping proband samples, enrichment correlation estimates may be
369 inflated and must be interpreted with caution. Finally, genetic correlation methods based on
370 GWAS summary data provided key motivations for the mixed-effects Poisson regression model
371 in our study. Built upon genetic correlations, a plethora of methods have been developed in the

372 GWAS literature to jointly model more than two GWAS⁵⁷, identify and quantify common factors
373 underlying multiple traits^{58,59}, estimate causal effects among different traits⁶⁰, and identify
374 pleiotropic genomic regions through hypothesis-free scans²³. Future directions of EncoreDNM
375 include using enrichment correlation to improve gene discovery, learning the directional effects
376 and the causal structure underlying multiple disorders, and dynamically searching for gene sets
377 and annotation classes with shared genetic effects without pre-specifying the hypothesis.

378

379 Taken together, we provide a new analytic approach to an important problem in DNM studies. We
380 believe EncoreDNM improves the statistical rigor in multi-disorder DNM modeling and opens up
381 many interesting future directions in both method development and follow-up analyses in WES
382 studies. As trio sample size in WES studies continues to grow, EncoreDNM will have broad
383 applications and can greatly benefit DNM research.

384

385

386 **Methods**

387

388 **Statistical Model**

389

390 For a single study, we assume that DNM counts in a given variant class (e.g., synonymous
391 variants) follow a mixed-effects Poisson model:

392

$$Y_i \sim \text{Poisson}(\lambda_i),$$

393

$$\log(\lambda_i) = \beta + \log(2Nm_i) + \phi_i,$$

394

$$\phi_i \sim N(0, \sigma^2), \quad \text{for } i = 1, \dots, G,$$

395

396 where Y_i is the DNM count in the i -th gene, N is the number of trios, m_i is the *de novo*
397 mutability for the i -th gene (i.e., mutation rate per chromosome per generation) which is known a
398 *priori*²⁶ (**Supplementary Table 15**), and G is the total number of genes in the study. The elevation
399 parameter β quantifies the global elevation of mutation rate compared to mutability estimates
400 based on genomic sequence alone. Gene-specific deviation from expected DNM rate is quantified
401 by random effect ϕ_i with a dispersion parameter σ . Here, the ϕ_i are assumed to be
402 independent across different genes, in which case the observed DNM counts of different genes
403 are independent.

403

404 Next, we describe how we expand this model to quantify the shared genetics of two disorders.

405 We assume DNM counts in a given variant class for two diseases follow:

406

$$\begin{bmatrix} Y_{i1} \\ Y_{i2} \end{bmatrix} \sim \text{Poisson} \left(\begin{bmatrix} \lambda_{i1} \\ \lambda_{i2} \end{bmatrix} \right),$$

407

$$\log \left(\begin{bmatrix} \lambda_{i1} \\ \lambda_{i2} \end{bmatrix} \right) = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} + \log \left(\begin{bmatrix} 2N_1 m_i \\ 2N_2 m_i \end{bmatrix} \right) + \begin{bmatrix} \phi_{i1} \\ \phi_{i2} \end{bmatrix},$$

408

$$\begin{bmatrix} \phi_{i1} \\ \phi_{i2} \end{bmatrix} \sim \text{MVN} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix} \right),$$

409 where Y_{i1}, Y_{i2} are the DNM counts for the i -th gene and N_1, N_2 are the trio sizes in two studies,
 410 respectively. Similar to the single-trait model, m_i is the mutability for the i -th gene. β_1, β_2 are
 411 the elevation parameters, and ϕ_{i1}, ϕ_{i2} are the gene-specific random effects with dispersion
 412 parameters σ_1, σ_2 , for two disorders respectively. ρ is the enrichment correlation which quantifies
 413 the concordance of the gene-specific DNM burden between two disorders. Here, $\beta_1, \beta_2, \sigma_1, \sigma_2, \rho$
 414 are unknown parameters. The gene specific effects for two disorders are assumed to be
 415 independent for different genes. We also assume that there is no shared sample for two disorders,
 416 in which case Y_{i1} is independent with Y_{i2} given $\begin{bmatrix} \lambda_{i1} \\ \lambda_{i2} \end{bmatrix}$.

417

418 Parameter estimation

419

420 We implement an MLE procedure to estimate unknown parameters. For single-trait analysis, the
 421 log-likelihood function can be expressed as follows:

$$422 \quad l(\beta, \sigma | \mathbf{Y}) = \sum_{i=1}^G \log \left[\int \exp(-\lambda_i) \lambda_i^{Y_i} * f(\phi_i) d\phi_i \right] + C,$$

423 where $\mathbf{Y} = [Y_1, \dots, Y_G]^T$, $\lambda_i = 2Nm_i \exp(\beta + \phi_i)$, $C = -\sum_{i=1}^G \log(Y_i!)$, and $f(\phi_i) =$
 424 $\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\phi_i^2}{2\sigma^2}\right)$. Note that there is no closed form for the integral in the log-likelihood function.

425 Therefore, we use Monte Carlo integration to evaluate the log-likelihood function. Let $\phi_{ij} = \sigma \xi_{ij}$,
 426 where the ξ_{ij} are independently and identically distributed random variables following a standard
 427 normal distribution. We have

$$428 \quad l(\beta, \sigma | \mathbf{Y}) \approx l'(\beta, \sigma | \mathbf{Y}) = \sum_{i=1}^G \log \left[\sum_{j=1}^M \exp(-\lambda_{ij}) \lambda_{ij}^{Y_i} \right] + C,$$

429 where $\lambda_{ij} = 2Nm_i \exp(\beta + \sigma \xi_{ij})$, and M is the Monte Carlo sample size which is set to be 1000.
 430 Then, we could obtain the MLE of β, σ through maximization of $l'(\beta, \sigma | \mathbf{Y})$. We obtain the
 431 standard error of the MLE through inversion of the observed Fisher information matrix.

432

433 The estimation procedure can be generalized to multi-trait analysis. Log-likelihood function can
 434 be expressed as follows:

$$435 \quad l(\beta_1, \beta_2, \sigma_1, \sigma_2, \rho | \mathbf{Y}_1, \mathbf{Y}_2) = \sum_{i=1}^G \log \left[\int \exp(-\lambda_{i1} - \lambda_{i2}) \lambda_{i1}^{Y_{i1}} \lambda_{i2}^{Y_{i2}} * f(\phi_{i1}, \phi_{i2}) d\phi_{i1} d\phi_{i2} \right] + C,$$

436 where $\mathbf{Y}_1 = [Y_{11}, \dots, Y_{G1}]^T$, $\mathbf{Y}_2 = [Y_{12}, \dots, Y_{G2}]^T$, $\lambda_{i1} = 2N_1 m_i \exp(\beta_1 + \phi_{i1})$, $\lambda_{i2} = 2N_2 m_i \exp(\beta_2 +$
 437 $\phi_{i2})$, $C = -\sum_{i=1}^G [\log(Y_{i1}!) + \log(Y_{i2}!)]$, and $f(\phi_{i1}, \phi_{i2}) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left[-\frac{1}{2\sqrt{1-\rho^2}} \left(\frac{\phi_{i1}^2}{\sigma_1^2} + \frac{\phi_{i2}^2}{\sigma_2^2} -$
 438 $\frac{2\rho\phi_{i1}\phi_{i2}}{\sigma_1\sigma_2}\right)\right]$. We use Monte Carlo integration to evaluate the log-likelihood function. Let $\phi_{i1j} =$

439 $\sigma_1 \xi_{i1j}$ and $\phi_{i2j} = \sigma_2 (\rho \xi_{i1j} + \sqrt{1 - \rho^2} \xi_{i2j})$, where the ξ_{i1j} and ξ_{i2j} are independently and
440 identically distributed random variables following a standard normal distribution. We have

$$441 \quad l(\beta_1, \beta_2, \sigma_1, \sigma_2, \rho | \mathbf{Y}_1, \mathbf{Y}_2) \approx l'(\beta_1, \beta_2, \sigma_1, \sigma_2, \rho | \mathbf{Y}_1, \mathbf{Y}_2) = \sum_{i=1}^G \log \left[\sum_{j=1}^M \exp(-\lambda_{i1j} - \lambda_{i2j}) \lambda_{i1j}^{Y_{i1j}} \lambda_{i2j}^{Y_{i2j}} \right] + C,$$

442 where $\lambda_{i1j} = 2N_1 m_i \exp(\beta_1 + \sigma_1 \xi_{i1j})$ and $\lambda_{i2j} = 2N_2 m_i \exp[\beta_2 + \sigma_2 (\rho \xi_{i1j} + \sqrt{1 - \rho^2} \xi_{i2j})]$. Then,
443 we obtain the MLE of $\beta_1, \beta_2, \sigma_1, \sigma_2, \rho$ through maximization of $l'(\beta_1, \beta_2, \sigma_1, \sigma_2, \rho | \mathbf{Y}_1, \mathbf{Y}_2)$. Standard
444 error of MLE can be obtained through inversion of the observed Fisher information matrix.

445

446 **DNM data and variant annotation**

447

448 We obtained DNM data from published studies (**Supplementary Table 1**). DNM data for EP from
449 the original release¹³ were not in an editable format and were instead collected from denovo-db⁶¹.
450 We used ANNOVAR⁶² to annotate all DNMs. Synonymous variants were determined based on
451 the 'synonymous SNV' annotation in ANNOVAR; Variants with 'startloss', 'stopgain', 'stoploss',
452 'splicing', 'frameshift insertion', 'frameshift deletion', or 'frameshift substitution' annotations were
453 classified as LoF; Dmis variants were defined as nonsynonymous SNVs predicted to be
454 deleterious by MetaSVM⁶³; nonsynonymous SNVs predicted to be tolerable by MetaSVM were
455 classified as Tmis. Other DNMs which did not fall into these categories were removed from the
456 analysis. For each variant class, we estimated the mutability of each gene using a sequence-
457 based mutation model²⁶ while adjusting for the sequencing coverage factor based on control trios
458 as previously described¹² (**Supplementary Table 15**). We included 18,454 autosomal protein-
459 coding genes in our analysis. *TTN* was removed due to its substantially larger size.

460

461 **Implementation of mTADA**

462

463 The software mTADA requires the following parameters as inputs: proportion of risk genes (π_1^S, π_2^S),
464 mean relative risks ($\bar{\gamma}_1^S, \bar{\gamma}_2^S$), and dispersion parameters ($\bar{\beta}_1^S, \bar{\beta}_2^S$) for both disorders. We used
465 extTADA¹⁰ to estimate these parameters as suggested by the mTADA paper⁹. mTADA reported
466 the estimated proportion of shared risk genes π_3 (posterior mode of π_3) and its corresponding
467 95% credible interval [LB, UB]. We considered $\pi_1^S * \pi_2^S$ as the expected proportion of shared risk
468 genes, and there is significant genetic sharing between two disorders when $LB > \pi_1^S * \pi_2^S$. P-
469 value for π_3 was calculated by comparing $\pi_1^S * \pi_2^S$ to the posterior distribution of π_3 . Number of
470 MCMC chain was set as 2 and number of iterations was set as 10,000.

471

472 **Simulation settings**

473

474 We assessed the performance of EncoreDNM under the mixed-effects Poisson model. We
475 performed simulations for two variant classes: Tmis and LoF variants, which have the largest and

476 the smallest median mutability values across all genes. First, we performed single-trait
477 simulations to assess estimation precision of elevation parameter β and dispersion parameter
478 σ . We set the true values of β to be -0.5, -0.25, and 0, and the true values of σ to be 0.5, 0.75,
479 and 1. These values were chosen based on the estimated parameters in real DNM data analyses
480 and ensured simulation settings to be realistic. Next, we performed simulations for cross-trait
481 analysis to assess estimation precision of enrichment correlation ρ , whose true values were set
482 to be 0, 0.2, 0.4, 0.6, and 0.8. Sample size for each disorder was set to be 5,000. Coverage rate
483 was calculated as the percentage of simulations that the 95% Wald confidence interval covered
484 the true parameter value. Each parameter setting was repeated 100 times.

485
486 We also carried out simulations to compare the performance of EncoreDNM and mTADA. Type I
487 error and statistical power for EncoreDNM were calculated as the proportion of simulation repeats
488 that p-value for enrichment correlation ρ was smaller than 0.05. and the proportion of simulation
489 repeats that p-value for estimated proportion of shared risk genes π_3 was smaller than 0.05 was
490 used for mTADA. We aggregated all variant classes together, so mutability for each gene was
491 determined as the sum of mutabilities across four variant classes (i.e. LoF, Dmis, Tmis, and
492 synonymous).

493
494 First, we simulated DNM data under the mixed-effects Poisson model. To see whether two
495 methods would produce false positive findings, we performed simulations under the null
496 hypothesis that the enrichment correlation ρ is zero. We compared two methods under a range
497 of parameter combinations of (β, σ, N) for both disorders: (-0.25, 0.75, 5000) for the baseline
498 setting, (-1, 0.75, 5000) for a setting with small β , (-0.25, 0.5, 5000) for a setting with small σ ,
499 and (-0.25, 0.75, 1000) for a setting with small sample size. We also assessed the statistical
500 power of two methods under the alternative hypothesis. True value of enrichment correlation ρ
501 was set to be 0.05, 0.1, 0.15, and 0.2. In the power analysis, parameters (β, σ, N) were fixed at (-
502 0.25, 0.75, 5000) as in the baseline setting when both methods had well-controlled type-I error.

503
504 To ensure a fair comparison, we also compared EncoreDNM and mTADA under a multinomial
505 model, which is different from the data generation processes for the two approaches. For each
506 disorder ($k = 1, 2$), we randomly selected causal genes of proportion π_k^S . A proportion (i.e., π_3) of
507 causal genes overlap between two disorders. We assumed that the total DNM count to follow a
508 Poisson distribution: $C_k \sim \text{Poisson}(u_k * 2N_k \sum_{i=1}^G m_i)$, where u_k represents an elevation factor to
509 represent systematic bias in the data. Let \mathbf{Y}_k denote the vector of DNMs counts in the exome,
510 \mathbf{m} denote the vector of mutability values for all genes, and $\mathbf{m}_{causal,k}$ denote the vector of
511 mutability with values set to be 0 for non-causal genes of disorder k . We assumed that a
512 proportion p_k of the probands could be attributed to DNMs burden in causal genes, and $1 - p_k$
513 of the probands obtained DNMs by chance:

$$\begin{aligned} \mathbf{Y}_k &= \mathbf{Y}_{causal,k} + \mathbf{Y}_{background,k}, \\ \mathbf{Y}_{causal,k} &\sim \text{Multinomial}(p_k C_k, \mathbf{m}_{causal,k}), \\ \mathbf{Y}_{background,k} &\sim \text{Multinomial}((1 - p_k) C_k, \mathbf{m}). \end{aligned}$$

517 To check whether false positive findings could arise, we performed simulations under the null
518 hypothesis that $\pi_3 = \pi_1^S * \pi_2^S$ across a range of parameter combinations of (u, p, N, π^S) for both
519 disorders: (0.95, 0.25, 5000, 0.1) for the baseline setting, (0.75, 0.25, 5000, 0.1) for a setting with
520 small u (i.e., reduced total mutation count), (0.95, 0.15, 5000, 0.1) for a setting with small p
521 (fewer probands explained by DNMs), and (0.95, 0.25, 1000, 0.1) for a setting with smaller sample
522 size. We also assessed the statistical power of two methods under the alternative hypothesis that
523 $\pi_3 > \pi_1^S * \pi_2^S$. In power analysis, (u, p, N, π^S) were fixed at (0.95, 0.25, 5000, 0.1) as in the baseline
524 setting when type-I error for both methods were well-calibrated.

525

526 **Comparison to the fixed-effects Poisson model**

527

528 For single-trait analysis, the fixed-effects Poisson model assumes that

$$529 \quad Y_i \sim \text{Poisson}(\lambda_i),$$
$$530 \quad \log(\lambda_i) = \beta + \log(2Nm_i), \quad \text{for } i = 1, \dots, G.$$

531 Note that the fixed-effects Poisson model is a special case of our proposed mixed-effects Poisson
532 model when $\sigma = 0$. We compared the two models using likelihood ratio test. Under the null
533 hypothesis that $\sigma = 0$, $2(l_{alt} - l_{null}) \sim \frac{1}{2}\chi_1^2$ asymptotically, where l_{alt} and l_{null} represent the
534 log likelihood of the fitted mixed-effects and fixed-effects Poisson models respectively.

535

536 **Recurrent genes and DNMs**

537

538 We used FUMA⁶⁴ to perform GO enrichment analysis for genes harboring LoF DNMs in multiple
539 disorders. Due to potential sample overlap between the studies of DD³ and ID², we excluded ID
540 from the analysis of recurrent DNMs. We calculated the probability of observing two identical
541 DNMs in two disorders using a Monte Carlo simulation method. For each disorder, we simulated
542 exome-wide DNMs profile from a multinomial distribution, where the size was fixed at the
543 observed DNM count and the per-base mutation probability was determined by the tri-nucleotide
544 base context. We repeated the simulation procedure 100,000 times to evaluate the significance
545 of recurrent DNMs. Lollipop plots for recurrent mutations were generated using MutationMapper
546 on the cBio Cancer Genomics Portal⁶⁵.

547

548 **Implementation of cross-trait LD score regression**

549

550 We used cross-trait LDSC¹⁸ to estimate genetic correlations between disorders. LD scores were
551 computed using European samples from the 1000 Genomes Project Phase 3 data⁶⁶. Only
552 HapMap 3 SNPs were used as observations in the explanatory variable with the `--merge-alleles`
553 flag. Intercepts were not constrained in the analyses.

554

555 **Estimating enrichment correlation in gene sets**

556

557 Genes with a high/low probability of intolerance to LoF variants (high-pLI/low-pLI) were defined
558 as the 4,614 genes in the upper/lower quartiles of pLI scores⁴⁵. Genes with high/low brain
559 expression (HBE/LBE) were defined as the 4,614 genes in the upper/lower quartiles of expression
560 in the human fetal brain⁴⁶. Genes with high/low heart expression (HHE/LHE) were defined as the
561 4,614 genes in the upper/lower quartiles of expression in the developing heart of embryonic
562 mouse⁶⁷. Five biological pathways have been reported to be involved in CHD: chromatin
563 remodeling, Notch signaling, cilia function, sarcomere structure and function, and RAS signaling⁵¹.
564 We extracted 1730 unique genes that belong to these five pathways from the gene ontology
565 database⁶⁸ and referred to the union set as CHD-related genes. We repeated EncoreDNM
566 enrichment correlation analysis in these gene sets. One-sided Kolmogorov-Smirnov test was
567 used to assess the statistical difference between enrichment correlation signal strength in different
568 gene sets.

569

570 **URLs**

571

572 GWAS summary statistics data of ASD, SCZ, and TD were downloaded on the PGC website,
573 <https://www.med.unc.edu/pgc/download-results/>; Summary statistics of cognitive performance
574 were downloaded on the SSGAC website, <https://www.thessgac.org/data>; Summary statistics of
575 epilepsy were downloaded on the epiGAD website, <http://www.epigad.org/>; pLI scores were
576 downloaded from gnomAD v3.1 repository <https://gnomad.broadinstitute.org/downloads>; mTADA,
577 <https://github.com/hoangtn/mTADA>; denovo-db, [https://denovo-db.gs.washington.edu/denovo-](https://denovo-db.gs.washington.edu/denovo-db/)
578 [db/](https://denovo-db.gs.washington.edu/denovo-db/); MutationMapper on cBioPortal, https://www.cbioportal.org/mutation_mapper; LDSC,
579 <https://github.com/bulik/ldsc>.

580

581

582 **Code availability**

583

584 EncoreDNM software is available at <https://github.com/ghm17/EncoreDNM>.

585

586 **Acknowledgements**

587 LH acknowledges research support from the National Science Foundation of China (Grant No.
588 12071243) and Shanghai Municipal Science and Technology Major Project (Grant No.
589 2017SHZDZX01). QL acknowledges research support from the University of Wisconsin-Madison
590 Office of the Chancellor and the Vice Chancellor for Research and Graduate Education with
591 funding from the Wisconsin Alumni Research Foundation and the Waisman Center pilot grant
592 program at University of Wisconsin-Madison. HZ acknowledges research support from the
593 National Institutes of Health (Grant No. R03HD100883)
594

595 **Author contribution**

596 H.G., L.H., and Q.L. designed the study.
597 H.G. performed data analysis and implemented the software.
598 Y.S. implemented an early version of the method.
599 S.C.J., X.Z., and B.L. assisted DNM and mutability data preparation.
600 R.P.L and M.B. advised on disease biology, data interpretation, and genetic issues.
601 H.Z. and Q.L. advised on statistical issues.
602 H.G., L.H., and Q.L. wrote the manuscript.
603 All authors contributed in manuscript editing and approved the manuscript.
604

605 **Competing financial interests**

606 The authors declare no competing financial interests.
607

608 References

- 609 1. Veltman, J.A. & Brunner, H.G. De novo mutations in human genetic disease. *Nature Reviews*
610 *Genetics* **13**, 565-575 (2012).
- 611 2. Lelieveld, S.H. *et al.* Meta-analysis of 2,104 trios provides support for 10 new genes for intellectual
612 disability. *Nature neuroscience* **19**, 1194-1196 (2016).
- 613 3. Kaplanis, J. *et al.* Evidence for 28 genetic disorders discovered by combining healthcare and
614 research data. *Nature* **586**, 757-762 (2020).
- 615 4. Satterstrom, F.K. *et al.* Large-scale exome sequencing study implicates both developmental and
616 functional changes in the neurobiology of autism. *Cell* **180**, 568-584. e23 (2020).
- 617 5. Hoischen, A., Krumm, N. & Eichler, E.E. Prioritization of neurodevelopmental disease genes by
618 discovery of new mutations. *Nature neuroscience* **17**, 764 (2014).
- 619 6. Fromer, M. *et al.* De novo mutations in schizophrenia implicate synaptic networks. *Nature* **506**, 179-
620 184 (2014).
- 621 7. Homsy, J. *et al.* De novo mutations in congenital heart disease with neurodevelopmental and other
622 congenital anomalies. *Science* **350**, 1262-1266 (2015).
- 623 8. Li, J. *et al.* Genes with de novo mutations are shared by four neuropsychiatric disorders discovered
624 from NPdenovo database. *Molecular psychiatry* **21**, 290-297 (2016).
- 625 9. Nguyen, T.-H. *et al.* mTADA is a framework for identifying risk genes from de novo mutations in
626 multiple traits. *Nature Communications* **11**, 2929 (2020).
- 627 10. Nguyen, H.T. *et al.* Integrated Bayesian analysis of rare exonic variants to identify risk genes for
628 schizophrenia and neurodevelopmental disorders. *Genome medicine* **9**, 114 (2017).
- 629 11. Willsey, A.J. *et al.* The psychiatric cell map initiative: a convergent systems biological approach to
630 illuminating key molecular pathways in neuropsychiatric disorders. *Cell* **174**, 505-520 (2018).
- 631 12. Jin, S.C. *et al.* Contribution of rare inherited and de novo variants in 2,871 congenital heart disease
632 probands. *Nature genetics* **49**, 1593 (2017).
- 633 13. Allen, A.S. *et al.* De novo mutations in epileptic encephalopathies. *Nature* **501**, 217-221 (2013).
- 634 14. Jin, S.C. *et al.* Mutations disrupting neuritogenesis genes confer risk for cerebral palsy. *Nature*
635 *Genetics* **52**, 1046-1056 (2020).
- 636 15. Howrigan, D.P. *et al.* Exome sequencing in schizophrenia-affected parent-offspring trios reveals
637 risk conferred by protein-coding de novo mutations. *Nature Neuroscience* **23**, 185-193 (2020).
- 638 16. Zhang, Y. *et al.* Comparison of methods for estimating genetic correlation between complex traits
639 using GWAS summary statistics. *Briefings in bioinformatics* (2021).
- 640 17. Lee, S.H., Yang, J., Goddard, M.E., Visscher, P.M. & Wray, N.R. Estimation of pleiotropy between
641 complex diseases using single-nucleotide polymorphism-derived genomic relationships and
642 restricted maximum likelihood. *Bioinformatics* **28**, 2540-2542 (2012).
- 643 18. Bulik-Sullivan, B. *et al.* An atlas of genetic correlations across human diseases and traits. *Nature*
644 *genetics* **47**, 1236 (2015).
- 645 19. Lu, Q. *et al.* A Powerful Approach to Estimating Annotation-Stratified Genetic Covariance via GWAS
646 Summary Statistics. *Am J Hum Genet* **101**, 939-964 (2017).
- 647 20. Ning, Z., Pawitan, Y. & Shen, X. High-definition likelihood inference of genetic correlations across

- 648 human complex traits. *Nature genetics* **52**, 859-864 (2020).
- 649 21. Shi, H., Mancuso, N., Spendlove, S. & Pasaniuc, B. Local genetic correlation gives insights into the
650 shared genetic architecture of complex traits. *The American Journal of Human Genetics* **101**, 737-
651 751 (2017).
- 652 22. Brainstorm, C. *et al.* Analysis of shared heritability in common disorders of the brain. *Science*
653 **360**(2018).
- 654 23. Guo, H., Li, J.J., Lu, Q. & Hou, L. Detecting local genetic correlations with scan statistics. *Nature*
655 *Communications* **12**, 2033 (2021).
- 656 24. Zhang, Y. *et al.* SUPERGNOVA: local genetic correlation analysis reveals heterogeneous etiologic
657 sharing of complex traits. *Genome biology* **22**, 1-30 (2021).
- 658 25. Wei, Q. *et al.* A Bayesian framework for de novo mutation calling in parents-offspring trios.
659 *Bioinformatics* **31**, 1375-1381 (2015).
- 660 26. Samocha, K.E. *et al.* A framework for the interpretation of de novo mutation in human disease.
661 *Nature genetics* **46**, 944-950 (2014).
- 662 27. Willsey, A.J. *et al.* De novo coding variants are strongly associated with Tourette disorder. *Neuron*
663 **94**, 486-499. e9 (2017).
- 664 28. Jin, S.C. *et al.* Exome sequencing implicates genetic disruption of prenatal neuro-gliogenesis in
665 sporadic congenital hydrocephalus. *Nature medicine* **26**, 1754-1765 (2020).
- 666 29. Krumm, N. *et al.* Excess of rare, inherited truncating mutations in autism. *Nature genetics* **47**, 582-
667 588 (2015).
- 668 30. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285-291
669 (2016).
- 670 31. Hormozdiari, F., Penn, O., Borenstein, E. & Eichler, E.E. The discovery of integrated gene networks
671 for autism and related disorders. *Genome research* **25**, 142-154 (2015).
- 672 32. Christensen, D. *et al.* Prevalence of cerebral palsy, co-occurring autism spectrum disorders, and
673 motor functioning—Autism and Developmental Disabilities Monitoring Network, USA, 2008.
674 *Developmental Medicine & Child Neurology* **56**, 59-65 (2014).
- 675 33. Reid, S.M., Meehan, E.M., Arnup, S.J. & Reddihough, D.S. Intellectual disability in cerebral palsy:
676 a population-based retrospective study. *Developmental Medicine & Child Neurology* **60**, 687-694
677 (2018).
- 678 34. Garne, E. *et al.* Cerebral palsy and congenital malformations. *European Journal of Paediatric*
679 *Neurology* **12**, 82-88 (2008).
- 680 35. Lumenta, C.B. & Skotarczak, U. Long-term follow-up in 233 patients with congenital hydrocephalus.
681 *Child's Nervous System* **11**, 173-175 (1995).
- 682 36. O’Roak, B.J. *et al.* Sporadic autism exomes reveal a highly interconnected protein network of de
683 novo mutations. *Nature* **485**, 246-250 (2012).
- 684 37. Cardozo, T. & Pagano, M. The SCF ubiquitin ligase: insights into a molecular machine. *Nature*
685 *reviews Molecular cell biology* **5**, 739-751 (2004).
- 686 38. Jansen, S. *et al.* De novo variants in FBXO11 cause a syndromic form of intellectual disability with

- 687 behavioral problems and dysmorphisms. *European Journal of Human Genetics* **27**, 738-746 (2019).
- 688 39. Gregor, A. *et al.* De novo variants in the F-box protein FBXO11 in 20 individuals with a variable
689 neurodevelopmental disorder. *The American Journal of Human Genetics* **103**, 305-316 (2018).
- 690 40. Grove, J. *et al.* Identification of common genetic risk variants for autism spectrum disorder. *Nature*
691 *genetics* **51**, 431 (2019).
- 692 41. Ripke, S., Walters, J.T., O'Donovan, M.C. & Consortium, S.W.G.o.t.P.G. Mapping genomic loci
693 prioritises genes and implicates synaptic biology in schizophrenia. *MedRxiv* (2020).
- 694 42. Lee, J.J. *et al.* Gene discovery and polygenic prediction from a genome-wide association study of
695 educational attainment in 1.1 million individuals. *Nature genetics* **50**, 1112-1121 (2018).
- 696 43. Yu, D. *et al.* Interrogating the genetic determinants of Tourette's syndrome and other tic disorders
697 through genome-wide association studies. *American Journal of Psychiatry* **176**, 217-227 (2019).
- 698 44. Abou-Khalil, B. *et al.* Genome-wide mega-analysis identifies 16 loci and highlights diverse
699 biological mechanisms in the common epilepsies. *Nature Communications* **9**, 5269 (2018).
- 700 45. Karczewski, K.J. *et al.* The mutational constraint spectrum quantified from variation in 141,456
701 humans. *Nature* **581**, 434-443 (2020).
- 702 46. Werling, D.M. *et al.* Whole-genome and RNA sequencing reveal variation and transcriptomic
703 coordination in the developing human prefrontal cortex. *Cell reports* **31**, 107489 (2020).
- 704 47. Kielinen, M., Rantala, H., Timonen, E., Linna, S.-L. & Moilanen, I. Associated medical disorders
705 and disabilities in children with autistic disorder: a population-based study. *Autism* **8**, 49-60 (2004).
- 706 48. Kilincaslan, A. & Mukaddes, N.M. Pervasive developmental disorders in individuals with cerebral
707 palsy. *Developmental Medicine & Child Neurology* **51**, 289-294 (2009).
- 708 49. Kume, T. *et al.* The forkhead/winged helix gene Mf1 is disrupted in the pleiotropic mouse mutation
709 congenital hydrocephalus. *Cell* **93**, 985-996 (1998).
- 710 50. Cao, M. & Wu, J.I. Camk2a-Cre-mediated conditional deletion of chromatin remodeler Brg1 causes
711 perinatal hydrocephalus. *Neuroscience letters* **597**, 71-76 (2015).
- 712 51. Zaidi, S. & Brueckner, M. Genetics and genomics of congenital heart disease. *Circulation research*
713 **120**, 923-940 (2017).
- 714 52. Gulisano, M. *et al.* Cardiovascular safety of aripiprazole and pimozide in young patients with
715 Tourette syndrome. *Neurological Sciences* **32**, 1213-1217 (2011).
- 716 53. Lee, S.H. *et al.* Genetic relationship between five psychiatric disorders estimated from genome-
717 wide SNPs. *Nature genetics* **45**, 984 (2013).
- 718 54. Gratten, J., Wray, N.R., Keller, M.C. & Visscher, P.M. Large-scale genomics unveils the genetic
719 architecture of psychiatric disorders. *Nature neuroscience* **17**, 782 (2014).
- 720 55. Lee, P.H. *et al.* Genomic relationships, novel loci, and pleiotropic mechanisms across eight
721 psychiatric disorders. *Cell* **179**, 1469-1482. e11 (2019).
- 722 56. Wang, Q., Yang, C., Gelernter, J. & Zhao, H. Pervasive pleiotropy between psychiatric disorders
723 and immune disorders revealed by integrative analysis of multiple GWAS. *Human genetics* **134**,
724 1195-1209 (2015).
- 725 57. Turley, P. *et al.* Multi-trait analysis of genome-wide association summary statistics using MTAG.
726 *Nature genetics* **50**, 229 (2018).
- 727 58. Grotzinger, A.D. *et al.* Genomic structural equation modelling provides insights into the multivariate

- 728 genetic architecture of complex traits. *Nature human behaviour* **3**, 513 (2019).
- 729 59. Grotzinger, A.D. *et al.* Genetic Architecture of 11 Major Psychiatric Disorders at Biobehavioral,
730 Functional Genomic, and Molecular Genetic Levels of Analysis. *medRxiv* (2020).
- 731 60. Pickrell, J.K. *et al.* Detection and interpretation of shared genetic influences on 42 human traits.
732 *Nature genetics* **48**, 709 (2016).
- 733 61. Turner, T.N. *et al.* denovo-db: A compendium of human de novo variants. *Nucleic acids research*
734 **45**, D804-D811 (2017).
- 735 62. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic variants from high-
736 throughput sequencing data. *Nucleic acids research* **38**, e164-e164 (2010).
- 737 63. Dong, C. *et al.* Comparison and integration of deleteriousness prediction methods for
738 nonsynonymous SNVs in whole exome sequencing studies. *Human molecular genetics* **24**, 2125-
739 2137 (2015).
- 740 64. Watanabe, K., Taskesen, E., Van Bochoven, A. & Posthuma, D. Functional mapping and annotation
741 of genetic associations with FUMA. *Nature communications* **8**, 1-11 (2017).
- 742 65. Cerami, E. *et al.* The cBio cancer genomics portal: an open platform for exploring multidimensional
743 cancer genomics data. (AACR, 2012).
- 744 66. Consortium, G.P. A global reference for human genetic variation. *Nature* **526**, 68-74 (2015).
- 745 67. Zaidi, S. *et al.* De novo mutations in histone-modifying genes in congenital heart disease. *Nature*
746 **498**, 220-223 (2013).
- 747 68. Ashburner, M. *et al.* Gene ontology: tool for the unification of biology. *Nature genetics* **25**, 25-29
748 (2000).
- 749