

Genetic modifiers and ascertainment drive variable expressivity of complex disorders

Matthew Jensen^{1,2*}, Corrine Smolen^{1,2*}, Anastasia Tyryshkina^{1*}, Lucilla Pizzo^{1*}, Deepto Banerjee¹, Matthew Oetjens³, Hermela Shimelis³, Cora M. Taylor³, Vijay Kumar Pounraja^{1,2}, Hyebin Song⁴, Laura Rohan¹, Emily Huber¹, Laila El Khattabi⁵, Ingrid van de Laar⁶, Rafik Tadros⁶, Connie Bezzina⁶, Marjon van Slegtenhorst⁶, Janneke Kammeraad⁶, Paolo Prontera⁷, Jean-Hubert Caberg⁸, Harry Fraser⁹, Siddhartha Banka^{9,10}, Anke Van Dijck¹¹, Charles Schwartz¹², Els Voorhoeve¹³, Patrick Callier¹⁴, Anne-Laure Mosca-Boidron¹⁴, Nathalie Marle¹⁴, Mathilde Lefebvre¹⁵, Kate Pope¹⁶, Penny Snell¹⁶, Amber Boys¹⁶, Paul J. Lockhart^{16,17}, Myla Ashfaq¹⁸, Elizabeth McCready¹⁹, Margaret Nowaczyk¹⁹, Lucia Castiglia²⁰, Ornella Galesi²⁰, Emanuela Avola²⁰, Teresa Mattina²⁰, Marco Fichera^{20,21}, Maria Grazia Bruccheri²⁰, Giuseppa Maria Luana Mandarà²², Francesca Mari²³, Flavia Privitera²³, Ilaria Longo²³, Aurora Currò²³, Alessandra Renieri²³, Boris Keren²⁴, Perrine Charles²⁴, Silvestre Cuinat²⁵, Mathilde Nizon²⁵, Olivier Pichon²⁵, Claire Bénétteau²⁵, Radka Stoeva²⁵, Dominique Martin-Coignard²⁶, Sophia Blesson²⁷, Cedric Le Caignec^{28,29}, Sandra Mercier²⁷, Marie Vincent²⁷, Christa Martin³, Katrin Mannik^{30,31}, Alexandre Reymond³², Laurence Faivre^{14,15}, Erik Sistermans¹³, R. Frank Kooy¹¹, David J. Amor¹³, Corrado Romano^{20,21}, Joris Andrieux³³, and Santhosh Girirajan^{1,2,34}

1. Department of Biochemistry and Molecular Biology, Pennsylvania State University, University Park, PA 16802, USA.
2. Bioinformatics and Genomics Graduate program, Pennsylvania State University, University Park, PA 16802, USA.
3. Autism & Developmental Medicine Institute, Geisinger, Lewisburg, PA 17837, USA.
4. Department of Statistics, Pennsylvania State University, University Park, PA 16802, USA.
5. Institut Cochin, Inserm U1016, CNRS UMR8104, Université Paris Cité, CARPEM, Paris, France.
6. Department of Clinical Genetics, Erasmus MC, Univ. Medical Center Rotterdam, Rotterdam, The Netherlands.
7. Medical Genetics Unit, Hospital Santa Maria della Misericordia, Perugia, Italy.
8. Centre Hospitalier Universitaire de Liège. Domaine Universitaire du Sart Tilman, Liège, Belgium.
9. Division of Evolution and Genomic Sciences, School of Biological Sciences, Faculty of Biology, Medicine and Health, University of Manchester, Manchester, UK.
10. Manchester Centre for Genomic Medicine, St. Mary's Hospital, Central Manchester University Hospitals, NHS Foundation Trust Manchester Academic Health Sciences Centre, Manchester, UK.
11. Department of Medical Genetics, University and University Hospital Antwerp, Antwerp, Belgium.
12. Greenwood Genetic Center, Greenwood, SC 29646, USA.

- 37 13. Department of Clinical Genetics, Amsterdam UMC, Amsterdam, The Netherlands.
38 14. Center for Rare Diseases and Reference Developmental Anomalies and Malformation Syndromes,
39 CHU Dijon, Dijon, France.
40 15. Laboratoire de Genetique Chromosomique et Moleculaire, CHU Dijon, France.
41 16. Bruce Lefroy Centre, Murdoch Children's Research Institute, Melbourne, Australia.
42 17. Department of Paediatrics, University of Melbourne, Melbourne, Australia.
43 18. Department of Pediatrics, McGovern Medical School, University of Texas Health Science Center,
44 Houston, TX 77030, USA.
45 19. Department of Pathology and Molecular Medicine, McMaster University, Hamilton, Ontario, Canada.
46 20. Research Unit of Rare Diseases and Neurodevelopmental Disorders, Oasi Research Institute-IRCCS,
47 Troina, Italy.
48 21. Section of Clinical Biochemistry and Medical Genetics, Department of Biomedical and
49 Biotechnological Sciences, University of Catania School of Medicine, Catania, Italy.
50 22. Medical Genetics, ASP Ragusa, Ragusa, Italy.
51 23. Laboratory of Clinical Molecular Genetics and Cytogenetics, IRCCS San Raffaele Scientific Institute,
52 Milan, Italy.
53 24. Département de Génétique, Hôpital Pitié-Salpêtrière, Assistance Publique-Hôpitaux de Paris,
54 Sorbonne Université, 75019 Paris, France.
55 25. CHU Nantes, Medical Genetics Department, Nantes, France.
56 26. Service de Cytogenetique, CHU de Le Mans, Le Mans, France.
57 27. Department of Genetics, Bretonneau University Hospital, Tours, France.
58 28. CHU Toulouse, Department of Medical Genetics, Toulouse, France.
59 29. Toulouse Neuro Imaging, Center, Inserm, UPS, Université de Toulouse, Toulouse, France.
60 30. Institute of Genomics, University of Tartu, Estonia.
61 31. Health2030 Genome Center, Fondation Campus Biotech, Geneva, Switzerland.
62 32. Center for Integrative Genomics, Faculty of Biology and Medicine, University of Lausanne,
63 Switzerland.
64 33. Institut de Genetique Medicale, Hopital Jeanne de Flandre, CHRU de Lille, Lille, France.
65 34. Department of Anthropology, Pennsylvania State University, University Park, PA 16802, USA.
66

67

68 **Correspondence:**

69 Santhosh Girirajan

70 205A Life Sciences Building

71 Pennsylvania State University

72 University Park, PA 16803

73 Email: sxg47@psu.edu

74 **SUMMARY**

75 Variable expressivity of disease-associated variants implies a role for secondary variants that
76 modify clinical features. We assessed the effects of modifier variants towards clinical outcomes
77 of 2,252 individuals with primary variants. Among 132 families with the 16p12.1 deletion,
78 distinct rare and common variant classes conferred risk for specific developmental features,
79 including short tandem repeats for neurological defects and SNVs for microcephaly, while
80 additional disease-associated variants conferred multiple genetic diagnoses. Within disease and
81 population cohorts of 773 individuals with the 16p12.1 deletion, we found opposing effects of
82 secondary variants towards clinical features across ascertainment. Additional analysis of 1,479
83 probands with other primary variants, such as 16p11.2 deletion and *CHD8* variants, and 1,084
84 without primary variants, showed that phenotypic associations differed by primary variant
85 context and were influenced by synergistic interactions between primary and secondary variants.
86 Our study provides a paradigm to dissect the genomic architecture of complex disorders towards
87 personalized treatment.

88

89 INTRODUCTION

90 As large-scale sequencing studies uncover increasingly complex links between genomic variants
91 and heritable disorders, identifying the genetic etiology in an affected individual has become
92 more challenging¹. In contrast to Mendelian disorders caused by single genes, many disorders
93 are characterized by extensive phenotypic heterogeneity, where individuals with the same variant
94 exhibit a range of phenotypes with variable penetrance and expressivity^{2,3}. Some instances of
95 phenotypic heterogeneity can be explained by multiple genetic diagnoses, where more than one
96 pathogenic variant contributes to largely independent disorders in the same individual⁴. These
97 variants can even synergistically contribute to new phenotypes not associated with the individual
98 variants, such as seizures observed among individuals with variants in both *MKSI* (associated
99 with Meckel-Gruber syndrome) and *BBS1* (associated with Bardet-Biedl syndrome)⁵. Other cases
100 of phenotypic heterogeneity could occur when the clinical features of causal variants are
101 modified by secondary variants that do not cause disease themselves⁶. For example, rare variants
102 in histone modifier genes were enriched among individuals with the 22q11.2 deletion who
103 exhibited variably expressive congenital heart defects⁷. The complexity increases exponentially
104 for neurodevelopmental disorders, where the combined effects of primary and secondary variants
105 with differing frequency and effect sizes explain their broad heterogeneity^{8,9}. For example, recent
106 studies have found significant contributions of polygenic risk from common variants towards
107 phenotypes of individuals with pathogenic copy-number variants (CNVs)^{10,11}, such as
108 schizophrenia risk in individuals with 22q11.2 deletion¹²⁻¹⁷. Further, variable expressivity could
109 also be explained by cohort ascertainment, as many pathogenic variants are enriched among
110 individuals across disease ascertainment and lead to medical consequences in the general
111 population or healthy-biased cohorts¹⁸⁻²². For example, the autism-associated 16p11.2 deletion²³
112 is also associated with obesity, musculoskeletal, pulmonary, hematologic, and renal features in
113 the general population^{24,25}. This complex interplay necessitates a systematic assessment to fully
114 understand which modifier variants contribute to specific phenotypes when ascertained for the
115 same primary variant.

116 Rare recurrent CNVs represent excellent models to study variable expressivity, as the
117 large number of duplicated or deleted genes increases the likelihood of interactions with
118 modifiers elsewhere in the genome^{3,26}. For example, the rare heterozygous 520-kbp 16p12.1
119 deletion (hg18/NCBI36 when originally reported; currently maps to 16p12.2 based on

120 hg19/GRCh37) is enriched among children with severe neurodevelopmental features and is
121 inherited in >90% of cases from a parent who manifests milder psychiatric and cognitive
122 features^{27–29}. The phenotypic manifestation among individuals with this deletion differs based on
123 cohort ascertainment. For instance, the deletion was originally described in children with
124 developmental delay²⁷, but studies from the general population also identified associations with
125 multiple psychiatric and cognitive features^{22,30–32}. Thus, the 16p12.1 deletion is an ideal
126 paradigm for studying the effects of modifier variants on the clinical trajectory of a primary
127 variant. We previously found that severely affected children with the deletion have a global
128 increase in rare variant burden compared to parents with the deletion, and these trends are
129 consistent for other primary variants^{27,28,33}. Our findings suggested a “multi-hit” model for
130 complex disease etiology, where a primary variant sensitizes an individual for disease and the
131 clinical outcome is determined by other “hits” elsewhere in the genome³. However, it is not
132 completely understood how specific variant classes of differing effect size and frequency modify
133 clinical features across different ascertainment and primary variant contexts.

134 Here, we performed deep clinical and quantitative phenotyping and comprehensive
135 analysis of genomic data for 2,252 individuals with primary variants from diverse disease and
136 population-based cohorts (**Fig. 1**). We systematically dissected the roles of multiple secondary
137 variant classes towards developmental features in 132 families with the 16p12.1 deletion and
138 expanded our analysis to uncover phenotypic associations in 773 16p12.1 deletion individuals
139 from disease cohorts and healthy populations as well as 1,479 autism probands who carry a range
140 of other primary variants and 1,084 autism probands without primary variants. Our results show
141 that variant-phenotype associations are dependent on both the primary and secondary variant
142 context as well as cohort ascertainment (**Fig. 1**), allowing for more accurate dissection of the
143 genetic etiology of variably expressive traits associated with pathogenic variants.

144

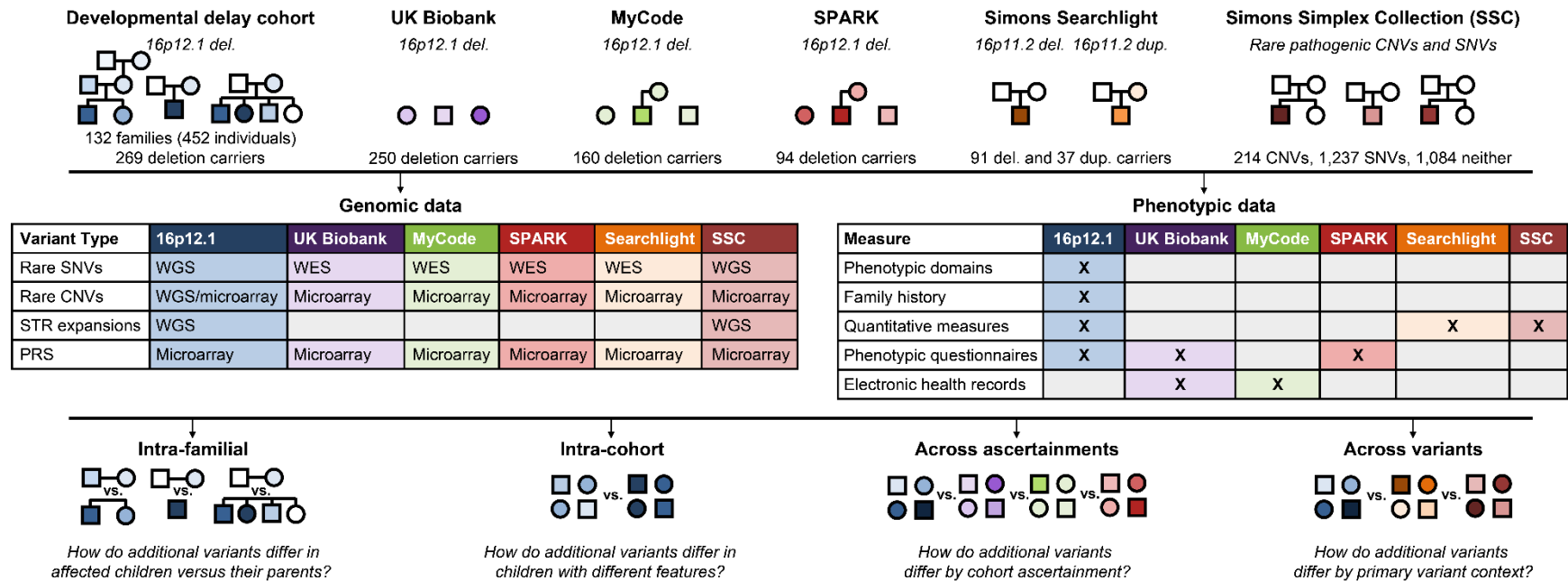


Figure 1. Overview of variant-phenotype analyses in 2,252 individuals with primary pathogenic variants. We assessed associations between variant classes and clinical phenotypes in six cohorts of individuals and families with primary variants. We directly recruited and assessed 132 families with the 16p12.1 deletion primarily ascertained for children with developmental delay (DD) (also including ten individuals from eight families from the Estonian Biobank not ascertained for DD). We further assessed 16p12.1 deletion carriers from cohorts with different ascertainment, including healthy volunteer-biased (UK Biobank), clinically-derived (MyCode), and single-disorder (SPARK, for autism) ascertainment. We finally assessed probands ascertained for autism with various primary pathogenic variants, including the 16p11.2 deletion or duplication (Simons Searchlight) and other large CNVs or rare SNVs in neurodevelopmental genes (SSC). We note that 100 probands in SSC have both pathogenic SNVs and CNVs and are included in both categories. Within and across these cohorts, we identified associations between up to 17 classes of rare and common variants (identified from WGS, WES, and microarrays) with phenotypic features from deep clinical datasets and electronic health records.

146 RESULTS

147 Variability of clinical features in 16p12.1 deletion

148 We recruited a cohort of 442 individuals from 124 families with the 16p12.1 deletion
149 (“DD cohort”), including multi-generational families, primarily ascertained for having children
150 with developmental delay (DD) (**Fig. 1**). We analyzed phenotypes from medical records, family
151 interviews, and online assessments for quantitative traits, such as non-verbal IQ³⁴ and social
152 responsiveness scores for autism-related social traits (SRS³⁵) (**Table S1A**). In total, 93% of
153 probands (84/90) inherited the deletion from a parent, with a slight bias towards maternal
154 inheritance (48/84, 57%), and 70% (87/124) of probands were male. Probands with the deletion
155 exhibited clinical features grouped across six broadly defined phenotypic domains, including
156 intellectual disability/developmental delay (ID/DD), behavioral, psychiatric, nervous system,
157 congenital, and growth/skeletal abnormalities (**Fig. 2A, Table S1A-B**). Probands also showed a
158 higher number of childhood developmental and behavioral features (i.e., increased complexity
159 scores, see *Methods*) compared to their siblings and cousins, while carrier siblings and cousins
160 manifested more features than non-carriers (**Fig. 2A**). Parents who transmitted the deletion
161 (“carrier parents”) often manifested milder cognitive or psychiatric phenotypes (**Fig. 2B**).
162 Probands had a 1.98 SD decrease in non-verbal IQ ($p=2.13\times 10^{-5}$) and a 1.91 SD increase in SRS
163 ($p=2.59\times 10^{-7}$) compared with their carrier parents (**Fig. 2C, Table S2A**). The average IQ score
164 among 16p12.1 deletion probands was 1.06 SD lower than all probands ascertained for autism
165 from the Simons Simplex Collection³⁴ (SSC) ($p=0.004$). The average SRS of 16p12.1 deletion
166 probands was 0.96 SD higher than probands with 16p11.2 deletions or duplications from the
167 Simons Searchlight cohort ($p=8.31\times 10^{-6}$) and 0.38 SD higher than SSC probands ($p=6.39\times 10^{-3}$)
168 (**Fig. 2C, Table S2A**). Beyond psychiatric traits, 16p12.1 deletion probands also showed
169 decreased head size ($p=0.001$) and increased body mass index (BMI, $p=0.009$) (**Fig. 2C, Table**
170 **S2A**). Finally, consistent with their ascertainment, probands exhibited significant delays in
171 several developmental milestones³⁶ compared to their siblings ($p<0.05$) (**Fig. 2D, Table S2B**).
172 Thus, our cohort represents families ascertained for probands who exhibit a range of
173 developmental features, including more severe IQ and social responsiveness deficits than
174 probands ascertained for autism or the 16p11.2 deletion.

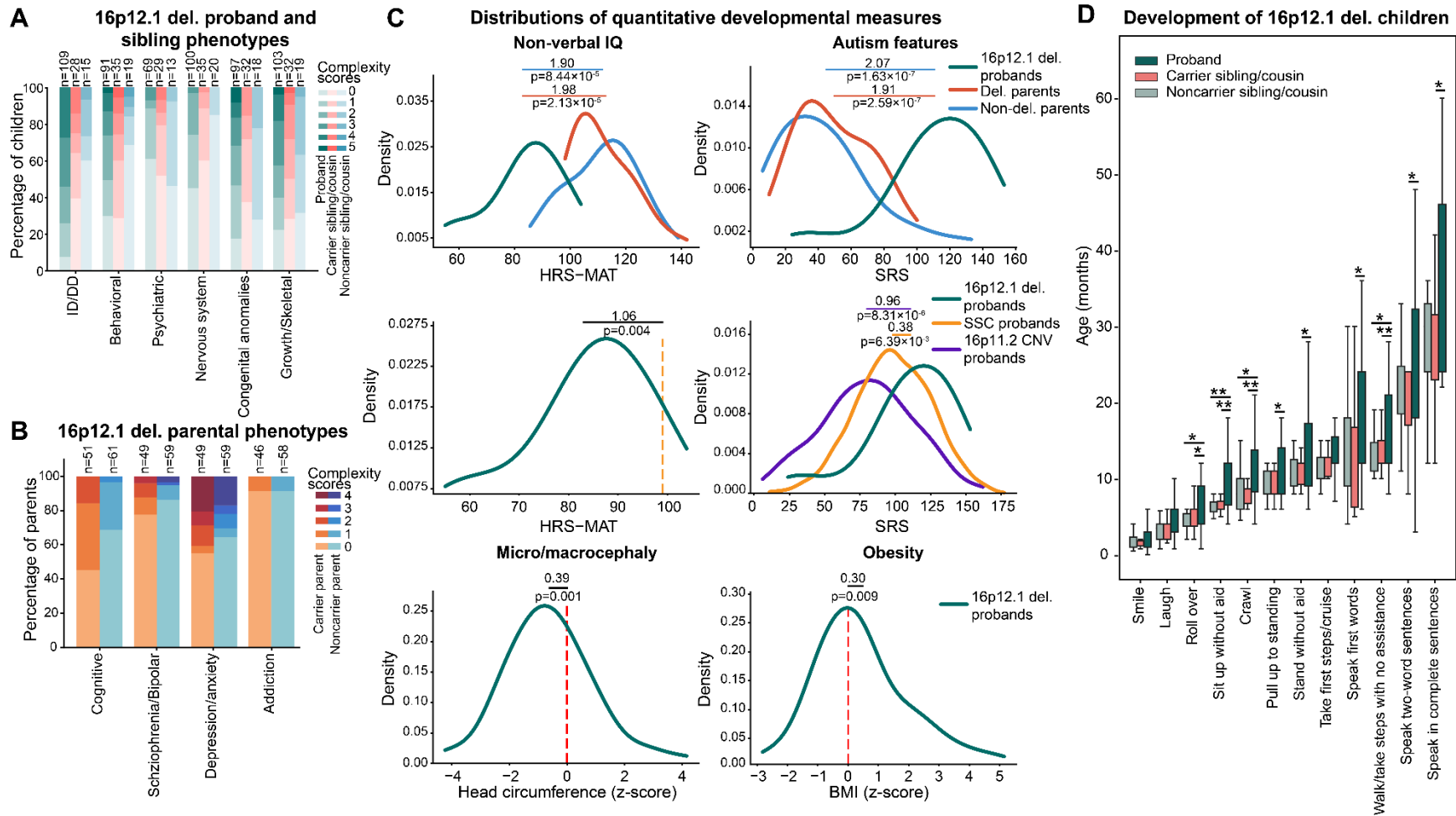


Figure 2. Variably expressive phenotypes of family members with the 16p12.1 deletion. (A) Distribution of complexity scores for six phenotypic domains in probands (n=69-109), carrier siblings and cousins (n=28-35), and noncarrier siblings and cousins (n=13-20) in 16p12.1 deletion families (numbers vary due to clinical data availability). Complexity scores were determined by identifying the number of clinical features manifested within each phenotypic domain (see Methods). (B) Distribution of complexity scores for four phenotypic domains in carrier parents (orange, n=46-51, orange) and non-carrier parents (blue, n=58-61) of 16p12.1 deletion probands. (C) Distributions of quantitative phenotypes observed in 16p12.1 deletion probands. Top plots

show the distribution of non-verbal IQ (HRS-MAT) and social responsiveness scores (SRS) in probands (green, n=10-27) compared to carrier (red, n=17-21) and non-carrier parents (blue, n=20-26). Middle plots compare the same scores in probands to the score for probands in the SSC cohort (SRS n=2,844, yellow; HRS-MAT mean derived from ³⁴) and probands with the 16p11.2 deletion or duplication from Simons Searchlight (n=139, purple). Bottom plots show the distribution of head circumference (n=64) and BMI z-scores (n=67) in deletion probands; red vertical lines represent the general population mean (i.e. z-score=0). P-values from Mann Whitney tests or one-sample t-tests. Individual scores for 16p12.1 deletion probands and parents are listed in **Table S1A. (D)** Distribution of the age of attainment for developmental milestones in probands (n=13-33), carrier siblings and cousins (n=16-18), and noncarrier siblings and cousins (n=11-15). One-tailed t-test, * $p \leq 0.05$, **Benjamini-Hochberg $FDR \leq 0.05$.

177 **Patterns of secondary variants within and across families**

178 Using WGS and microarray data, we evaluated 17 classes of secondary variants, including rare
179 coding SNVs (missense and splice variants with CADD Phred scores ≥ 25 and LOF variants),
180 non-coding SNVs (5' untranslated region [UTR], promoter [1kb upstream of transcription start
181 site], and enhancer [variants in regions with enhancer chromatin signatures in fetal brain]
182 variants), rare CNVs (deletions and duplications), and short tandem repeat expansions (STRs,
183 defined as repeat length $\geq 2SD$ than the cohort mean), a subset of which disrupted genes under
184 evolutionary constraint³⁷ (defined as LOEUF <0.35 and referred to as “(LF)” variants) (**Table**
185 **S1A**). We also calculated polygenic risk scores (PRS) for four psychiatric features, including
186 schizophrenia, intelligence, educational attainment, and autism^{38–41}. These secondary variants
187 could contribute to independent diagnoses from the 16p12.1 deletion, additively contribute to the
188 same phenotypes as the deletion, or synergistically modify the phenotypes of the deletion. We
189 first assessed whether probands carried additional pathogenic CNVs³³ or secondary variants that
190 were also present in ClinVar⁴² or in genes present in clinically relevant databases, such as Online
191 Mendelian Inheritance in Man⁴³ (OMIM), Developmental Brain Disorder database⁴⁴ (DBD), and
192 SFARI Gene⁴⁵. Overall, 58% of probands (57/99) had at least one such variant, including 19%
193 (19/99) of probands who had ClinVar-defined pathogenic variants (**Fig. S1A, Table S1D**). A
194 subset of these cases represented probands with multiple genetic diagnoses⁴. For instance, one
195 proband had a loss-of-function (LOF) variant in *KMT2A* and manifested Wiedemann-Steiner
196 syndrome features, including ID/DD, dysmorphic features, and hypertrichosis⁴⁶. Another
197 proband with an LOF variant in *DMD* showed expected hypotonia and muscular abnormalities as
198 well as ID/DD and craniofacial defects²⁸. Additionally, we found 17 probands with STR
199 expansions in spinocerebellar ataxia genes⁴⁷ such as *ATXN7* and *CACNA1A*⁴⁵. Although these
200 probands had fewer repeats than the clinical threshold for ataxia, 64.7% (11/17) of them
201 manifested nervous system phenotypes. We also identified a deleterious missense variant in
202 *POLR3E* on the non-deleted allele in a proband with global developmental delay and multiple
203 congenital defects (such as bilateral club feet and natal teeth).

204 We next sought to identify patterns of secondary variants in probands compared to their
205 parents (**Fig. S1B**). Probands carried more coding (LF) SNVs (union of missense, LOF, and
206 splice variants) ($p=0.041$) and missense (LF) variants ($p=0.017$), as well as increases in non-
207 coding SNVs in 5' UTRs ($p=0.045$), compared to their carrier parents (**Fig. 3A, Fig. S1C, Table**

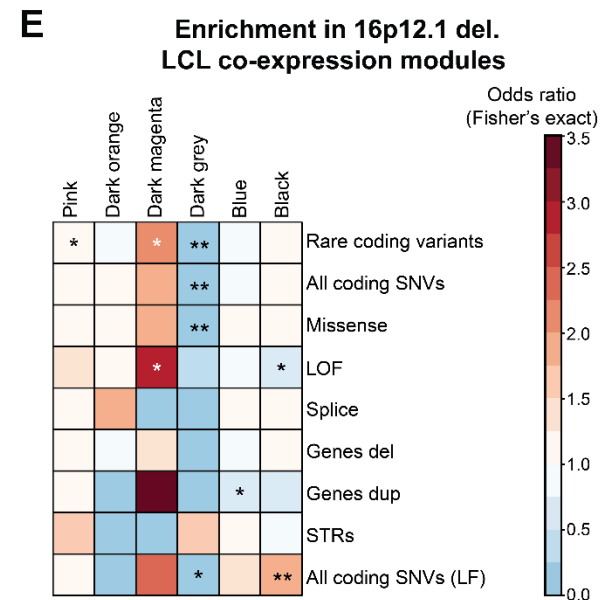
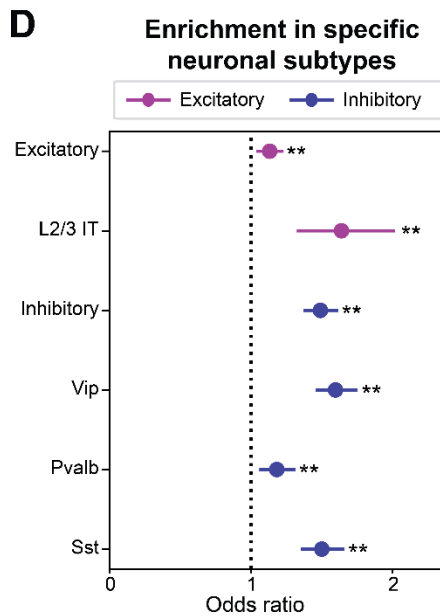
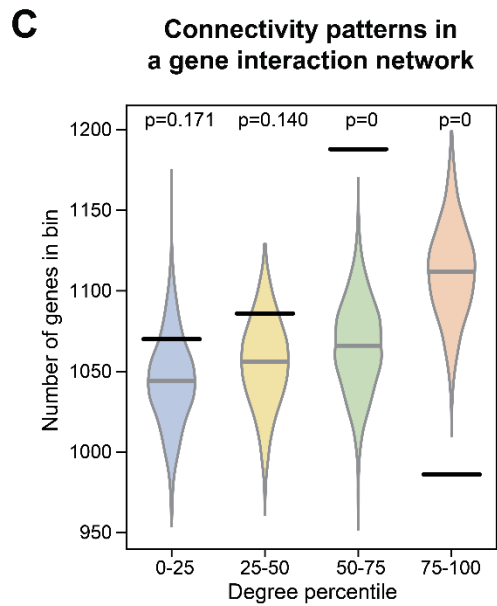
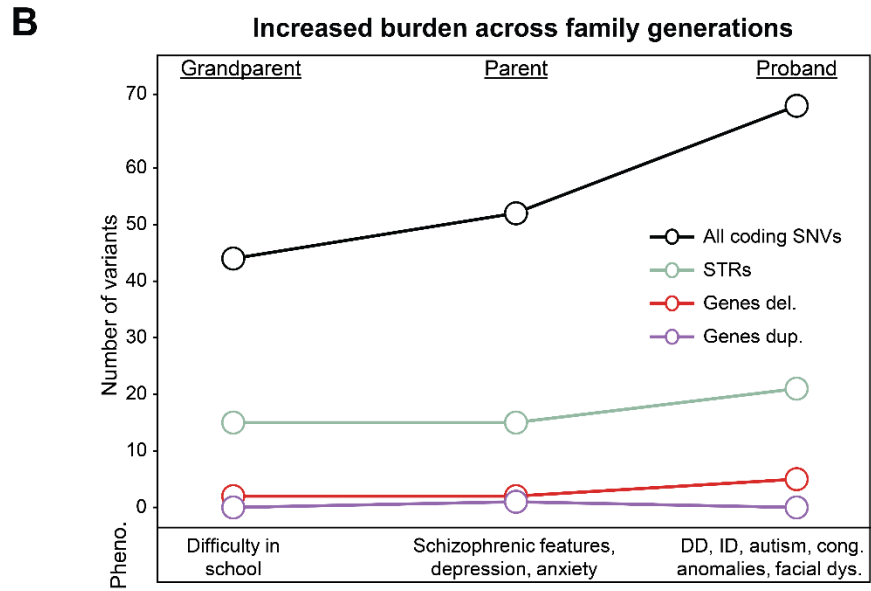
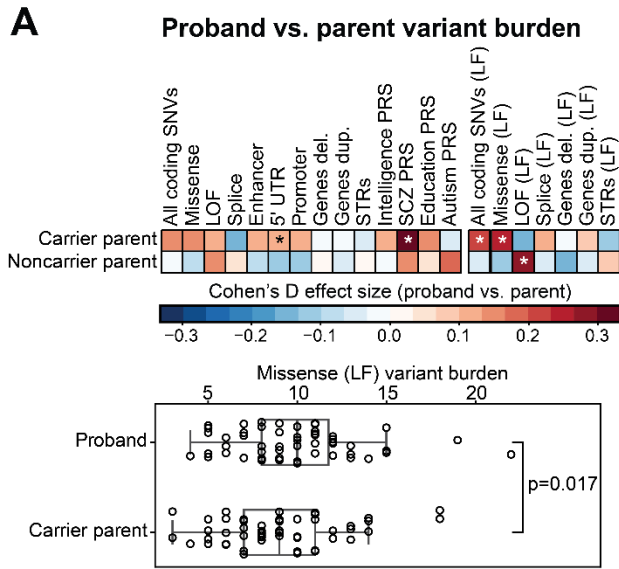


Figure 3. Secondary variants contribute to phenotypic variability within 16p12.1 deletion families. (A) Cohen's D effect sizes (top) for changes in secondary variant burden (i.e. rare variant burden or PRS) between probands and their carrier or noncarrier parents (n=49-54 pairs). * $p \leq 0.05$, paired one-tailed (rare variant classes) or two-tailed (PRS) t-test. Red indicates increased burden in probands relative to their parents. Boxplot (bottom) highlights increased burden of missense (LF) variants between probands and carrier parents. (B) Increased burden of rare variants corresponds with more severe clinical features across successive generations of 16p12.1 deletion carriers in a multi-generational family. (C) Distribution of genes by average connectivity (degree) within a brain-specific interaction network, binned into quartiles from 1000 simulations of randomly selected gene sets. Black lines represent the observed number of genes with secondary variants in 16p12.1 deletion probands in each degree quartile. Empirical p-values derived from simulation distributions. (D) Enrichment of genes with secondary SNVs in 16p12.1 deletion probands for genes preferentially expressed in neuronal classes (excitatory and inhibitory) and sub-classes (colored by main class) in the adult motor cortex. Fisher's exact test, **Benjamini-Hochberg $FDR \leq 0.05$. Full results are listed in **Table S2E**. (E) Enrichment of genes with secondary variants in probands for six gene co-expression modules identified from WGCNA analysis of lymphoblastoid cell lines (LCL) from individuals with the 16p12.1 deletion⁵⁴. Fisher's exact test, * $p \leq 0.05$, **Benjamini-Hochberg $FDR \leq 0.05$.

210 **S2C**). Proband also carried higher schizophrenia polygenic risk than their carrier parents
211 ($p=0.009$), showing that polygenic risk may also contribute to the features observed among
212 16p12.1 deletion probands (**Fig. 3A, Fig. S1C, Table S2C**). Except for an increase in LOF (LF)
213 variants ($p=0.039$), no differences across other variant classes were observed in probands
214 compared to non-carrier parents (**Fig. 3A, Table S2C**). This is consistent with a model in which
215 the deletion and secondary variants are transmitted in specific patterns that lead to different
216 clinical outcomes in probands. In fact, assessment of multi-generational families showed that
217 variant burden tends to compound over generations towards increased phenotypic severity. For
218 example, in one multi-generational family, the carrier grandparent had mild cognitive features,
219 while the carrier parent manifested psychiatric features and the proband had neurodevelopmental
220 features (**Fig. 3B**). This increase in phenotypic severity corresponded with an increase in the
221 burden of multiple rare variant classes across generations, akin to the genetic anticipation
222 observed for certain Mendelian disorders.

223 We further profiled the putative functions of secondary variants and found that missense
224 variants were enriched for brain-expressed, constrained, and post-synaptic density genes
225 ($FDR \leq 7.41 \times 10^{-9}$)^{37,48}, while genes with LOF variants were depleted for these gene sets
226 ($FDR \leq 0.007$) (**Fig. S2A, Table S2G**). This suggests that LOF variants in essential genes may not
227 be tolerated, particularly in the presence of the deletion, while less severe variants in these genes
228 may contribute to neurodevelopmental disorders seen in probands²⁶. As further evidence of this,
229 secondary variants were enriched (empirical $p=0.000$) for genes with intermediate connectivity
230 within a brain-specific gene interaction network^{49,50} but depleted for genes with high network
231 connectivity, which typically represent essential genes across species⁵¹ (empirical $p=0.000$; **Fig.**
232 **3C, Fig. S2B, Table S2D**). Secondary variant genes were also preferentially expressed in several
233 brain regions during early and mid-fetal development⁵², including the frontal cortex ($FDR \leq 0.05$)
234 and hippocampus ($FDR = 1.73 \times 10^{-5}$) (**Fig. S2C, Table S2H**). SNVs in particular were enriched
235 for genes preferentially expressed across multiple neuronal classes in the adult motor cortex⁵³,
236 including excitatory ($FDR = 0.013$) and inhibitory ($FDR = 3.82 \times 10^{-23}$) neurons (**Fig. 3D, Table**
237 **S2E**). All coding SNVs (LF) were also enriched for genes co-expressed with 16p12.1 deletion
238 genes (black module, $FDR = 0.016$) in lymphoblastoid cell lines (LCL) derived from a subset of
239 19 individuals with the deletion⁵⁴ (**Fig. 3E, Table S2F**). Overall, our results indicate that a

240 diverse range of biologically relevant modifiers contribute to variable phenotypes in probands
241 with 16p12.1 deletion.

242

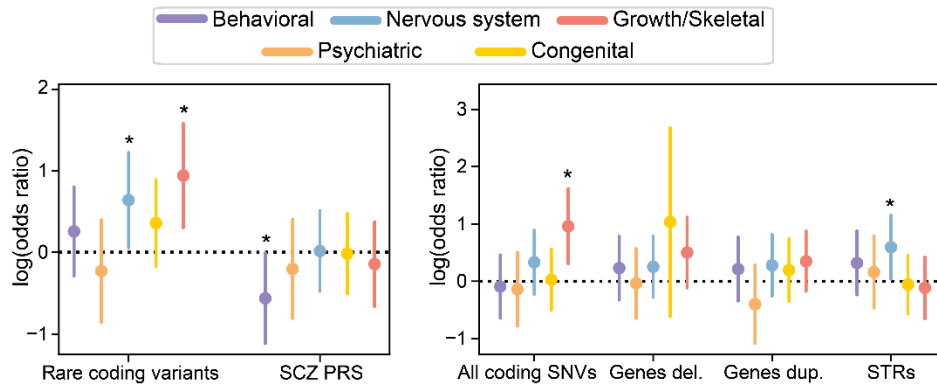
243 **Distinct secondary variant classes contribute to specific phenotypic outcomes**

244 We next used a series of logistic regression models to assess contributions of rare variant classes
245 and PRS towards individual phenotypic domains. Overall, rare coding variants contributed to
246 nervous system (logOR=0.640, p=0.032) and growth/skeletal features (logOR=0.941, p=0.004)
247 (**Fig. 4A, Table S3A**). Specifically, STRs were associated with nervous system features
248 (logOR=0.596, p=0.036) while SNVs were associated with growth/skeletal features
249 (logOR=0.962, p=0.004) (**Fig. 4A, Fig. S3A, Table S3A**). In contrast, schizophrenia PRS was
250 negatively associated with behavioral phenotypes (logOR=-0.563, p=0.046) (**Fig. 4A, Table**
251 **S3A**). Combined variant models explained an average of 8% variance (McFadden's pseudo-R²;
252 range 2% to 14%) for each phenotypic domain, and showed better performance than models built
253 using individual variant classes (average of 2% explained variance) (**Fig. S3B, Table S3A**).
254 These estimates further highlighted the specificity of variant-phenotype associations; for
255 example, STRs (LF) explained 12% of variance in nervous system defects but less than 4% of
256 variance for other features (**Fig. S3B, Table S3A**). Orthogonal burden tests also identified fewer
257 rare variants in enhancers, promoters, and 5' UTR elements (p≤0.012) as well as increased
258 autism PRS (p=0.028) among probands with psychiatric features (**Fig. S3C, Table S3D**). These
259 results suggest that the modifying roles of different secondary variant classes vary across specific
260 phenotypes, with PRS in particular modulating behavioral features.

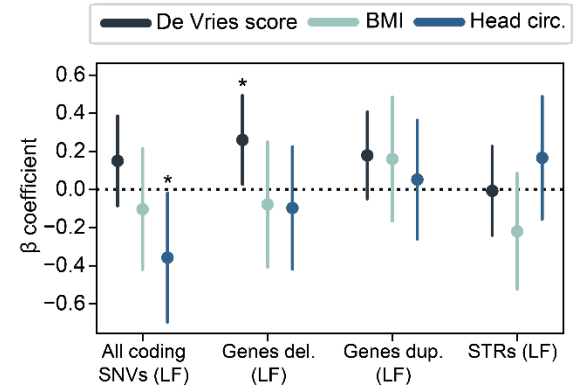
261 Linear regression models testing specific variant classes towards quantitative traits
262 revealed negative associations of head circumference z-scores with SNVs (LF) (β =-0.357,
263 p=0.039) (**Fig. 4B, Fig. S3A, Table S3A**). Secondary CNVs were associated with increased de
264 Vries scores, a quantitative assessment for global developmental features⁵⁵ (deletions: β =0.288,
265 p=0.013; duplications: β =0.246, p=0.030) (**Fig. 4B, Fig. S3A, Table S3A**). Correlation analyses
266 revealed that intelligence and educational attainment PRS were positively correlated with head
267 circumference (education r=0.318, p=0.026; intelligence r=0.287, p=0.045), while duplications
268 (LF) were negatively correlated with social responsiveness deficits (r=-0.605, FDR=0.030) (**Fig.**
269 **S3D, Table S3B**).

270

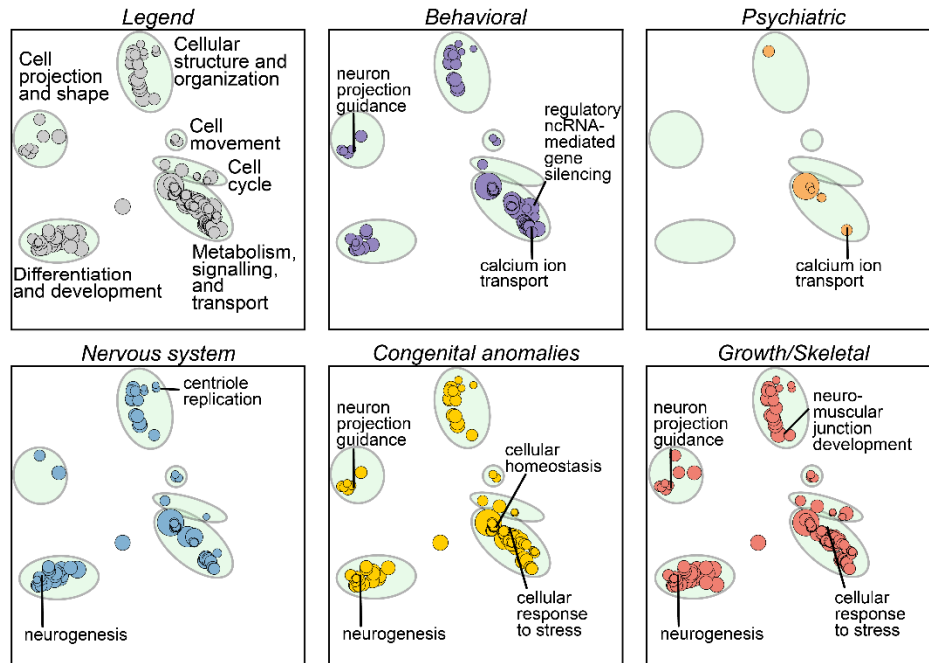
A Contributions of specific variant classes toward proband phenotypes



B Contributions of variant classes towards quantitative phenotypes



C Biological functions of secondary variants in probands by phenotype



D Disease gene burden by phenotype

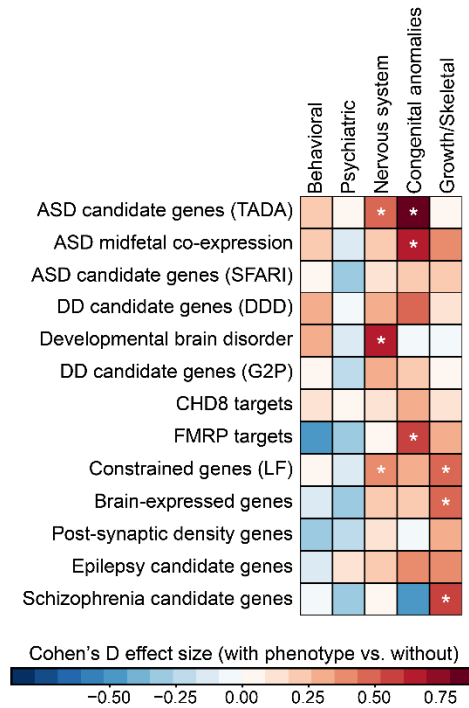


Figure 4. Secondary variant associations for phenotype domains of the 16p12.1 deletion. (A) Forest plots show log-scaled odds ratios from logistic regression models for secondary variant burden in 16p12.1 deletion probands with higher complexity scores for five phenotypic domains, compared with probands with lower complexity scores (n=47-71). * $p \leq 0.05$. Model results for variants (LF) are shown in **Fig. S3A**. (B) Forest plots show β coefficients from linear regression models for secondary variant burden in genes under evolutionary constrain (LF genes) towards quantitative phenotypes in deletion probands (n=43-76). * $p \leq 0.05$. Model results for variants without LF filter are shown in **Fig. S3A**. (C) Gene Ontology (GO) biological process terms enriched among secondary variants in probands with each phenotypic domain. Circles represent individual GO terms, clustered based on semantic similarity into broad categories (green ovals, as defined in the “legend” plot). Size of each circle represents the number of genes in each term, such that broader terms are larger. Colored circles in each plot indicate significant enrichment of the GO terms for the given phenotype. (D) Changes in burden of secondary variants disrupting sets of genes involved with neurodevelopmental disease and related functions (see Methods) in probands with phenotypic domains (n=23-67) compared to probands without each domain (n=12-36). * $p \leq 0.05$, one-tailed t-test.

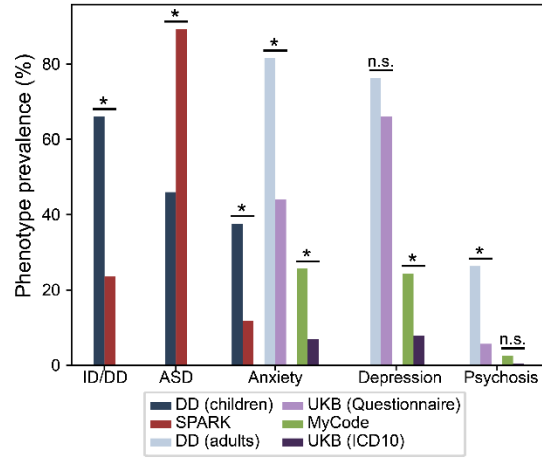
273 Secondary variants in probands ascertained for the same phenotypes showed specific
274 enrichments for biological function, including neuromuscular junction development genes in
275 probands with growth/skeletal defects and axonogenesis-related genes in probands with
276 behavioral and nervous system features (**Fig. 4C, Table S3B**). Additionally, probands with
277 specific phenotypes showed increased burden of rare variants in key neuronal genes, such as
278 *FMRP*-binding targets⁴⁸ ($p=0.050$) in probands with congenital anomalies and candidate autism⁴⁵
279 ($p=0.014$) and developmental brain disorder genes⁴⁴ ($p=0.001$) in probands with nervous system
280 defects (**Fig. 4D, Table S3C**). Overall, our findings indicate that the disruption of distinct
281 biological functions and molecular pathways by secondary variants may underlie specific
282 phenotypic features of individuals with the deletion.

283

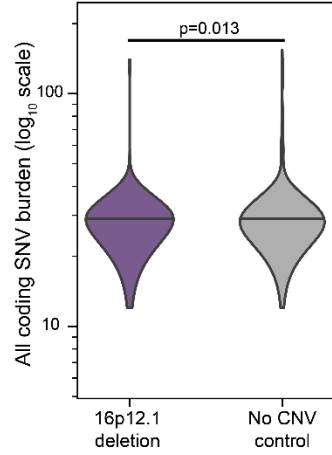
284 **Differing ascertainment confer distinct genotype-phenotype patterns**

285 Clinical outcomes associated with the same genetic variant may vary across cohorts with
286 different ascertainment, especially for cohorts composed of affected individuals compared to
287 those drawn from the general population^{56,57}. We sought to compare the phenotypic effects of the
288 16p12.1 deletion in 757 individuals with complete phenotypic data, including 253 pediatric and
289 adult carriers from the DD cohort and three independent cohorts with distinct ascertainment:
290 SPARK ($n=94$), where families were ascertained for probands with autism features⁵⁸, and two
291 population-based cohorts^{59,60}, the healthy-biased UK Biobank (UKB; $n=250$) and the hospital-
292 derived Geisinger MyCode (MyCode; $n=160$) (**Fig. 1**). Assessment of UKB individuals with the
293 deletion showed enrichment for a variety of clinical phenotypes within electronic health record
294 (EHR) data, including circulatory, endocrine, and genitourinary features ($n=3,488$, $FDR \leq 0.004$)
295 (**Fig. S4A, Table S5H**). PheWAS analysis further revealed enrichment of obesity- and kidney-
296 related features, including type 2 diabetes and hypertension ($n=99,363-255,262$, $p \leq 3.48 \times 10^{-7}$)
297 (**Fig. S4B, Table S5I**). This pattern of obesity-related features is in line with the pattern of
298 increased BMI we observed in probands in the DD cohort (**Fig. 2C**). To more directly compare
299 phenotype prevalence across cohorts, we next harmonized EHR and clinical questionnaire
300 responses (**Table S4A**). As expected, the prevalence of neuropsychiatric phenotypes in both
301 pediatric and adult deletion carriers varied across the cohorts (**Fig. 5A, Fig. S4C-D, Table S4B**).
302 For example, we found increased anxiety symptoms in adults from the DD cohort compared to

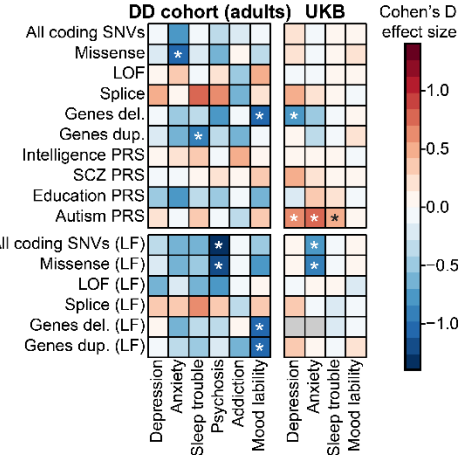
A Phenotype prevalence by ascertainment



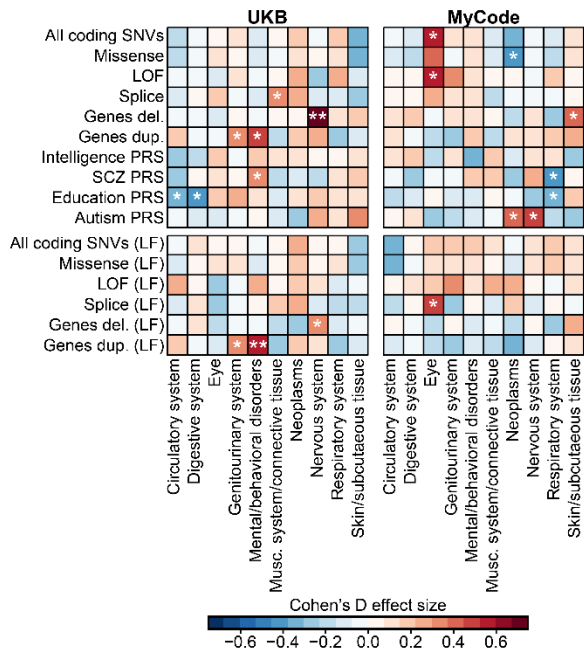
B Decreased secondary variant burden in UKB



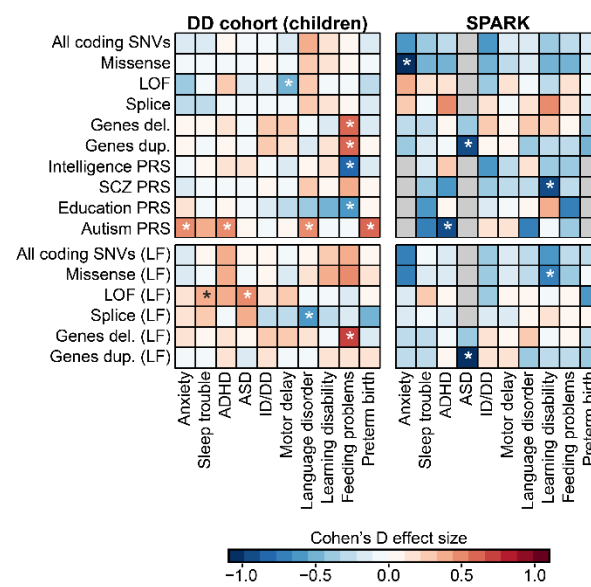
C Variant contributions to psychiatric features in adults



D Variant contributions to clinical features in adults



E Variant contributions to features in children



F Joint variant contributions to child phenotypes

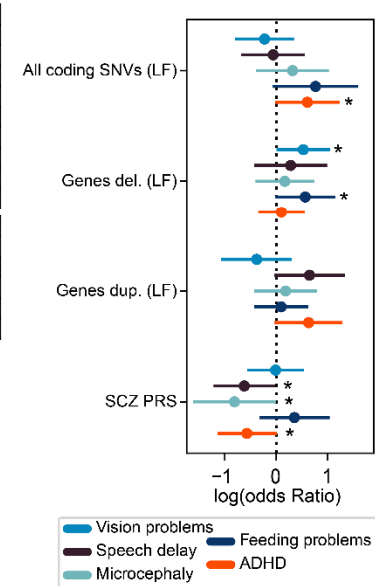


Figure 5. Effects of ascertainment on associations of 16p12.1 deletion. (A) Prevalence of phenotypes among adults and children with 16p12.1 deletion from four ascertainments: DD cohort (adults n=38, children n=93-151), SPARK (n=51-56), UK Biobank (UKB; questionnaire n=50-53, ICD10 n=217), and MyCode (n=160). Fisher's exact test, *p≤0.05. (B) Distribution of rare secondary SNVs in UKB individuals with 16p12.1 deletion (n=240, left) and age and sex-matched controls without large rare (>500kb) CNVs (n=2,640, right). P-value from two-tailed t-test. (C) Associations of secondary variant burden with select psychiatric phenotypes derived from clinical questionnaires in 16p12.1 deletion adults from DD cohort (n=24-31) and UKB (n=46-249). Two-tailed t-test, *p<0.05. (D) Associations of secondary variant burden with select clinical phenotypes derived from EHR data (ICD10 codes) in 16p12.1 deletion individuals from UKB (n=187-218) and MyCode (n=143-159). Two-tailed t-test, *p≤0.05. **Benjamini-Hochberg FDR≤0.05. (E) Associations of secondary variant burden and select developmental phenotypes in children with 16p12.1 deletion from the DD cohort (n=67-125) and SPARK (n=27-56). Two-tailed t-test, *p≤0.05. (F) Associations of secondary variant burden and developmental phenotypes from joint logistic models of 16p12.1 deletion children from the DD and SPARK cohorts (n=98-125). Joint models for non-LF are shown in **Fig. S4J**, and joint models for adults are shown in **Fig. S4H-I**. *p≤0.05.

305 UKB ($p=4.18\times 10^{-4}$) (**Fig. 5A, Table S4B**), likely reflecting biases due to ascertainment for
306 severely affected family members in the DD cohort compared to healthy volunteers in UKB⁶¹.

307 We next investigated how patterns of secondary variants differed across cohorts. Adult
308 deletion carriers in UKB showed decreased burden of additional rare coding SNVs compared to
309 individuals without large CNVs ($p=0.013$) (**Fig. 5B, Table S4C**). Reduced secondary variant
310 burden in 16p12.1 deletion carriers may explain the less severe features observed in the UKB
311 compared to those typically observed among deletion carriers in the clinic. This trend was
312 reversed in SPARK, where individuals with the deletion had an increased burden of SNVs (LF)
313 compared to individuals without large CNVs ($p=0.048$) (**Fig. S4E, Table S4C**). Thus, we
314 observed a higher rare variant burden in deletion carriers compared to controls in cohorts with
315 more severe disease ascertainment (SPARK) and reduced burden in cohorts with less severe
316 ascertainment (UKB). We also directly compared the variant burden between deletion carriers in
317 the DD cohort to identically processed data from eight deletion carriers in the Estonian
318 Biobank⁶². As expected, Estonian Biobank carriers showed a depletion of missense ($FDR\leq 0.012$)
319 and non-coding SNVs ($FDR\leq 0.019$) compared with probands and carrier parents in the DD
320 cohort (**Fig. S4F, Table S4J**).

321 We next assessed how the relationship of secondary variants and phenotypes varies by
322 ascertainment by assessing the burden of variant classes in individuals with and without specific
323 phenotypes across cohorts. We found similar trends for psychiatric features based on self-
324 reported questionnaires for adults in both the DD cohort and UKB (**Fig. 5C, Table S4D**). For
325 example, autism PRS was associated with depression, anxiety, and sleep disturbance ($p\leq 0.037$)
326 in UKB (**Fig. 5C, Table S4D**). Conversely, rare variants were negatively associated with
327 psychiatric features in both UKB and DD cohorts, including SNVs (LF) with anxiety in UKB
328 ($p=0.010$) and psychosis in DD ($p=0.003$), and deletions with depression in UKB ($p=0.016$) and
329 mood lability in DD ($p=0.027$). These data suggest potential opposing roles of PRS and rare
330 variants towards specific psychiatric features. We further compared secondary variant profiles
331 for broader groups of clinical features represented by ICD10 chapters in UKB and MyCode, in
332 contrast to assessment of specific psychiatric features from questionnaires. In UKB, nervous
333 system features were associated with deletions ($FDR=0.007$), while mental health features were
334 associated with duplications (LF) ($FDR=0.035$) and schizophrenia PRS ($p=0.045$) (**Fig. 5D,**
335 **Table S4E**). In MyCode, eye phenotypes were associated with SNVs ($p=0.016$), while autism

336 PRS was associated with both nervous system defects and cancer ($p \leq 0.039$) (**Fig. 5D, Table**
337 **S4E**). These differences in genotype-phenotype patterns reflect ascertainment differences
338 between the cohorts, potentially due to healthcare system differences, phenotyping modalities, or
339 biases stemming from healthy volunteers versus clinical patients.

340 We further observed differences in children with the deletion from the DD cohort and
341 those in SPARK. In general, both rare variants and PRS were associated with increased risk for
342 neurodevelopmental features in DD probands (for example, duplications and deletions for
343 feeding difficulty; $p \leq 0.028$) but decreased risk in SPARK probands (for example, decreased
344 missense variants in individuals with anxiety; $p = 0.025$) (**Fig. 5E, Table S4F**). In fact,
345 individuals with ADHD had increased autism PRS in the DD cohort ($p = 0.017$) but decreased
346 autism PRS risk in SPARK ($p = 0.035$) (**Fig. 5E, Table S4F**). This trend reflects the influence of
347 ascertainment towards variant-phenotype associations, where secondary variants may not show
348 the expected associations in highly ascertained cohorts due to saturated genetic risk for the
349 ascertained phenotype, such as ADHD and autism PRS in SPARK (**Fig. S4G**). Overall, we found
350 marked differences in secondary variant-phenotype associations between cohorts (**Fig. 6C-E**),
351 which potentially explains the variable phenotypic trajectories of the deletion observed across
352 ascertainment.

353 We finally combined individuals with the 16p12.1 deletion across cohorts and developed
354 logistic regression models to identify variant-phenotype associations independent of
355 ascertainment bias. We found 12 associations among children in the combined SPARK and DD
356 cohorts ($n = 84-125$), including SNVs with preterm birth ($\log OR = 1.24$, $p = 0.039$), deletions with
357 vision problems ($\log OR = 0.878$, $p = 0.018$), and SNVs (LF) with ADHD ($\log OR = 0.602$, $p = 0.049$)
358 (**Fig. 5F, Fig. S4J, Table S4G**). Across adults in the DD, UKB, and MyCode cohorts ($n = 331$),
359 duplications were associated with anxiety ($\log OR = 0.092$, $p = 0.004$) (**Fig. S4H, Table S4G**).
360 When examining EHR-derived features in UKB and MyCode ($n = 321$), duplications were
361 associated with circulatory system features ($\log OR = 0.278$, $p = 0.037$) and deletions were
362 associated with nervous system ($\log OR = 0.352$, $p = 0.016$) and skin/subcutaneous tissue
363 ($\log OR = 0.336$, $p = 0.020$) phenotypes (**Fig. S4I, Table S4G**). Thus, leveraging combined data
364 from multiple cohorts allowed for the increased statistical power necessary for deriving robust
365 variant-phenotype associations across ascertainment.

366

367 **Differing contributions of secondary variants by primary variant ascertainment**

368 To extend our findings beyond the 16p12.1 deletion, we assessed contributions of secondary
369 variants towards developmental, cognitive, and behavioral features among 1,479 probands with
370 different rare pathogenic CNVs or SNVs in known neurodevelopmental genes who were
371 ascertained for the same disorder, i.e. autism (**Fig. 1**). We first assessed 128 probands with
372 reciprocal 16p11.2 deletions (n=91) and duplications (n=37) in the Simon Searchlight cohort⁶³
373 and found eight variant-phenotype associations using linear regression models (**Fig. 7A, Table**
374 **S6A**). Among 16p11.2 deletion probands, schizophrenia PRS contributed to higher full-scale IQ
375 ($\beta=0.343$, $p=0.020$), while secondary deletions contributed to decreased IQ ($\beta=-0.283$, $p=0.040$),
376 similar to previous findings¹¹ (**Fig. 6A, Fig. S5A, Table S5A**). Different trends were observed in
377 16p11.2 duplication probands; deletions and duplications were negatively associated with autism
378 behavioral features (BSI; duplications: $\beta=-0.497$, $p=0.022$; deletions: $\beta=-0.432$, $p=0.037$) and
379 duplications (LF) were negatively associated with SRS ($\beta=-0.701$, $p=0.002$) (**Fig. 6A, Fig. S5A,**
380 **Table S5A**). Orthogonal correlation analyses revealed several other trends, including opposing
381 effects of secondary duplications towards BSI (16p11.2 deletion individuals: $r=0.241$, $p=0.023$;
382 16p11.2 duplication individuals: $r=-0.391$, $p=0.019$) (**Fig. S5B, Table S5E**).

383 We next assessed 214 probands with a more heterogeneous set of large CNVs, including
384 pathogenic deletions and duplications³³ at 15q13.3, 3q29, and 16p13.11, from SSC⁶⁴. Among
385 probands with large deletions, linear regression models uncovered negative associations between
386 secondary duplications (LF) with BMI ($\beta = -0.275$, $p=0.049$), while secondary deletions (LF)
387 were associated with decreased IQ in probands with large duplications ($\beta = -0.255$, $p=0.021$)
388 (**Fig. 6A, Fig. S6A, Table S5A**). Correlation analyses revealed additional associations, including
389 SNVs with coordination impairment in probands with large duplications (DCDQ, $r=0.178$,
390 $p=0.042$) and decreased SRS in those with large deletions ($r=-0.246$, $p=0.035$) (**Fig. S6B, Table**
391 **S5F**). We further assessed 1,237 SSC probands with SNVs disrupting canonical
392 neurodevelopmental genes⁴⁴, such as *CHD8*, *DYRK1A*, *SCN1A*, and *PTEN*. We again identified a
393 negative association for deletions (LF) with IQ ($\beta = -0.154$, $p=5.23 \times 10^{-5}$), while STRs were
394 associated with increased IQ ($\beta=0.093$, $p=0.015$) (**Fig. 6A, Fig. S6A, Table S5A**). Correlation
395 analyses uncovered additional effects, such as externalizing behavior ($r=-0.111$, $p=0.001$) and
396 repetitive behavior ($r=-0.083$, $p=0.017$) negatively correlating with educational attainment PRS
397 (**Fig. S6B, Table S5F**). Thus, we found some consistent patterns across primary variant

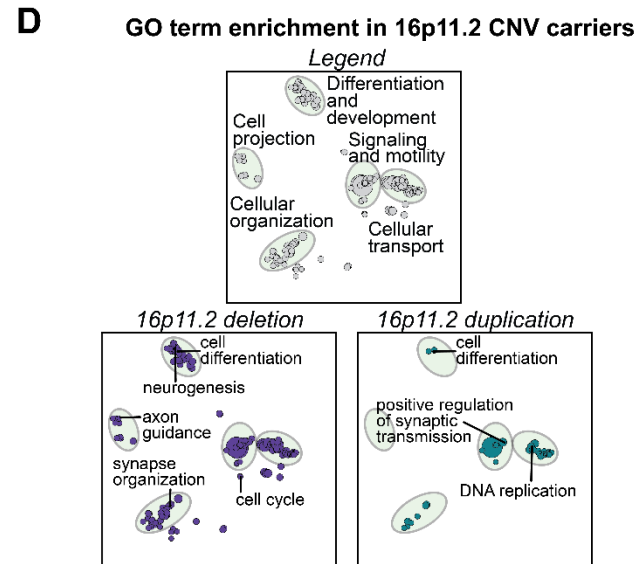
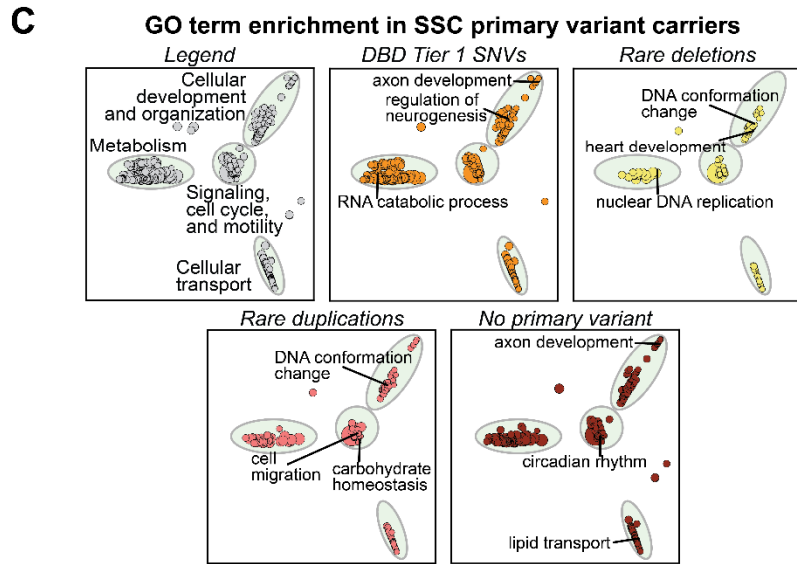
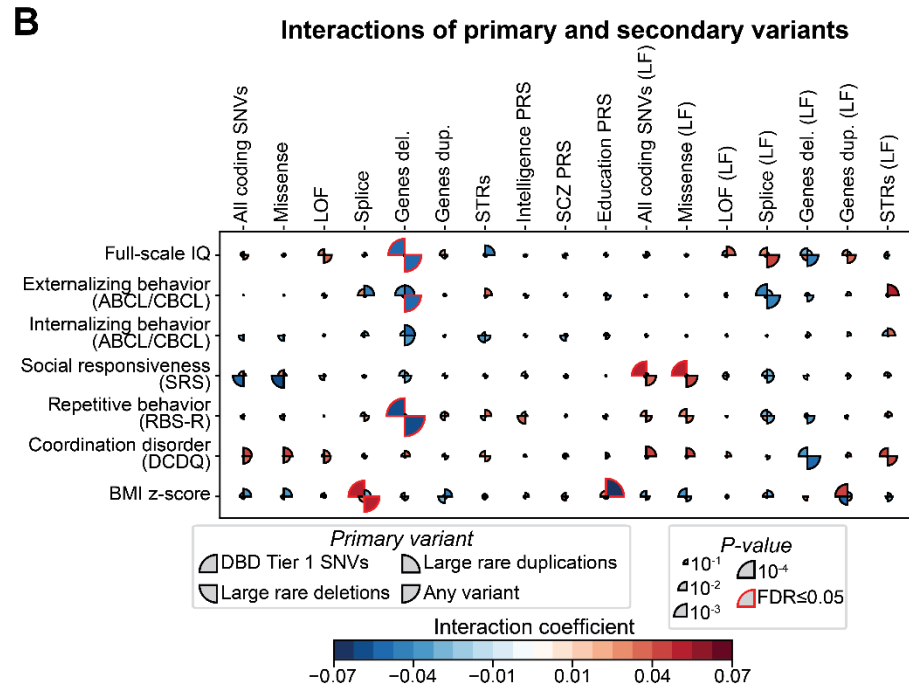
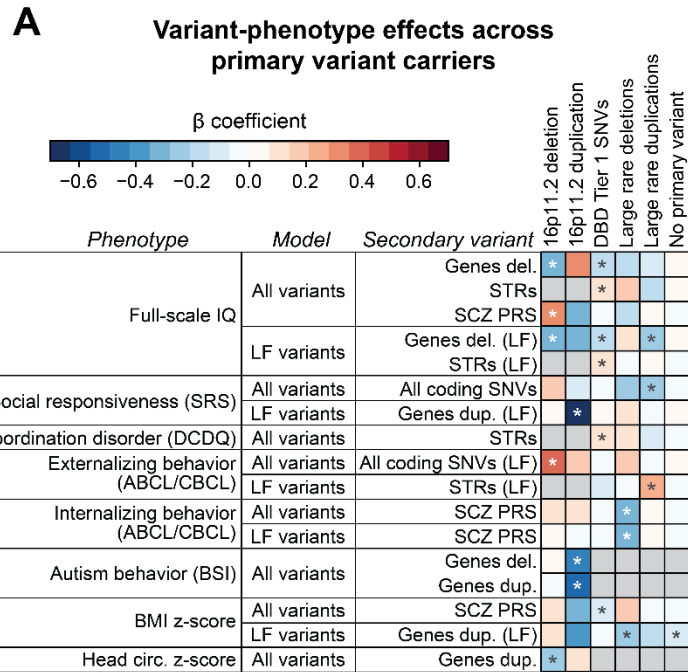


Figure 6. Secondary variant associations in probands with primary variants. (A) Heatmap shows β coefficients from select linear regression models for secondary variant burden (y-axis, third column) towards quantitative developmental phenotypes (y-axis, first column) in probands from SSC and Simons Searchlight cohorts with different classes of primary variants (x-axis) (n=21-660). * $p \leq 0.05$. (B) β coefficients from linear regression models examining interactions between primary variants (pie chart slices) and specific secondary variant classes (x-axis) towards quantitative phenotypes (y-axis) in SSC probands (n=1,597-2,591). Color of pie chart slices indicate interaction coefficients, and size of pie chart slices indicate p-value for strength of interaction coefficient. Red highlights indicate Benjamini-Hochberg FDR ≤ 0.05 (C-D) Gene Ontology (GO) biological process terms enriched among secondary variants observed in (C) probands with different classes of primary variants from the SSC cohort and (D) probands with 16p11.2 deletions and duplications from the Searchlight cohort. Circles represent individual GO terms, clustered based on semantic similarity into broad categories (green ovals, as defined in the two “legend” plots). Size of each circle represents the number of genes in each term, such that broader terms are larger. Colored circles in each plot indicate significant enrichment of the GO terms for the given primary variant.

400 ascertainments, such as negative effects of rare deletions on IQ, while secondary variant
401 effects on other features were more dependent on the primary variant context.

402 We additionally examined secondary variants in 1,084 SSC probands without primary
403 variants in the above categories to assess the role of the genetic background outside of a primary
404 variant context. The only observed association from regression analysis was for duplications
405 (LF) and lower BMI ($\beta=-0.087$, $p=0.031$) (**Fig. 6A, Table S5A**). The paucity of associations in
406 the absence of primary variants suggests that secondary variant classes mostly exert their effects
407 through interactions with primary variants instead of contributing directly towards disease
408 phenotypes. To assess this, we used linear models to identify interactions between primary and
409 secondary variants towards clinical features. We found ten instances of multiplicative
410 interactions, including primary SNVs and secondary SNVs (LF) towards SRS ($\beta=0.054$,
411 $FDR=0.034$) as well as primary SNVs and secondary deletions towards full-scale IQ ($\beta=-0.058$,
412 $FDR=0.020$) and repetitive behavior ($\beta=-0.062$, $FDR=0.011$) (**Fig. 6B, Table S5B**). Notably,
413 these interactions tended to be primary variant-specific, further supporting the hypothesis that
414 secondary variant effects are influenced by primary variant context.

415 The relevance of primary variant context was further evident when we assessed the
416 biological functions of secondary variants (**Fig. 6C-D**). For example, secondary variants in
417 probands with primary SNVs showed specific enrichments for neuronal development and cell-to-
418 cell signaling, while probands with primary deletions showed enrichments for DNA repair and
419 replication (**Table S5C**). Secondary variants in 16p11.2 duplication probands were enriched for
420 DNA replication and synaptic transmission regulation, while variants in probands with the
421 reciprocal deletion were enriched for cell cycle regulation and synapse organization (**Table**
422 **S5D**). In fact, multiple genes within the 16p11.2 deletion have similar molecular functions (i.e.
423 *MAPK3* and cell cycle regulation⁶⁵), and many of the same GO enrichments, including neuronal
424 differentiation and projection, were identified among differentially expressed genes in animal
425 models of genes within the 16p11.2 region⁶⁶ (**Table S5D**). We therefore posit that modifier
426 variants influence developmental features by acting additively or synergistically in molecular
427 pathways disrupted by the primary variant, further underscoring the importance of primary
428 variant context⁶⁷.

429 **DISCUSSION**

430 Our comprehensive analysis of 2,252 individuals with primary variants from several diverse
431 cohorts allowed us to find strong evidence that modifier variants confer distinct risks towards
432 different developmental and clinical phenotypes. These effects are contingent upon the context
433 of the primary variant, secondary variant class and function, phenotype of interest, and cohort
434 ascertainment. Our results emphasize the importance of assessing a full spectrum of genomic
435 variants in a variety of contexts to unravel the etiology of heterogeneous clinical features
436 observed in individuals with the same primary variants.

437 Several of our results expand on previous work to identify mechanisms for variable
438 expressivity of pathogenic variants and may help refine the broader roles of modifier variants in
439 complex disorders. First, we identified roles for a wide set of rare variants towards
440 developmental features of the 16p12.1 deletion, including noncoding variants and STRs. These
441 findings expand previous definitions of secondary variants beyond additional CNVs or rare
442 coding SNVs^{28,33}. Second, we expanded the role for PRS towards various developmental and
443 psychiatric phenotypes in individuals with pathogenic CNVs. These findings are in line with
444 recent studies that have identified roles for PRS as modifiers of specific phenotypes of
445 pathogenic CNVs, such as BMI z-scores for 16p11.2 CNVs¹⁰. Third, we observed cases of
446 compounding variant burden across generations of 16p12.1 deletion carriers, which could
447 explain our previous findings correlating rare variant burden with family history of psychiatric
448 disorders^{28,54}. This phenomenon could be attributed to assortative mating on psychiatric features
449 between deletion and non-deletion parents; in fact, we recently reported that 16p12.1 deletion
450 spouse pairs show strong correlations for psychiatric disorders, which may lead to increasing
451 genetic risk over generations⁶⁸. Fourth, while disease-relevant secondary variants contributed to
452 multiple genetic diagnoses, we did not find any instance where a single variant solely accounted
453 for all phenotypes observed in a proband, suggesting that secondary variants modify effects of
454 the deletion.

455 Ascertainment bias can preclude a more thorough assessment of a full spectrum of
456 phenotypes due to primary variants, including subtler or progressive effects, as deeper
457 evaluations are typically restricted to individuals with specific disorders^{69,70}. We therefore
458 examined the effects of cohort ascertainment within a single primary variant-specific context.
459 While the 16p12.1 deletion contributed to clinical outcomes across disease-ascertained and

460 general population cohorts, the specific phenotypic trajectories and the influence of secondary
461 variants both differed substantially across ascertainment. Our findings have several
462 implications, as management of primary variant-related symptoms may differ in individuals
463 evaluated for severe developmental features versus those with other medical features, where the
464 variant may first present as an incidental finding¹⁹. Thus, a shift in treatment focus from just the
465 primary variant to all variants in an affected individual could allow for more effective
466 management of individuals who carry primary variants.

467 The observed variability among secondary variant-phenotype patterns in each cohort
468 makes it difficult to identify consistent patterns across ascertainment, limiting the
469 generalizability of variant association studies conducted in a single cohort. In fact, the only
470 consistent trends we observed across primary variants, such as 16p11.2 deletion and rare disease-
471 associated SNVs, were for reduced IQ correlating with increased rare secondary variants,
472 mirroring previous studies¹¹. Joint models that integrate data across cohorts with appropriate
473 covariates can be used to overcome this ascertainment bias; for example, we found several
474 significant associations using joint models of 16p12.1 deletion carriers, including rare SNV
475 burden towards ADHD and speech delay in affected children. More broadly across primary
476 variant and ascertainment, we observed general trends for more PRS effects towards psychiatric
477 features, such as autism PRS and ADHD in DD children and SCZ PRS for mental/behavioral
478 disorders in UKB, and more rare variant effects towards cognitive features, such as rare deletions
479 (LF) with full-scale IQ in SSC probands with disease-associated SNVs and 16p11.2 deletion
480 probands. Both trends mirror previous variant-phenotype associations in individuals with autism
481 outside of a primary variant context⁷¹. However, some exceptions to these patterns exist: SSC
482 probands with rare duplications and 16p11.2 duplication probands show higher effects of rare
483 secondary variants towards psychiatric phenotypes, including negative associations of splice
484 variants with externalizing behavior in both groups. Further, 16p12.1 deletion carriers from
485 SPARK and 16p11.2 deletion probands show greater effects of PRS towards cognitive features,
486 such as the positive association of SCZ PRS and full-scale IQ in 16p11.2 deletion probands and a
487 negative association of SCZ PRS with learning disability in SPARK, potentially a facet of
488 autism-specific genetic etiology for cognitive features. Therefore, when describing genotype-
489 phenotype associations in a primary variant context, future studies should strive to assess

490 multiple independent cohorts with different ascertainties to determine the extent that
491 ascertainment could bias their results.

492 Despite assessing the contributions of multiple secondary variant classes towards specific
493 clinical features of pathogenic variants, much of the genetic etiology for these features is still not
494 accounted for in our study. Several factors could account for the unexplained variance, including
495 environmental factors or additional variant classes such as inversions, as well as those that could
496 explain ascertainment variability, including population-specific effects in the genetic
497 background. Another under-studied source of the unexplained variance could be non-additive
498 interactions between variants. Only a small number of synergistic variant interactions have been
499 identified to date in complex genomic disorders⁷², and very large cohorts will be required to
500 quantify the effects of these interactions towards clinical features⁷³. Here, we identified non-
501 additive effects of primary and secondary variants among children ascertained for autism.
502 Molecular studies could help unravel the mechanisms by which modifier variants interact with
503 primary variants to influence their phenotypes. For instance, we previously found 11 cases where
504 secondary variants synergistically altered the expression of genes dysregulated by the deletion in
505 patient-derived LCL models⁵⁴. While the overall effects of such interactions will likely explain
506 only a portion of the unexplained variance, they may play an outsized role in CNV disorders due
507 to the potential for interactions among multiple genes within the primary variant²⁶.

508 Overall, we identified family-, phenotype-, ascertainment-, and primary variant-specific
509 patterns of secondary variants that influence the variable expressivity of the 16p12.1 deletion and
510 other primary variants. Our study emphasizes the complexity of neurodevelopmental disorders
511 even after a causal variant is identified, suggestive of an oligogenic model for disease
512 pathogenicity⁷⁴. For researchers and clinicians alike, our study highlights the importance of
513 understanding the influence of cohort ascertainment and thoroughly investigating genomic
514 variants with smaller effect sizes. The complexity of the 16p12.1 deletion and other genomic
515 disorders calls for personalized medicine approaches that fully account for individual-level
516 phenotypic presentation, family history, and genome-wide variant profile towards counseling,
517 management, and potential treatment.

518
519
520

521 **Limitations of this study**

522 One limitation of our study is the relatively low sample size of families with the 16p12.1 deletion
523 and other primary variants. While this study represents one of the largest cohorts of individuals
524 with the same pathogenic variant to date and is well-powered for assessing changes in rare
525 coding SNV burden among deletion carriers, the study is under-powered for identifying
526 enrichments of individual variants or genes towards specific phenotypes. Additionally, while our
527 study captures major themes regarding variable expressivity of pathogenic variants, some
528 specific associations have only nominal significance. Larger cohorts will allow for further study
529 of these trends and could uncover roles for specific genes or molecular pathways towards clinical
530 features. Finally, while we were able to leverage data from 773 16p12.1 deletion individuals
531 from multiple ascertainment, differences in genotyping and phenotyping methods precluded
532 direct comparisons between cohorts.

533

534 **AUTHOR CONTRIBUTIONS**

535 M.J., C.S., A.T., L.P., and S.G. designed the study and analyses. C.S., L.P., E.H., and L.R.
536 recruited patients and conducted interviews, harmonized phenotypic data from interviews and
537 medical records, and extracted DNA from blood for WGS sequencing. C.T. and C.L.M. assisted
538 with collection and analysis of quantitative phenotypic data. Other authors provided de-identified
539 DNA, blood samples, or genomic and phenotypic data of 16p12.1 deletion families to the study.
540 M.J., C.S., A.T., D.B., and V.K.P. designed bioinformatics pipelines to identify variants,
541 uniformly processed sequencing data from 16p12.1 deletion and external cohorts, and performed
542 all statistical, enrichment, and modeling analysis. V.R. and H.S. assisted with design of the PRS
543 calculations and modeling approaches. M.J., C.S., A.T., L.P., and S.G. wrote the manuscript with
544 approval from all authors.

545

546 **ACKNOWLEDGEMENTS**

547 This work was supported by NIH R01-GM121907 and resources from the Huck Institutes of the
548 Life Sciences to S.G. M.J. and C.S. were supported by NIH T32-GM102057. A.T. was supported
549 by NIH T32-LM012415. L.P. was supported by Fulbright Commission Uruguay-ANII. A.R. was
550 supported by grants from the Swiss National Science Foundation 31003A_182632. S.B. was
551 supported by the NIHR Manchester Biomedical Research Centre (NIHR203308). We thank

552 Craig Praul and the Penn State Genomics Core Facility for assistance with designing the WGS
553 sequencing strategy, Abby Hare-Harris (Geisinger ADMI) for assistance with RedCap analysis
554 of quantitative phenotypic data, Veera Rajagopal (Regenron Pharmaceuticals) and Bertrand
555 Isidor (CHU Nantes) for useful comments on the manuscript, and Jianyu Yang, Sarah Dwiekat,
556 and Edmundo Torres-Rodriguez (Penn State) for assistance with curating annotation data for
557 variant analysis. We are grateful to all of the families in each cohort (DD, SPARK, Simons
558 Searchlight, SSC, MyCode, and UK Biobank) as well as clinical sites and staff who participated
559 in the study. We thank the SSC principal investigators (A. Beaudet, R. Bernier, J. Constantino,
560 E. Cook, E. Fombonne, D. Geschwind, R. Goin-Kochel, E. Hanson, D. Grice, A. Klin, D.
561 Ledbetter, C. Lord, C. Martin, D. Martin, R. Maxim, J. Miles, O. Ousley, K. Pelphrey, B.
562 Peterson, J. Piggot, C. Saulnier, M. State, W. Stone, J. Sutcliffe, C. Walsh, Z. Warren, E.
563 Wijsman) as well as the Simons Searchlight Consortium. We appreciate obtaining access to
564 genomic and phenotypic data on SFARI Base. Approved researchers can obtain the SPARK,
565 SSC, and Simons Searchlight population data sets described in this study by applying at
566 <https://www.base.sfari.org>. This research has been conducted using data from UK Biobank.
567 More information about UK Biobank is available at <https://www.ukbiobank.ac.uk/>.

568

569 **DECLARATION OF INTERESTS**

570 The authors declare no competing interests.

571 **METHODS**

572 **16p12.1 deletion developmental delay cohort description**

573 We analyzed genomic and phenotypic data from a cohort of 452 individuals belonging to 132
574 families with the 16p12.1 deletion, which we refer to as the Developmental Delay cohort (“DD
575 cohort”). (**Fig. 1, Table S1A**). These families comprised single probands (n=14), parent-child
576 pairs (n=13), complete trios (n=97), and extended families, including eight families with three
577 generations and 22 with multiple affected children. The deletion was identified through prior
578 clinical diagnostic tests for ID/DD or other developmental disorders in probands. We note that
579 10 individuals from eight families with the 16p12.1 deletion representing an unselected
580 population from the Estonian BioBank were included as a comparison group within the DD
581 cohort⁶². Whole genome sequencing was performed on 287 individuals (107 probands),
582 including the eight Estonian BioBank families, and microarray experiments were performed for
583 368 individuals. Informed consent was obtained from families recruited directly according to a
584 protocol approved by the Pennsylvania State University Institutional Review Board (IRB
585 #STUDY00000278), while de-identified information was obtained from families recruited
586 through clinics according to another approved protocol (IRB #STUDY00017269). A list of all
587 individuals in the DD cohort, including familial membership and 16p12.1 deletion status, is
588 available in **Table S1A**. Note that some information that could be construed as personally
589 identifiable is summarized in **Table S1A** (i.e., age ranges instead of age values); full datasets are
590 available upon request from the authors for purposes of reproducibility (i.e., for the model
591 covariates). We also note that sample and family identifiers used in **Table S1A** are specific for
592 this study and were not known outside of the research group. Power calculations for detecting
593 changes in rare variant burden within deletion family members (**Fig. S1B**) were based on
594 estimated effect sizes of burden differences between 16p12.1 deletion probands and parents
595 (coding SNVs and CNVs) or between autism probands and parents (non-coding SNVs and
596 STRs) from previous studies^{28,75,76}.

597

598 **Phenotypic analysis**

599 Individual-level phenotypic data described below (phenotypic domain complexity scores,
600 quantitate phenotypes, and family history status) are available in **Table S1A**.

601 *(a) Collection and analysis of clinical features:* We collected detailed medical history
602 and used clinician-, guardian- (for children), or self- (for adults) reported standardized
603 questionnaires to assess developmental phenotypes in children (average age=10.1 years) and
604 psychiatric features in adults. Questionnaires for children assessed neuropsychiatric and
605 developmental features, anthropomorphic measures, congenital abnormalities in multiple organ
606 systems, and family history of medical or psychiatric disorders. Phone surveys were conducted
607 to fill in missing information, and families were recontacted every 3-4 years to track longitudinal
608 data and to note any later-onset clinical features. We analyzed phenotypes for each individual by
609 calculating “complexity scores” for clinical features within specific domains, as well as
610 measured specific quantitative features (see “Assessment of quantitative phenotypes” below).

611 *For children,* we first grouped clinical features into six broadly defined phenotypic
612 domains: (i) ID/DD, (ii) behavioral phenotypes, (iii) psychiatric features, (iv) nervous system
613 defects, (v) craniofacial and skeletal abnormalities, and (vi) congenital abnormalities (**Fig. 2A,**
614 **Table S1A-B**). We determined complexity scores ranging from 0 to 4 or 5 for each phenotypic
615 domain by assessing the total number of affected phenotypic sub-domains in each child. The full
616 list of phenotypes considered for each sub-domain is available in **Table S1B**. Presence of at least
617 one clinical feature within a sub-domain added an additional point of complexity to the total
618 score, but additional phenotypes within the same sub-domain added no additional complexity
619 score. For example, proband PIC_001 had three nervous system-related phenotypes (tremors,
620 abnormal gait, and abnormal brain morphology) grouped into two sub-domains (nervous system
621 abnormalities and nervous system morphology defects), and therefore received two points for
622 complexity. We note that younger probands were not assessed for psychiatric domains based on
623 the typical age of onset, such as for schizophrenia (**Table S1C**). As most probands (92%)
624 exhibited >1 feature within the ID/DD domain, we focused on the other five phenotypic domains
625 for downstream analysis.

626 *For adult family members,* we binned and calculated complexity scores in a similar
627 manner for clinical features within four domains, including (i) cognitive (ID, learning difficulty),
628 (ii) psychiatric (schizophrenia, bipolar disorder), (iii) depression/anxiety, and (iv) addiction
629 phenotypes (**Fig. 2B**). We note that most phenotypes for early-onset behavioral and
630 developmental features assessed in children were not examined for adult family members.

631 **(b) Assessment of quantitative phenotypes:** We performed online quantitative
632 assessments using the Hansen Research Services Matrix Adaptive Test (HRS-MAT) for non-
633 verbal IQ³⁴ and Social Responsiveness Scale (SRS) for autism-related social behavior³⁵ (**Fig.**
634 **2C**). HRS-MAT was self-administered to participants through an online platform, while the SRS
635 was administered through a RedCap-based survey platform maintained by the Geisinger Autism
636 and Developmental Medicine Institute. The SRS assessment was self-reported if the participant
637 was over 18 years or completed by parents or guardians for children under the age of 18 years.
638 Body Mass Index (BMI) and head circumference were obtained from medical records or
639 self/guardian-reports or, for BMI, calculated from height and weight data obtained from medical
640 records or self/guardian-reports. Both BMI and head circumference were converted into age- and
641 sex-adjusted z-scores^{77,78}. We further obtained SRS, BMI, and head circumference z-scores for
642 probands in the SRS and Simons Searchlight cohorts (see below), while the mean HRS-MAT
643 score in SSC probands was obtained from Hansen, 2019³⁴. Differences in phenotype
644 distributions between groups of 16p12.1 deletion family members and between sets of probands
645 with different primary variants were calculated using one- and two-tailed Mann Whitney tests,
646 respectively (**Table S2A**). Differences of proband scores from a defined mean were calculated
647 using one-tailed one-sample t-tests (**Table S2A**).

648 **(c) Developmental milestones:** We assessed the achievement of developmental
649 milestones in children from the DD cohort based on CDC guidelines³⁶. Parents/guardians of
650 children reported the ages at which children achieved 12 milestones, including age first smiled,
651 rolled over, crawled, walked, and spoke. Age of milestone attainment for all available samples is
652 reported in **Table S1A**. Differences in milestone achievement between probands and their
653 siblings/cousins were assessed using one-tailed t-tests (**Table S2B**).

654

655 **DNA extraction and whole-genome sequencing**

656 We performed DNA extraction and whole-genome sequencing on 287 individuals in the DD
657 cohort (**Table S1A**). Genomic DNA was extracted from peripheral blood samples from some
658 participants using the QIAamp DNA Blood Maxi extraction kit (Qiagen, Hilden, Germany)
659 while clinical collaborators sent isolated DNA from other participants. Illumina TruSeq DNA
660 PCR-free libraries (San Diego, CA, USA) were constructed for 150bp paired-end whole-genome
661 sequencing using Illumina HiSeq X by Macrogen Labs (Rockville, MD, USA). Samples were

662 sequenced at an average 35.7X coverage, or 716.2 M reads/sample, with 94.9% of reads
663 mapping to the human genome. After processing for quality control using Trimmomatic⁷⁹
664 (leading:5, trailing:5, and slidingwindow:4:20 parameters), sequences were aligned to the hg19
665 reference genome using BWA v.0.7.13⁸⁰, and sorted and indexed using Samtools v.1.9⁸¹. We
666 note that sequencing data from 163 individuals was described previously^{54,68}.

667 We used SNP microarrays for copy-number variant validation and genotyping
668 experiments (i.e. CNV calling and polygenic risk score calculation) for 368 individuals. Samples
669 were run on Illumina OmniExpress 24 v.1.1 microarrays at Northwest Genomics Center in the
670 University of Washington (Seattle, WA, USA). We note that microarray data of 208 individuals
671 was described previously^{28,54,68}.

672

673 **Identification and annotation of single-nucleotide variants**

674 We identified SNVs and small indels using the GATK Best Practices pipeline⁸², followed by
675 quality control and extensive variant and gene-level annotations. Duplicate removal with
676 PicardTools was followed by base-pair quality score recalibration and variant calling for each
677 sample using GATK HaplotypeCaller v.3.8. We then merged calls from all samples using GATK
678 GenotypeGVCFs v.4.0.11, performed variant quality score recalibration to finalize variant calls,
679 and used Vcfanno to annotate variants with GnomAD frequency^{83,84}. All calls were filtered for
680 QUAL ≥ 50 , allele balance between 0.25 and 0.75 or ≥ 0.9 , read depth ≥ 8 , QUAL/alternative read
681 depth ≥ 1.5 , GnomAD frequency $\leq 0.1\%$ (or not present), and intracohort frequency ≤ 10 to
682 account for technical differences between our data and GnomAD. We annotated coding and
683 noncoding variants within genes from GENCODE⁸⁵ v19 using ANNOVAR⁸⁶ and Vcfanno⁸³ as
684 follows: (a) *Coding SNVs*: Rare coding variants were filtered for loss-of-function (LOF),
685 missense, or splicing exonic variants in protein-coding genes, and annotated with CADD Phred-
686 like scores, presence in ClinVar database, the gene-level pathogenicity metric LOEUF, and gene-
687 level phenotype associations using OMIM^{37,42,43,87}. Missense and splice variants were filtered to
688 include only those with a CADD score ≥ 25 . (b) *Noncoding SNVs*: All rare variants located 1kbp
689 upstream of a gene transcription start site were classified as promoter variants, while genes
690 within the 5' UTR were classified as 5' UTR variants. Fetal brain-active enhancer regions were
691 identified using chromatin state data from the Roadmap Epigenomics consortium⁸⁸ (states 6, 7,

692 and 12 in the fetal brain), and rare variants in those regions were classified as enhancer variants.
693 Rare SNV burden for all available individuals are listed in **Table S1A**.

694

695 **Identification of copy-number variants and short tandem repeats**

696 CNVs were called from both microarray data using PennCNV⁸⁹ and WGS data using CNVnator
697 v.0.4.1, Lumpy-sv v.0.2.13 with Smoove v.0.2.5, Delly v.1, and Manta v.1.6.0⁹⁰⁻⁹³. For CNVs
698 >50 kbp in length, we used a union of PennCNV and CNVnator calls. For CNVs <50 kbp, we
699 used CNVs called by least two of CNVnator, Lumpy, Manta, or Delly, defined by 50%
700 reciprocal overlap. All CNVs were annotated for 50% reciprocal overlap with known pathogenic
701 CNVs³³. WGS-based CNV calls were filtered to remove calls with >50% overlap with
702 centromeres, segmental duplications, regions of low mappability, and V(D)J recombination
703 regions, while microarray CNVs were filtered to remove samples with >50% overlap with
704 centromeres, telomeres, and segmental duplications⁹⁴. Known pathogenic CNVs were excluded
705 from this filter³³. All CNVs were then filtered for GnomAD-SV frequency⁹⁵ <0.1% and
706 intracohort frequency ≤ 10 , and >50 kb CNVs were additionally filtered for <0.1% frequency in a
707 control cohort⁹⁶. All CNVs were finally filtered for those intersecting at least one protein-coding
708 gene, using gene annotations from GENCODE v19⁸⁵.

709 We identified STR expansions from WGS data using the GangSTR v.2.4 (reference file
710 v.13.1) and TRTools pipelines^{97,98}. We filtered calls with read depth >20 and <1000, excluding
711 reads that were not spanning and bounding the STR locus and calls with maximum likelihood
712 estimates not within the confidence interval using dumpSTR. After merging the calls with
713 mergeSTR, we ran dumpSTR with population level filters, including locus call rate >0.8 and
714 departure from Hardy Weinberg equilibrium (Fisher's exact p-value) >0.00001, and removed
715 loci that overlapped with segmental duplications. For chromosome X, the Hardy-Weinberg
716 equilibrium p-value was calculated from female samples only. For each family, we extracted the
717 STR loci that passed variant filtering, and used GangSTR v2.5 to call STR variants, which were
718 used in subsequent analyses. We defined STR expansions as STR variants with lengths >2SD
719 higher than the average of STR lengths among all individuals in our in-house cohort of
720 individuals with WGS data at a particular locus. STR expansions spanning protein coding
721 regions defined by GENCODE v19⁸⁵ were selected for downstream analysis. All CNV and STR
722 genes were further annotated for pathogenicity metrics (LOEUF score³⁷). The number of genes

723 affected by rare CNVs and the number of STR expansions for all individuals with available
724 WGS data are listed in **Table S1A**.

725

726 **Polygenic risk score calculations**

727 Using microarray data, we calculated polygenic risk scores for educational attainment³⁸,
728 intelligence³⁹, schizophrenia⁴⁰, and autism⁴¹ among the samples in the DD cohort, based on
729 standardized bioinformatics pipelines for quality control⁹⁹. We first downloaded summary
730 statistics from four recent GWAS studies of neuropsychiatric traits, and filtered SNPs for
731 imputation INFO scores >0.8 and removed duplicate and ambiguous SNPs. We then merged
732 SNP genotype data from different microarray batches together using PLINK. Initial quality
733 control removed SNPs with minor allele frequency <0.05, Hardy-Weinberg equilibrium <1.0×10⁻⁶,
734 and genotype rate <0.01, along with samples missing >1% of genotypes. We used the HRC-
735 1000 Genomes Imputation toolkit (<https://www.well.ox.ac.uk/~wrayner/tools>) to process PLINK
736 files into individual chromosomes for imputation, and VcfCooker

737 (<https://genome.sph.umich.edu/wiki/VcfCooker>) to convert PLINK files to VCF files.

738 Microarray-based SNPs were imputed using the TOPMed v.r2 imputation server using Eagle
739 v2.4 for phasing¹⁰⁰. After imputation, VCF files were converted back to PLINK format, and
740 SNPs were again filtered with identical QC filters. Additional QC filters included removing
741 samples with >±3SD of the mean heterozygosity rate. We also selected individuals with
742 European ancestry, based on imputed genetic ancestry (calculated using Peddy v.0.4.8 with 1000
743 Genomes-based population panel) or self-reported ancestry, for downstream analysis¹⁰¹. Finally,
744 we performed strand-flipping to match SNPs in microarray data with the GWAS summary
745 statistics. To calculate PRS, we used standardized pipelines for the LDpred2 software package,
746 which uses Bayesian approaches to optimize parameters for PRS calculation¹⁰². Briefly, we
747 filtered the four sets of GWAS summary statistics for SNPs present in the HapMap3 dataset¹⁰³,
748 and used 1000 Genomes datasets to calculate linkage disequilibrium matrices for the SNPs⁹⁹.
749 After regressing betas or odds ratios of GWAS SNPs according to linkage disequilibrium, we
750 used the LDpred2-auto model to calculate the four PRS values for all samples with available
751 genotype data. PRS for all available individuals are listed in **Table S1A**.

752

753 **Genotype-phenotype statistical and modeling analysis**

754 We performed multiple analyses to compare effects of rare variants and PRS towards different
755 phenotypic domains among 16p12.1 deletion probands or between probands and their carrier and
756 non-carrier parents. When assessing variant effects towards phenotypic domains, we binned
757 probands with different complexity scores for each phenotypic domain into binary groups of
758 roughly equal sizes (i.e., probands with higher versus lower complexity scores) for logistic
759 regression models and burden tests. Burden analysis (paired and independent T-tests) and
760 Pearson's correlation analyses were calculated in Python v.3.7 using the `scipy v1.13.1 ttest_rel`,
761 `ttest_ind`, or `pearsonr` functions, respectively. We note that paired t-tests for PRS burden between
762 probands and parents were two-tailed, due to the dual directionality of PRS for different
763 phenotypes, while other t-tests were one-tailed. Benjamini-Hochberg multiple testing correction
764 was performed using the `scipy v.1.13.1 false_discovery_control` function for all statistical
765 analyses unless otherwise stated. Multiple testing was performed separately for analyses with all
766 sets of rare variants and variants filtered for evolutionary constraint (defined by $LOEUF < 0.35$,
767 which are intolerant to loss-of-function variants in the general population³⁷; referred to as
768 "(LF)"). FDR values reported in the text are corrected for multiple testing, while p-values are not
769 corrected for or did not pass multiple testing. Sample sizes, test statistics, and corrected and
770 uncorrected p-values for all analyses are available in **Tables S2-S5**.

771 Logistic and linear regression models for phenotypic variation among probands were
772 performed using the *Logit* and *OLM* functions in `statsmodels v.0.14.2`, respectively (**Table S3A**).
773 For joint variant regression models, we used three different sets of genetic input variables to test
774 for effects towards phenotypes: (a) all rare variants (sum of SNV, STR, and CNV gene burden)
775 and schizophrenia PRS; (b) SNVs, STRs, duplications, and deletions; and (c) SNVs, STRs,
776 duplications, and deletions restricted to genes with $LOEUF$ scores < 0.35 (referred to in the
777 figures as "LF model"). Single variant regression models used only a single variant class as
778 input. All models also included sex as a covariate, while models *b* and *c* also included
779 schizophrenia PRS as a covariate. Additional covariates, such as age and genotype PCs, were not
780 included due to concerns regarding potential overfitting of models with lower sample sizes. The
781 variance explained by the models (R^2 for linear models and McFadden's pseudo- R^2 for logistic
782 models) was calculated for all models. Sample sizes, odds ratios, uncorrected p-values,

783 confidence intervals, and variance statistics for all regression models used in the DD cohort are
784 available in **Tables S3A**.

785

786 **Variant enrichment and pathogenicity analysis**

787 *Gene set enrichment:* We assessed enrichment of genes with secondary variants among
788 sets of neurodevelopmental disease genes and genes with neuronal function from several
789 previously published resources^{37,44,45,48,104–106}. We identified enrichment of variants in these gene
790 sets by performing Fisher's Exact tests against the whole genome for each gene list and
791 calculated odds ratios and p-values for genes with variants in the DD cohort using the
792 *contingency.odds_ratio* function from *scipy* v.1.13.1 (**Table S2G**).

793 *Gene ontology:* Gene ontology (GO) term enrichment was performed using the Panther
794 API and the GO-Slim Biological Process annotation dataset¹⁰⁷ (**Table S3B**). The GO term
795 network figures were created using GO term semantic similarity calculated from *rrvgo*¹⁰⁸ v1.10.0
796 in R v4.2.3.

797 *Spatio-temporal brain expression:* We assessed variant enrichment in genes
798 preferentially expressed in specific brain tissues using the BrainSpan Atlas⁵² and in genes
799 preferentially expressed in specific cell types using single-cell RNA-seq expression data in the
800 M1 motor cortex⁵³. We defined preferentially expressed genes as those with expression >2SD
801 higher than the median expression across all tissues or all cell types for that gene. We used
802 Fisher's exact tests as described above to find the odds that a gene both carries a variant in the
803 DD cohort and is expressed in a specific brain region or cell type (**Table S2E, S2H**).

804 *16p12.1 differentially and co-expressed genes:* We used Fisher's Exact tests as described
805 to calculate enrichment of secondary variants in gene co-expression modules previously
806 identified using WGCNA in lymphoblastoid cell line (LCL) models of the 16p12.1 deletion⁵⁴
807 (**Table S2F**). We specifically assessed six co-expression modules whose genes showed
808 differential expression between deletion carriers and controls, one of which (black module) also
809 contained two genes within the 16p12.1 deletion region (*POLR3E*, *MOSMO*).
810 For all enrichments, we applied Benjamini-Hochberg multiple testing correction as described
811 above. Sample sizes, test statistics and corrected and uncorrected p-values for enrichments are
812 listed in **Tables S2** and **S3**.

813 *Pathogenic variant analysis:* We defined pathogenic SNVs as those that are “Pathogenic”
814 or “Likely pathogenic” in ClinVar for neurodevelopmental phenotypes⁴², or loss-of-function
815 variants in genes that (a) have dominant or (if male) X-linked neurodevelopmental OMIM
816 phenotype⁴³, (b) are a Tier S or Tier 1 SFARI gene, which represent strong candidate autism
817 genes⁴⁵, or (c) are in the Tier 1 or Tier 2 gene list from the Developmental Brain Disorder Gene
818 Database, which represent genes with well-documented connections to neurodevelopmental
819 disease⁴⁴ (**Table S1D**). Pathogenic CNVs were identified from 50% reciprocal overlap with a
820 previously published list of CNVs³³.

821

822 **Network analysis**

823 We assessed the connectivity of genes within a previously described brain-specific interaction
824 network. In brief, the network was built using a machine-learning model trained to predict the
825 likelihood of interactions between pairs of genes using brain-specific gene co-expression,
826 protein-protein interaction, and regulatory sequence datasets^{49,50}. We restricted the network to
827 genetic interactions with weighted probabilities >2 and calculated the degree for each gene (or
828 number of connections between a particular gene and other genes) as a descriptor of connectivity
829 to other genes in the network. We then binned each gene into one of four quantiles based on the
830 gene’s degree of connectivity and counted the number of times that genes with coding variants
831 fell into each quantile. To calculate empirical p-values, we compared the resulting values to 1000
832 simulations in which we randomly selected the same number of genes from the genome and
833 counted the number of randomly selected genes that fell into each quantile (**Table S2D**).

834

835 **Genotype-phenotype analysis in 16p12.1 deletion samples from other ascertainment**

836 (a) *Description of cohorts:* We identified individuals carrying 16p12.1 deletions from three
837 additional cohorts, each representing a distinct ascertainment. Individuals in the Simons
838 Powering Autism Research for Knowledge (SPARK) cohort (n=94) were ascertained for families
839 with autism⁵⁸, while the Geisinger MyCode Community Health Initiative (MyCode) (n=160)
840 represents a health care-based cohort⁶⁰ and the UK Biobank⁵⁹ (UKB) (n=250) consists of
841 individuals with a “healthy volunteer” bias⁶¹. Combined with samples from the DD cohort (after
842 excluding samples with incomplete phenotypic information), we assessed a total of 757 deletion
843 carriers. De-identified data from these cohorts were obtained and analyzed according to a

844 protocol approved by the Pennsylvania State University Institutional Review Board (IRB
845 #STUDY00011008). Data from the UK Biobank was accessed under application 45023.
846 Individuals from the MyCode cohort were recruited during primary care or specialty clinic visits
847 to Geisinger Health System locations, independent of condition, diagnosis, or demographic
848 characteristic. Written informed consent was obtained from adult patients and from the parents or
849 guardians of pediatric patients. The study was conducted with approval from the Geisinger
850 institutional review board.

851 *(b) CNV calling:* Carriers of the 16p12.1 deletion in each cohort were identified based on
852 CNV calls from microarray data. Samples from MyCode were genotyped using the Illumina
853 Global Screening Array and Illumina OmniExpressExome-8 Kit. SNP log-r ratio and b-allele
854 frequencies for SPARK samples were downloaded through the SFARI Base portal
855 (<https://www.base.sfari.org>), while signals for the UK Biobank were accessed through Data
856 Fields 22437 and 22431. CNVs for all cohorts were called using the PennCNV pipeline
857 described above for the DD cohort⁸⁹. Additionally, 2,640 and 356 additional samples without any
858 large (>500kb), rare (<0.1% population frequency) CNVs were identified as controls for
859 additional genetic analysis from the UK Biobank and SPARK, respectively.

860 *(c) SNV calling from sequencing data:* SNVs for all three cohorts were identified from
861 whole exome sequencing (WES) data. NimbleGen (SeqCap VCRome) and xGEN probes from
862 Integrated DNA Technologies (IDT) were used for target sequence capture in the MyCode
863 cohort^{109,110}. Sequencing was performed by paired end 75bp reads on an Illumina NovaSeq or
864 HiSeq at coverage >20x at >80% of the targeted bases. Alignments and variant calling were
865 based on GRCh38 human genome reference sequence. Variants were called with the WeCall
866 variant caller version 1.1.2 (<https://github.com/Genomicsplc/wecall>). Whole exome VCFs for
867 SPARK samples were downloaded through the SFARI Base portal (<https://www.base.sfari.org>).
868 WES VCFs from both cohorts were processed using the same pipeline described above for the
869 DD cohort. WES data from UKB individuals was available as multi-sample project VCFs¹¹¹ in
870 the UK Biobank Research Analysis Platform. After splitting multi-allelic records, we applied the
871 following set of quality control filters using Hail in the DNAnexus platform: (a) variant call rate
872 >90%, (b) Hardy Weinberg equilibrium p-value >10⁻¹⁵, (c) minimum read depth > 10, and (d) at
873 least one sample passing the allelic balance threshold of 0.2. Next, we filtered for variants with
874 an intracohort frequency <0.1% and present in at least two samples. The remaining variants were

875 then annotated using Variant Effect Predictor¹¹² (VEP) v.109 and dbNSFP¹¹³ v.4 to identify their
876 effects on gene transcripts. We specifically annotated variants based on VEP annotations as LOF
877 (transcript ablation, stop gained, frameshift variant, stop lost, and start lost), missense, or splice
878 (splice acceptor variant and splice donor variant). Missense variants were further filtered for
879 those predicted to be deleterious by at least five of nine selected tools (SIFT, LRT, FATHMM,
880 PROVEAN, MetaSVM, MetaLR, PrimateAI, DEOGEN2, and MutationAssessor) available
881 through the dbNSFP database¹¹³.

882 *(c) Phenotype analysis:* We assessed phenotypic data in these cohorts using ICD10 codes
883 derived from electronic health records (MyCode and UKB) and self-reported questionnaire
884 responses (UKB and SPARK) (**Table S4A**). Phenotypic information from SPARK was
885 downloaded from the Simons Foundation through the SFARI Base portal
886 (<https://www.base.sfari.org>). Electronic health records were available from participants in
887 MyCode⁶⁰. Electronic health records for UKB were identified from Data Fields 41202 and 41204
888 (main and secondary inpatient ICD10 codes), while questionnaire data was identified from
889 additional Data Fields (**Table S4A**). For harmonization of phenotypic data across cohorts, ICD10
890 codes were matched to phenotypes from questionnaires, details of which are provided in **Table**
891 **S4A**.

892 *(d) Secondary variant associations:* Comparisons of secondary variant burden between
893 16p12.1 deletion carriers and age and sex matched controls in the UK Biobank and SPARK were
894 assessed using two-tailed t-tests. The relationship of secondary variant burden and phenotypes in
895 all three cohorts, and comparison with adults and children in the DD cohort, were assessed using
896 two-tailed t-tests. We note that phenotypes in these cohorts were only assessed if they were
897 present in at least five individuals or 10% of the cohort, whichever was larger, and the cohort had
898 at least five or 10% of the cohort controls. We further assessed the effects of secondary variant
899 on phenotypes across cohorts using logistic regression using the model structures *(b)* and *(c)*
900 described above for the DD cohort, without the inclusion of STR variants. Phenotypic logistic
901 regression was performed on main and secondary ICD10 codes with age and sex included as
902 covariates. PheWAS was performed using the *PheWAS* v.0.99.6-1 package in R¹¹⁴ on all
903 available samples from the UK Biobank using Phecodes derived from ICD9 and ICD10 codes,
904 identified from Data Fields 41202, 41204, 41203, and 41205 (main and secondary inpatient
905 ICD9/10 codes), while correcting for sex, age, and the top four genetic principal components.

906 Sample sizes, test statistics, uncorrected and corrected p-values, confidence intervals, and
907 variance statistics for all analyses are available in **Table S4**.

908

909 **Genotype-phenotype analysis in other neurodevelopmental disease cohorts**

910 We assessed genomic and phenotypic data from individuals from the Simons Searchlight
911 project⁶³, ascertained for probands with 16p11.2 deletions and duplications, and Simons Simplex
912 Collection (SSC), ascertained for families with simplex cases of autism⁶⁴. Within the Simons
913 Searchlight cohort, we assessed 128 probands with the 16p11.2 duplication (n=37) or deletion
914 (n=91). Within the SSC cohort, we assessed genomic data of 2,435 total probands, and classified
915 probands with the following primary variant classes for downstream analysis: (i) 1,237 probands
916 with rare, deleterious variants (<0.1% frequency, loss-of-function or missense variants with
917 CADD Phred-like scores >25); (ii) 79 probands with large rare deletions (<0.1% population
918 frequency, >500kbp); (iii) 148 probands with large rare duplications (<0.1% population
919 frequency, >500kbp); and (iv) 1,084 probands who did not carry any of these variants or any
920 other known pathogenic CNVs³³. We note that groups with primary variants have overlapping
921 samples. Additionally, we assessed phenotypic data for an additional 419 SSC probands and 32
922 Simons Searchlight probands to compare SRS distributions with 16p12.1 deletion probands (a
923 total of 2,844 total SSC probands and 139 Simons Searchlight probands were used in these
924 analyses). De-identified data from these cohorts were obtained and analyzed according to a
925 protocol approved by the Pennsylvania State University Institutional Review Board (IRB
926 #STUDY00011008). In sum, we assessed secondary variants and phenotypes in 2,252
927 individuals with primary variants from six cohorts: DD (n=269), SPARK (n=94), MyCode
928 (n=160), UKB (n=250), Searchlight (n=128), and SSC (n=1,351). We also assessed data from
929 1,084 SSC probands without primary variants. Data from an additional 311,980 control
930 individuals was also included, including non-carrier samples from DD (n=183), age and sex-
931 matched controls without CNVs from SPARK (n=356), PheWAS and age and sex-matched
932 controls without CNVs from UKB (n=310,990), 16p11.2 CNV samples from Searchlight and
933 SSC probands without genetic data but with SRS data (n=32 and 419, respectively).

934 Exome sequencing VCFs and raw microarray data for Searchlight cohorts, whole genome
935 sequencing VCFs and STR calls for SSC⁷⁵, and all phenotype data were downloaded through the
936 SFARI Base portal (<https://www.base.sfari.org>). We processed and filtered exome and WGS-

937 based SNVs and indels for the same quality control filters used to process the DD cohort, and
938 then annotated variants using our standardized pipeline. Short tandem repeat calls from SSC
939 were previously published by Mitra and colleagues⁷⁵ and were processed and filtered with the
940 same pipeline as our cohort. CNV calls from microarrays for SSC were previously published by
941 Sanders and colleagues¹¹⁵, while CNV calls from microarrays for the Searchlight cohort were
942 processed using PennCNV⁸⁹ as previously described²⁸. For this manuscript, genes within CNVs
943 were reannotated using GENCODE v19⁸⁵, but otherwise used as-is without additional
944 processing. Primary variant SNVs and CNVs were removed from secondary variant lists for
945 downstream processing. We further processed microarray data and calculated PRS for both
946 cohorts using the same pipelines as the DD cohort, except that autism PRS was not calculated in
947 SSC samples, as the underlying GWAS summary statistics were calculated in part using SSC
948 samples⁴¹. Finally, we curated results of quantitative phenotypic assessments for each cohort
949 from SFARI Base, including full-scale IQ, internalizing and externalizing behavior
950 (ABCL/CBCL), social responsiveness (SRS), autism-related behaviors (BSI, Searchlight only),
951 repetitive behavior (RBS-R, SSC only), coordination disorder (DCDQ, SSC only), BMI z-score,
952 and head circumference z-score (Searchlight only).

953 Linear regression models for assessing variation in these phenotypes were constructed
954 using the *OLS* function from statsmodels v0.14.2, using the same model structures (*b*) and (*c*)
955 described above for the DD cohort (note that Searchlight models did not include STRs). For
956 models investigating the interactions of primary and secondary variants, Benjamini-Hochberg
957 FDR correction was performed using the statsmodels v.0.14.2 *false_discovery_control* function.
958 All other correlations, statistical analyses, GO enrichments, and multiple testing corrections for
959 comparing variant classes and quantitative phenotypes were performed in the same manner as for
960 the DD cohort. Sample sizes, test statistics, uncorrected and corrected p-values, confidence
961 intervals, and variance statistics for all analyses are available in **Table S5**.

962

963 **Data and code availability**

964 Whole genome sequencing and SNP microarray data generated in this study are available at
965 NCBI dbGaP phs002450.v2.p1. All code generated for this project, including pipelines for
966 running bioinformatic software and custom analysis scripts, are available at
967 https://github.com/girirajanlab/16p12_WGS_project. Statistical analyses and experimental

968 results for the data presented in **Figs. 2-7** and associated supplementary figures are available in
969 **Tables S2-S6.**

970 **REFERENCES**

- 971 1. Claussnitzer, M., Cho, J.H., Collins, R., Cox, N.J., Dermitzakis, E.T., Hurles, M.E.,
972 Kathiresan, S., Kenny, E.E., Lindgren, C.M., MacArthur, D.G., et al. (2020). A brief
973 history of human disease genetics. *Nature* 577, 179–189. [https://doi.org/10.1038/s41586-](https://doi.org/10.1038/s41586-019-1879-7)
974 019-1879-7.
975
- 976 2. Kingdom, R., and Wright, C.F. (2022). Incomplete Penetrance and Variable Expressivity:
977 From Clinical Studies to Population Cohorts. *Front Genet* 13, 920390.
978 <https://doi.org/10.3389/fgene.2022.920390>.
979
- 980 3. Girirajan, S., and Eichler, E.E. (2010). Phenotypic variability and genetic susceptibility to
981 genomic disorders. *Hum Mol Genet* 19, R176-87. <https://doi.org/10.1093/hmg/ddq366>.
982
- 983 4. Posey, J.E., Harel, T., Liu, P., Rosenfeld, J.A., James, R.A., Coban Akdemir, Z.H.,
984 Walkiewicz, M., Bi, W., Xiao, R., Ding, Y., et al. (2017). Resolution of Disease
985 Phenotypes Resulting from Multilocus Genomic Variation. *N Engl J Med* 376, 21–31.
986 <https://doi.org/10.1056/NEJMoa1516767>.
987
- 988 5. Leitch, C.C., Zaghoul, N.A., Davis, E.E., Stoetzel, C., Diaz-Font, A., Rix, S., Alfadhel,
989 M., Lewis, R.A., Eyaid, W., Banin, E., et al. (2008). Hypomorphic mutations in
990 syndromic encephalocele genes are associated with Bardet-Biedl syndrome. *Nat Genet* 40,
991 443–448. <https://doi.org/10.1038/ng.97>.
992
- 993 6. Riordan, J.D., and Nadeau, J.H. (2017). From Peas to Disease: Modifier Genes, Network
994 Resilience, and the Genetics of Health. *Am J Hum Genet* 101, 177–191.
995 <https://doi.org/10.1016/j.ajhg.2017.06.004>.
996
- 997 7. Guo, T., Chung, J.H., Wang, T., McDonald-McGinn, D.M., Kates, W.R., Hawuła, W.,
998 Coleman, K., Zackai, E., Emanuel, B.S., and Morrow, B.E. (2015). Histone Modifier
999 Genes Alter Conotruncal Heart Phenotypes in 22q11.2 Deletion Syndrome. *Am J Hum*
1000 *Genet* 97, 869–877. <https://doi.org/10.1016/j.ajhg.2015.10.013>.

- 1001
- 1002 8. Parenti, I., Rabaneda, L.G., Schoen, H., and Novarino, G. (2020). Neurodevelopmental
1003 Disorders: From Genetics to Functional Pathways. *Trends Neurosci* *43*, 608–621.
1004 <https://doi.org/10.1016/j.tins.2020.05.004>.
- 1005
- 1006 9. Jacquemont, S., Huguet, G., Klein, M., Chawner, S.J.R.A., Donald, K.A., van den Bree,
1007 M.B.M., Sebat, J., Ledbetter, D.H., Constantino, J.N., Earl, R.K., et al. (2022). Genes To
1008 Mental Health (G2MH): A Framework to Map the Combined Effects of Rare and
1009 Common Variants on Dimensions of Cognition and Psychopathology. *Am J Psychiatry*
1010 *179*, 189–203. <https://doi.org/10.1176/appi.ajp.2021.21040432>.
- 1011
- 1012 10. Oetjens, M.T., Kelly, M.A., Sturm, A.C., Martin, C.L., and Ledbetter, D.H. (2019).
1013 Quantifying the polygenic contribution to variable expressivity in eleven rare genetic
1014 disorders. *Nat Commun* *10*, 4897. <https://doi.org/10.1038/s41467-019-12869-0>.
- 1015
- 1016 11. Hudac, C.M., Bove, J., Barber, S., Duyzend, M., Wallace, A., Martin, C.L., Ledbetter,
1017 D.H., Hanson, E., Goin-Kochel, R.P., Green-Snyder, L., et al. (2020). Evaluating
1018 heterogeneity in ASD symptomatology, cognitive ability, and adaptive functioning among
1019 16p11.2 CNV carriers. *Autism Res* *13*, 1300–1310. <https://doi.org/10.1002/aur.2332>.
- 1020
- 1021 12. Cleynen, I., Engchuan, W., Hestand, M.S., Heung, T., Holleman, A.M., Johnston, H.R.,
1022 Monfeuga, T., McDonald-McGinn, D.M., Gur, R.E., Morrow, B.E., et al. (2021). Genetic
1023 contributors to risk of schizophrenia in the presence of a 22q11.2 deletion. *Mol Psychiatry*
1024 *26*, 4496–4510. <https://doi.org/10.1038/s41380-020-0654-3>.
- 1025
- 1026 13. Tansey, K.E., Rees, E., Linden, D.E., Ripke, S., Chambert, K.D., Moran, J.L., McCarroll,
1027 S.A., Holmans, P., Kirov, G., Walters, J., et al. (2016). Common alleles contribute to
1028 schizophrenia in CNV carriers. *Mol Psychiatry* *21*, 1085–1089.
1029 <https://doi.org/10.1038/mp.2015.143>.
- 1030

- 1031 14. Davies, R.W., Fiksinski, A.M., Breetvelt, E.J., Williams, N.M., Hooper, S.R., Monfeuga,
1032 T., Bassett, A.S., Owen, M.J., Gur, R.E., Morrow, B.E., et al. (2020). Using common
1033 genetic variation to examine phenotypic expression and risk prediction in 22q11.2 deletion
1034 syndrome. *Nat Med* 26, 1912–1918. <https://doi.org/10.1038/s41591-020-1103-1>.
1035
- 1036 15. Alver, M., Mancini, V., Läll, K., Schneider, M., Romano, L., Estonian Biobank Research
1037 Team, Mägi, R., Dermitzakis, E.T., Eliez, S., and Reymond, A. (2022). Contribution of
1038 schizophrenia polygenic burden to longitudinal phenotypic variance in 22q11.2 deletion
1039 syndrome. *Mol Psychiatry* 27, 4191–4200. <https://doi.org/10.1038/s41380-022-01674-9>.
1040
- 1041 16. Bergen, S.E., Ploner, A., Howrigan, D., CNV Analysis Group and the Schizophrenia
1042 Working Group of the Psychiatric Genomics Consortium, O’Donovan, M.C., Smoller,
1043 J.W., Sullivan, P.F., Sebat, J., Neale, B., and Kendler, K.S. (2019). Joint Contributions of
1044 Rare Copy Number Variants and Common SNPs to Risk for Schizophrenia. *Am J*
1045 *Psychiatry* 176, 29–35. <https://doi.org/10.1176/appi.ajp.2018.17040467>.
1046
- 1047 17. Alver, M., Mancini, V., Läll, K., Schneider, M., Romano, L., Estonian Biobank Research
1048 Team, Mägi, R., Dermitzakis, E.T., Eliez, S., and Reymond, A. (2022). Contribution of
1049 schizophrenia polygenic burden to longitudinal phenotypic variance in 22q11.2 deletion
1050 syndrome. *Mol Psychiatry* 27, 4191–4200. <https://doi.org/10.1038/s41380-022-01674-9>.
1051
- 1052 18. Banerjee, D., and Girirajan, S. (2023). Pathogenic Variants and Ascertainment:
1053 Neuropsychiatric Disease Risk in a Health System Cohort. *American Journal of*
1054 *Psychiatry* 180, 11–13. <https://doi.org/10.1176/appi.ajp.20220934>.
1055
- 1056 19. Shimelis, H., Oetjens, M.T., Walsh, L.K., Wain, K.E., Znidarsic, M., Myers, S.M.,
1057 Finucane, B.M., Ledbetter, D.H., and Martin, C.L. (2023). Prevalence and Penetrance of
1058 Rare Pathogenic Variants in Neurodevelopmental Psychiatric Genes in a Health Care
1059 System Population. *Am J Psychiatry* 180, 65–72.
1060 <https://doi.org/10.1176/APPI.AJP.22010062>.
1061

- 1062 20. Crawford, K., Bracher-Smith, M., Owen, D., Kendall, K.M., Rees, E., Pardiñas, A.F.,
1063 Einon, M., Escott-Price, V., Walters, J.T.R., O'Donovan, M.C., et al. (2019). Medical
1064 consequences of pathogenic CNVs in adults: Analysis of the UK Biobank. *J Med Genet*
1065 *56*, 131–138. <https://doi.org/10.1136/jmedgenet-2018-105477>.
1066
- 1067 21. Auwerx, C., Lepamets, M., Sadler, M.C., Patxot, M., Stojanov, M., Baud, D., Mägi, R.,
1068 Estonian Biobank Research Team, Porcu, E., Reymond, A., et al. (2022). The individual
1069 and global impact of copy-number variants on complex human traits. *Am J Hum Genet*
1070 *109*, 647–668. <https://doi.org/10.1016/j.ajhg.2022.02.010>.
1071
- 1072 22. Auwerx, C., Jõeloo, M., Sadler, M.C., Tesio, N., Ojavee, S., Clark, C.J., Mägi, R.,
1073 Estonian Biobank Research Team, Reymond, A., and Kutalik, Z. (2024). Rare copy-
1074 number variants as modulators of common disease susceptibility. *Genome Med* *16*, 5.
1075 <https://doi.org/10.1186/s13073-023-01265-5>.
1076
- 1077 23. Weiss, L.A., Shen, Y., Korn, J.M., Arking, D.E., Miller, D.T., Fossdal, R., Saemundsen,
1078 E., Stefansson, H., Ferreira, M.A.R., Green, T., et al. (2008). Association between
1079 Microdeletion and Microduplication at 16p11.2 and Autism. *New England Journal of*
1080 *Medicine* *358*, 667–675. <https://doi.org/10.1056/NEJMoa075974>.
1081
- 1082 24. Auwerx, C., Moix, S., Kutalik, Z., and Reymond, A. (2024). Disentangling mechanisms
1083 behind the pleiotropic effects of proximal 16p11.2 BP4-5 CNVs. medRxiv,
1084 2024.03.20.24304613. <https://doi.org/10.1101/2024.03.20.24304613>.
1085
- 1086 25. Walters, R.G., Jacquemont, S., Valsesia, A., de Smith, A.J., Martinet, D., Andersson, J.,
1087 Falchi, M., Chen, F., Andrieux, J., Lobbens, S., et al. (2010). A new highly penetrant form
1088 of obesity due to deletions on chromosome 16p11.2. *Nature* *463*, 671–675.
1089 <https://doi.org/10.1038/nature08727>.
1090

- 1091 26. Jensen, M., and Girirajan, S. (2019). An interaction-based model for neuropsychiatric
1092 features of copy-number variants. *PLoS Genet* 15, e1007879.
1093 <https://doi.org/10.1371/journal.pgen.1007879>.
1094
- 1095 27. Girirajan, S., Rosenfeld, J.A., Cooper, G.M., Antonacci, F., Siswara, P., Itsara, A., Vives,
1096 L., Walsh, T., McCarthy, S.E., Baker, C., et al. (2010). A recurrent 16p12.1 microdeletion
1097 supports a two-hit model for severe developmental delay. *Nat Genet* 42, 203–209.
1098 <https://doi.org/10.1038/ng.534>.
1099
- 1100 28. Pizzo, L., Jensen, M., Polyak, A., Rosenfeld, J.A., Mannik, K., Krishnan, A., McCready,
1101 E., Pichon, O., Le Caignec, C., Van Dijck, A., et al. (2019). Rare variants in the genetic
1102 background modulate cognitive and developmental phenotypes in individuals carrying
1103 disease-associated variants. *Genet Med* 21, 816–825. [https://doi.org/10.1038/s41436-018-](https://doi.org/10.1038/s41436-018-0266-3)
1104 [0266-3](https://doi.org/10.1038/s41436-018-0266-3).
1105
- 1106 29. Girirajan, S., Pizzo, L., Moeschler, J., and Rosenfeld, J. (1993). 16p12.2 Recurrent
1107 Deletion. In *GeneReviews®* (University of Washington, Seattle).
1108
- 1109 30. Rees, E., Walters, J.T.R., Chambert, K.D., O’Dushlaine, C., Szatkiewicz, J., Richards,
1110 A.L., Georgieva, L., Mahoney-Davies, G., Legge, S.E., Moran, J.L., et al. (2014). CNV
1111 analysis in a large schizophrenia sample implicates deletions at 16p12.1 and SLC1A1 and
1112 duplications at 1p36.33 and CGNL1. *Hum Mol Genet* 23, 1669–1676.
1113 <https://doi.org/10.1093/hmg/ddt540>.
1114
- 1115 31. Stefansson, H., Meyer-Lindenberg, A., Steinberg, S., Magnusdottir, B., Morgen, K.,
1116 Arnarsdottir, S., Bjornsdottir, G., Walters, G.B., Jonsdottir, G.A., Doyle, O.M., et al.
1117 (2014). CNVs conferring risk of autism or schizophrenia affect cognition in controls.
1118 *Nature* 505, 361–366. <https://doi.org/10.1038/nature12818>.
1119
- 1120 32. Rees, E., Kendall, K., Pardiñas, A.F., Legge, S.E., Pocklington, A., Escott-Price, V.,
1121 MacCabe, J.H., Collier, D.A., Holmans, P., O’Donovan, M.C., et al. (2016). Analysis of

- 1122 intellectual disability copy number variants for association with schizophrenia. *JAMA*
1123 *Psychiatry* 73, 963–969. <https://doi.org/10.1001/jamapsychiatry.2016.1831>.
- 1124
- 1125 33. Girirajan, S., Rosenfeld, J.A., Coe, B.P., Parikh, S., Friedman, N., Goldstein, A., Filipink,
1126 R.A., McConnell, J.S., Angle, B., Meschino, W.S., et al. (2012). Phenotypic heterogeneity
1127 of genomic disorders and rare copy-number variants. *N Engl J Med* 367, 1321–1331.
1128 <https://doi.org/10.1056/NEJMoa1200395>.
- 1129
- 1130 34. Hansen, J.A. (2019). Development and Psychometric Evaluation of the Hansen Research
1131 Services Matrix Adaptive Test: A Measure of Nonverbal IQ. *J Autism Dev Disord* 49,
1132 2721–2732. <https://doi.org/10.1007/s10803-016-2932-0>.
- 1133
- 1134 35. Constantino, J.N., Davis, S.A., Todd, R.D., Schindler, M.K., Gross, M.M., Brophy, S.L.,
1135 Metzger, L.M., Shoushtari, C.S., Splinter, R., and Reich, W. (2003). Validation of a brief
1136 quantitative measure of autistic traits: comparison of the social responsiveness scale with
1137 the autism diagnostic interview-revised. *J Autism Dev Disord* 33, 427–433.
1138 <https://doi.org/10.1023/a:1025014929212>.
- 1139
- 1140 36. Zubler, J., and Whitaker, T. (2022). CDC’s Revised Developmental Milestone Checklists.
1141 *Am Fam Physician* 106, 370–371.
- 1142
- 1143 37. Karczewski, K.J., Francioli, L.C., Tiao, G., Cummings, B.B., Alföldi, J., Wang, Q.,
1144 Collins, R.L., Laricchia, K.M., Ganna, A., Birnbaum, D.P., et al. (2020). The mutational
1145 constraint spectrum quantified from variation in 141,456 humans. *Nature* 581, 434–443.
1146 <https://doi.org/10.1038/s41586-020-2308-7>.
- 1147
- 1148 38. Lee, J.J., Wedow, R., Okbay, A., Kong, E., Maghzian, O., Zacher, M., Nguyen-Viet, T.A.,
1149 Bowers, P., Sidorenko, J., Karlsson Linnér, R., et al. (2018). Gene discovery and
1150 polygenic prediction from a genome-wide association study of educational attainment in
1151 1.1 million individuals. *Nat Genet* 50, 1112–1121. [https://doi.org/10.1038/s41588-018-](https://doi.org/10.1038/s41588-018-0147-3)
1152 0147-3.

- 1153
- 1154 39. Savage, J.E., Jansen, P.R., Stringer, S., Watanabe, K., Bryois, J., De Leeuw, C.A., Nagel,
1155 M., Awasthi, S., Barr, P.B., Coleman, J.R.I., et al. (2018). Genome-wide association meta-
1156 analysis in 269,867 individuals identifies new genetic and functional links to intelligence.
1157 *Nat Genet* 50, 912–919. <https://doi.org/10.1038/s41588-018-0152-6>.
- 1158
- 1159 40. Ripke, S., Neale, B.M., Corvin, A., Walters, J.T.R., Farh, K.H., Holmans, P.A., Lee, P.,
1160 Bulik-Sullivan, B., Collier, D.A., Huang, H., et al. (2014). Biological insights from 108
1161 schizophrenia-associated genetic loci. *Nature* 511, 421–427.
1162 <https://doi.org/10.1038/nature13595>.
- 1163
- 1164 41. Grove, J., Ripke, S., Als, T.D., Mattheisen, M., Walters, R.K., Won, H., Pallesen, J.,
1165 Agerbo, E., Andreassen, O.A., Anney, R., et al. (2019). Identification of common genetic
1166 risk variants for autism spectrum disorder. *Nat Genet* 51, 431–444.
1167 <https://doi.org/10.1038/s41588-019-0344-8>.
- 1168
- 1169 42. Landrum, M.J., Lee, J.M., Benson, M., Brown, G.R., Chao, C., Chitipiralla, S., Gu, B.,
1170 Hart, J., Hoffman, D., Jang, W., et al. (2018). ClinVar: Improving access to variant
1171 interpretations and supporting evidence. *Nucleic Acids Res* 46, D1062–D1067.
1172 <https://doi.org/10.1093/nar/gkx1153>.
- 1173
- 1174 43. Amberger, J.S., Bocchini, C.A., Scott, A.F., and Hamosh, A. (2019). OMIM.org:
1175 Leveraging knowledge across phenotype-gene relationships. *Nucleic Acids Res* 47,
1176 D1038–D1043. <https://doi.org/10.1093/nar/gky1151>.
- 1177
- 1178 44. Gonzalez-Mantilla, A.J., Moreno-De-Luca, A., Ledbetter, D.H., and Martin, C.L. (2016).
1179 A cross-disorder method to identify novel candidate genes for developmental brain
1180 disorders. *JAMA Psychiatry* 73, 275–283.
1181 <https://doi.org/10.1001/jamapsychiatry.2015.2692>.
- 1182

- 1183 45. Abrahams, B.S., Arking, D.E., Campbell, D.B., Mefford, H.C., Morrow, E.M., Weiss,
1184 L.A., Menashe, I., Wadkins, T., Banerjee-Basu, S., and Packer, A. (2013). SFARI Gene
1185 2.0: a community-driven knowledgebase for the autism spectrum disorders (ASDs). *Mol*
1186 *Autism* 4, 36. <https://doi.org/10.1186/2040-2392-4-36>.
1187
- 1188 46. Jones, W.D., Dafou, D., McEntagart, M., Woollard, W.J., Elmslie, F. V, Holder-
1189 Espinasse, M., Irving, M., Saggat, A.K., Smithson, S., Trembath, R.C., et al. (2012). De
1190 novo mutations in MLL cause Wiedemann-Steiner syndrome. *Am J Hum Genet* 91, 358–
1191 364. <https://doi.org/10.1016/j.ajhg.2012.06.008>.
1192
- 1193 47. Tang, H., Kirkness, E.F., Lippert, C., Biggs, W.H., Fabani, M., Guzman, E.,
1194 Ramakrishnan, S., Lavrenko, V., Kakaradov, B., Hou, C., et al. (2017). Profiling of Short-
1195 Tandem-Repeat Disease Alleles in 12,632 Human Whole Genomes. *Am J Hum Genet*
1196 101, 700–715. <https://doi.org/10.1016/j.ajhg.2017.09.013>.
1197
- 1198 48. Werling, D.M., Brand, H., An, J.-Y., Stone, M.R., Zhu, L., Glessner, J.T., Collins, R.L.,
1199 Dong, S., Layer, R.M., Markenscoff-Papadimitriou, E., et al. (2018). An analytical
1200 framework for whole-genome sequence association studies and its implications for autism
1201 spectrum disorder. *Nat Genet* 50, 727–736. <https://doi.org/10.1038/s41588-018-0107-y>.
1202
- 1203 49. Greene, C.S., Krishnan, A., Wong, A.K., Ricciotti, E., Zelaya, R.A., Himmelstein, D.S.,
1204 Zhang, R., Hartmann, B.M., Zaslavsky, E., Sealfon, S.C., et al. (2015). Understanding
1205 multicellular function and disease with human tissue-specific networks. *Nat Genet* 47,
1206 569–576. <https://doi.org/10.1038/ng.3259>.
1207
- 1208 50. Krishnan, A., Zhang, R., Yao, V., Theesfeld, C.L., Wong, A.K., Tadych, A., Volfovsky,
1209 N., Packer, A., Lash, A., and Troyanskaya, O.G. (2016). Genome-wide prediction and
1210 functional characterization of the genetic basis of autism spectrum disorder. *Nat Neurosci*
1211 19, 1454–1462. <https://doi.org/10.1038/nn.4353>.
1212

- 1213 51. Hou, J., van Leeuwen, J., Andrews, B.J., and Boone, C. (2018). Genetic Network
1214 Complexity Shapes Background-Dependent Phenotypic Expression. *Trends Genet* 34,
1215 578–586. <https://doi.org/10.1016/j.tig.2018.05.006>.
1216
- 1217 52. Miller, J.A., Ding, S.-L., Sunkin, S.M., Smith, K.A., Ng, L., Szafer, A., Ebbert, A., Riley,
1218 Z.L., Royall, J.J., Aiona, K., et al. (2014). Transcriptional landscape of the prenatal human
1219 brain. *Nature* 508, 199–206. <https://doi.org/10.1038/nature13185>.
1220
- 1221 53. Bakken, T.E., van Velthoven, C.T., Menon, V., Hodge, R.D., Yao, Z., Nguyen, T.N.,
1222 Graybuck, L.T., Horwitz, G.D., Bertagnolli, D., Goldy, J., et al. (2021). Single-cell and
1223 single-nucleus RNA-seq uncovers shared and distinct axes of variation in dorsal LGN
1224 neurons in mice, non-human primates, and humans. *Elife* 10.
1225 <https://doi.org/10.7554/eLife.64875>.
1226
- 1227 54. Jensen, M., Tyryshkina, A., Pizzo, L., Smolen, C., Das, M., Huber, E., Krishnan, A., and
1228 Girirajan, S. (2021). Combinatorial patterns of gene expression changes contribute to
1229 variable expressivity of the developmental delay-associated 16p12.1 deletion. *Genome*
1230 *Med* 13, 163. <https://doi.org/10.1186/s13073-021-00982-z>.
1231
- 1232 55. De Vries, B.B.A., White, S.M., Knight, S.J.L., Regan, R., Homfray, T., Young, I.D.,
1233 Super, M., McKeown, C., Splitt, M., Quarrell, O.W.J., et al. (2001). Clinical studies on
1234 submicroscopic subtelomeric rearrangements: A checklist. *J Med Genet* 38, 145–150.
1235 <https://doi.org/10.1136/jmg.38.3.145>.
1236
- 1237 56. Flannick, J., Beer, N.L., Bick, A.G., Agarwala, V., Molnes, J., Gupta, N., Burt, N.P.,
1238 Florez, J.C., Meigs, J.B., Taylor, H., et al. (2013). Assessing the phenotypic effects in the
1239 general population of rare variants in genes for a dominant Mendelian form of diabetes.
1240 *Nat Genet* 45, 1380–1385. <https://doi.org/10.1038/ng.2794>.
1241
- 1242 57. Gunther, D.F., Eugster, E., Zagar, A.J., Bryant, C.G., Davenport, M.L., and Quigley, C.A.
1243 (2004). Ascertainment bias in Turner syndrome: new insights from girls who were

- 1244 diagnosed incidentally in prenatal life. *Pediatrics* *114*, 640–644.
1245 <https://doi.org/10.1542/peds.2003-1122-L>.
1246
- 1247 58. Feliciano, P., Daniels, A.M., Green Snyder, L.A., Beaumont, A., Camba, A., Esler, A.,
1248 Gulsrud, A.G., Mason, A., Gutierrez, A., Nicholson, A., et al. (2018). SPARK: A US
1249 Cohort of 50,000 Families to Accelerate Autism Research. *Neuron* *97*, 488–493.
1250 <https://doi.org/10.1016/j.neuron.2018.01.015>.
1251
- 1252 59. Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L.T., Sharp, K., Motyer, A.,
1253 Vukcevic, D., Delaneau, O., O’Connell, J., et al. (2018). The UK Biobank resource with
1254 deep phenotyping and genomic data. *Nature* *562*, 203–209.
1255 <https://doi.org/10.1038/S41586-018-0579-Z>.
1256
- 1257 60. Carey, D.J., Fetterolf, S.N., Davis, F.D., Faucett, W.A., Kirchner, H.L., Mirshahi, U.,
1258 Murray, M.F., Smelser, D.T., Gerhard, G.S., and Ledbetter, D.H. (2016). The Geisinger
1259 MyCode community health initiative: an electronic health record-linked biobank for
1260 precision medicine research. *Genet Med* *18*, 906–913.
1261 <https://doi.org/10.1038/gim.2015.187>.
1262
- 1263 61. Fry, A., Littlejohns, T.J., Sudlow, C., Doherty, N., Adamska, L., Sprosen, T., Collins, R.,
1264 and Allen, N.E. (2017). Comparison of Sociodemographic and Health-Related
1265 Characteristics of UK Biobank Participants With Those of the General Population. *Am J*
1266 *Epidemiol* *186*, 1026–1034. <https://doi.org/10.1093/aje/kwx246>.
1267
- 1268 62. Männik, K., Mägi, R., Macé, A., Cole, B., Guyatt, A.L., Shihab, H.A., Maillard, A.M.,
1269 Alavere, H., Kolk, A., Reigo, A., et al. (2015). Copy number variations and cognitive
1270 phenotypes in unselected populations. *JAMA* *313*, 2044–2054.
1271 <https://doi.org/10.1001/jama.2015.4845>.
1272
- 1273 63. Spiro, J.E., and Chung, W.K. (2012). Simons Variation in Individuals Project (Simons
1274 VIP): A Genetics-First Approach to Studying Autism Spectrum and Related

- 1275 Neurodevelopmental Disorders. *Neuron* 73, 1063–1067.
1276 <https://doi.org/10.1016/j.neuron.2012.02.014>.
1277
- 1278 64. Fischbach, G.D., and Lord, C. (2010). The Simons Simplex Collection: A Resource for
1279 Identification of Autism Genetic Risk Factors. *Neuron* 68, 192–195.
1280 <https://doi.org/10.1016/j.neuron.2010.10.006>.
1281
- 1282 65. Pucilowska, J., Vithayathil, J., Tavares, E.J., Kelly, C., Karlo, J.C., and Landreth, G.E.
1283 (2015). The 16p11.2 deletion mouse model of autism exhibits altered cortical progenitor
1284 proliferation and brain cytoarchitecture linked to the ERK MAPK pathway. *J Neurosci* 35,
1285 3190–3200. <https://doi.org/10.1523/JNEUROSCI.4864-13.2015>.
1286
- 1287 66. Iyer, J., Singh, M.D., Jensen, M., Patel, P., Pizzo, L., Huber, E., Koerselman, H., Weiner,
1288 A.T., Lepanto, P., Vadodaria, K., et al. (2018). Pervasive genetic interactions modulate
1289 neurodevelopmental defects of the autism-associated 16p11.2 deletion in *Drosophila*
1290 *melanogaster*. *Nat Commun* 9, 2548. <https://doi.org/10.1038/s41467-018-04882-6>.
1291
- 1292 67. Veltman, J.A., and Brunner, H.G. (2010). Understanding variable expressivity in
1293 microdeletion syndromes. *Nat Genet* 42, 192–193. <https://doi.org/10.1038/ng0310-192>.
1294
- 1295 68. Smolen, C., Jensen, M., Dyer, L., Pizzo, L., Tyryshkina, A., Banerjee, D., Rohan, L.,
1296 Huber, E., El Khattabi, L., Prontera, P., et al. (2023). Assortative mating and parental
1297 genetic relatedness contribute to the pathogenicity of variably expressive variants. *Am J*
1298 *Hum Genet* 110, 2015–2028. <https://doi.org/10.1016/J.AJHG.2023.10.015>.
1299
- 1300 69. Martin, C.L., Wain, K.E., Oetjens, M.T., Tolwinski, K., Palen, E., Hare-Harris, A.,
1301 Habegger, L., Maxwell, E.K., Reid, J.G., Walsh, L.K., et al. (2020). Identification of
1302 Neuropsychiatric Copy Number Variants in a Health Care System Population. *JAMA*
1303 *Psychiatry* 77, 1276–1285. <https://doi.org/10.1001/jamapsychiatry.2020.2159>.
1304

- 1305 70. Männik, K., Mägi, R., Macé, A., Cole, B., Guyatt, A.L., Shihab, H.A., Maillard, A.M.,
1306 Alavere, H., Kolk, A., Reigo, A., et al. (2015). Copy number variations and cognitive
1307 phenotypes in unselected populations. *JAMA* 313, 2044–2054.
1308 <https://doi.org/10.1001/jama.2015.4845>.
1309
- 1310 71. Antaki, D., Guevara, J., Maihofer, A.X., Klein, M., Gujral, M., Grove, J., Carey, C.E.,
1311 Hong, O., Arranz, M.J., Hervas, A., et al. (2022). A phenotypic spectrum of autism is
1312 attributable to the combined effects of rare variants, polygenic risk and sex. *Nat Genet* 54,
1313 1284–1292. <https://doi.org/10.1038/s41588-022-01064-5>.
1314
- 1315 72. Mitra, I., Lavillaureix, A., Yeh, E., Traglia, M., Tsang, K., Bearden, C.E., Rauen, K.A.,
1316 and Weiss, L.A. (2017). Reverse Pathway Genetic Approach Identifies Epistasis in
1317 Autism Spectrum Disorders. *PLoS Genet* 13, e1006516.
1318 <https://doi.org/10.1371/journal.pgen.1006516>.
1319
- 1320 73. Hivert, V., Sidorenko, J., Rohart, F., Goddard, M.E., Yang, J., Wray, N.R., Yengo, L., and
1321 Visscher, P.M. (2021). Estimation of non-additive genetic variance in human complex
1322 traits from a large sample of unrelated individuals. *Am J Hum Genet* 108, 786–798.
1323 <https://doi.org/10.1016/j.ajhg.2021.02.014>.
1324
- 1325 74. Boyle, E.A., Li, Y.I., and Pritchard, J.K. (2017). An Expanded View of Complex Traits:
1326 From Polygenic to Omnigenic. *Cell* 169, 1177–1186.
1327 <https://doi.org/10.1016/j.cell.2017.05.038>.
1328
- 1329 75. Mitra, I., Huang, B., Mousavi, N., Ma, N., Lamkin, M., Yanicky, R., Shleizer-Burko, S.,
1330 Lohmueller, K.E., and Gymrek, M. (2021). Patterns of de novo tandem repeat mutations
1331 and their role in autism. *Nature* 589, 246–250. [https://doi.org/10.1038/s41586-020-03078-](https://doi.org/10.1038/s41586-020-03078-7)
1332 7.
1333
- 1334 76. Turner, T.N., Coe, B.P., Dickel, D.E., Hoekzema, K., Nelson, B.J., Zody, M.C.,
1335 Kronenberg, Z.N., Hormozdiari, F., Raja, A., Pennacchio, L.A., et al. (2017). Genomic

- 1336 Patterns of De Novo Mutation in Simplex Autism. *Cell* 171, 710–722.
1337 <https://doi.org/10.1016/j.cell.2017.08.047>.
1338
- 1339 77. Kuczmariski, R.J., Ogden, C.L., Guo, S.S., Grummer-Strawn, L.M., Flegal, K.M., Mei, Z.,
1340 Wei, R., Curtin, L.R., Roche, A.F., and Johnson, C.L. (2002). 2000 CDC Growth Charts
1341 for the United States: methods and development. *Vital Health Stat* 11, 1–190.
1342
- 1343 78. Rollins, J.D., Collins, J.S., and Holden, K.R. (2010). United States head circumference
1344 growth reference charts: birth to 21 years. *J Pediatr* 156, 907-913.e2.
1345 <https://doi.org/10.1016/j.jpeds.2010.01.009>.
1346
- 1347 79. Bolger, A.M., Lohse, M., and Usadel, B. (2014). Trimmomatic: A flexible trimmer for
1348 Illumina sequence data. *Bioinformatics* 30, 2114–2120.
1349 <https://doi.org/10.1093/bioinformatics/btu170>.
1350
- 1351 80. Li, H., and Durbin, R. (2010). Fast and accurate long-read alignment with Burrows-
1352 Wheeler transform. *Bioinformatics* 26, 589–595.
1353 <https://doi.org/10.1093/bioinformatics/btp698>.
1354
- 1355 81. Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis,
1356 G., and Durbin, R. (2009). The Sequence Alignment/Map format and SAMtools.
1357 *Bioinformatics* 25, 2078–2079. <https://doi.org/10.1093/bioinformatics/btp352>.
1358
- 1359 82. Poplin, R., Ruano-Rubio, V., DePristo, M.A., Fennell, T.J., Carneiro, M.O., Auwera, G.A.
1360 Van der, Kling, D.E., Gauthier, L.D., Levy-Moonshine, A., Roazen, D., et al. (2017).
1361 Scaling accurate genetic variant discovery to tens of thousands of samples. *bioRxiv*,
1362 201178. <https://doi.org/10.1101/201178>.
1363
- 1364 83. Pedersen, B.S., Layer, R.M., and Quinlan, A.R. (2016). Vcfanno: fast, flexible annotation
1365 of genetic variants. *Genome Biol* 17, 118. <https://doi.org/10.1186/s13059-016-0973-5>.
1366

- 1367 84. Lek, M., Karczewski, K.J., Minikel, E. V., Samocha, K.E., Banks, E., Fennell, T.,
1368 O'Donnell-Luria, A.H., Ware, J.S., Hill, A.J., Cummings, B.B., et al. (2016). Analysis of
1369 protein-coding genetic variation in 60,706 humans. *Nature* 536, 285–291.
1370 <https://doi.org/10.1038/nature19057>.
1371
- 1372 85. Frankish, A., Diekhans, M., Jungreis, I., Lagarde, J., Loveland, J.E., Mudge, J.M., Sisu,
1373 C., Wright, J.C., Armstrong, J., Barnes, I., et al. (2021). GENCODE 2021. *Nucleic Acids*
1374 *Res* 49, D916–D923. <https://doi.org/10.1093/nar/gkaa1087>.
1375
- 1376 86. Wang, K., Li, M., and Hakonarson, H. (2010). ANNOVAR: functional annotation of
1377 genetic variants from high-throughput sequencing data. *Nucleic Acids Res* 38, e164.
1378 <https://doi.org/10.1093/nar/gkq603>.
1379
- 1380 87. Rentzsch, P., Witten, D., Cooper, G.M., Shendure, J., and Kircher, M. (2019). CADD:
1381 predicting the deleteriousness of variants throughout the human genome. *Nucleic Acids*
1382 *Res* 47, D886–D894. <https://doi.org/10.1093/nar/gky1016>.
1383
- 1384 88. Roadmap Epigenomics Consortium, Kundaje, A., Meuleman, W., Ernst, J., Bilenky, M.,
1385 Yen, A., Heravi-Moussavi, A., Kheradpour, P., Zhang, Z., Wang, J., et al. (2015).
1386 Integrative analysis of 111 reference human epigenomes. *Nature* 518, 317–329.
1387 <https://doi.org/10.1038/NATURE14248>.
1388
- 1389 89. Wang, K., Li, M., Hadley, D., Liu, R., Glessner, J., Grant, S.F.A., Hakonarson, H., and
1390 Bucan, M. (2007). PennCNV: An integrated hidden Markov model designed for high-
1391 resolution copy number variation detection in whole-genome SNP genotyping data.
1392 *Genome Res* 17, 1665–1674. <https://doi.org/10.1101/gr.6861907>.
1393
- 1394 90. Abyzov, A., Urban, A.E., Snyder, M., and Gerstein, M. (2011). CNVnator: An approach
1395 to discover, genotype, and characterize typical and atypical CNVs from family and
1396 population genome sequencing. *Genome Res* 21, 974–984.
1397 <https://doi.org/10.1101/gr.114876.110>.

- 1398
- 1399 91. Chen, X., Schulz-Trieglaff, O., Shaw, R., Barnes, B., Schlesinger, F., Källberg, M., Cox,
1400 A.J., Kruglyak, S., and Saunders, C.T. (2016). Manta: Rapid detection of structural
1401 variants and indels for germline and cancer sequencing applications. *Bioinformatics* 32,
1402 1220–1222. <https://doi.org/10.1093/bioinformatics/btv710>.
- 1403
- 1404 92. Layer, R.M., Chiang, C., Quinlan, A.R., and Hall, I.M. (2014). LUMPY: A probabilistic
1405 framework for structural variant discovery. *Genome Biol* 15, R84.
1406 <https://doi.org/10.1186/gb-2014-15-6-r84>.
- 1407
- 1408 93. Rausch, T., Zichner, T., Schlattl, A., Stütz, A.M., Benes, V., and Korbel, J.O. (2012).
1409 DELLY: Structural variant discovery by integrated paired-end and split-read analysis.
1410 *Bioinformatics* 28, 333–339. <https://doi.org/10.1093/bioinformatics/bts378>.
- 1411
- 1412 94. Brandler, W.M., Antaki, D., Gujral, M., Kleiber, M.L., Whitney, J., Maile, M.S., Hong,
1413 O., Chapman, T.R., Tan, S., Tandon, P., et al. (2018). Paternally inherited cis-regulatory
1414 structural variants are associated with autism. *Science* 360, 327–331.
1415 <https://doi.org/10.1126/science.aan2261>.
- 1416
- 1417 95. Collins, R.L., Brand, H., Karczewski, K.J., Zhao, X., Alföldi, J., Francioli, L.C., Khera, A.
1418 V., Lowther, C., Gauthier, L.D., Wang, H., et al. (2020). A structural variation reference
1419 for medical and population genetics. *Nature* 581, 444–451.
1420 <https://doi.org/10.1038/s41586-020-2287-8>.
- 1421
- 1422 96. Coe, B.P., Witherspoon, K., Rosenfeld, J.A., Van Bon, B.W.M., Vulto-Van Silfhout, A.T.,
1423 Bosco, P., Friend, K.L., Baker, C., Buono, S., Vissers, L.E.L.M., et al. (2014). Refining
1424 analyses of copy number variation identifies specific genes associated with developmental
1425 delay. *Nat Genet* 46, 1063–1071. <https://doi.org/10.1038/ng.3092>.
- 1426

- 1427 97. Mousavi, N., Margoliash, J., Pusarla, N., Saini, S., Yanicky, R., and Gymrek, M. (2021).
1428 TRTools: A toolkit for genome-wide analysis of tandem repeats. *Bioinformatics* 37, 731–
1429 733. <https://doi.org/10.1093/bioinformatics/btaa736>.
1430
- 1431 98. Mousavi, N., Shleizer-Burko, S., Yanicky, R., and Gymrek, M. (2019). Profiling the
1432 genome-wide landscape of tandem repeat expansions. *Nucleic Acids Res* 47, e90.
1433 <https://doi.org/10.1093/nar/gkz501>.
1434
- 1435 99. Choi, S.W., Mak, T.S.H., and O'Reilly, P.F. (2020). Tutorial: a guide to performing
1436 polygenic risk score analyses. *Nat Protoc* 15, 2759–2772. <https://doi.org/10.1038/s41596-020-0353-1>.
1437
1438
- 1439 100. Das, S., Forer, L., Schönherr, S., Sidore, C., Locke, A.E., Kwong, A., Vrieze, S.I., Chew,
1440 E.Y., Levy, S., McGue, M., et al. (2016). Next-generation genotype imputation service
1441 and methods. *Nat Genet* 48, 1284–1287. <https://doi.org/10.1038/ng.3656>.
1442
- 1443 101. Pedersen, B.S., and Quinlan, A.R. (2017). Who's Who? Detecting and Resolving Sample
1444 Anomalies in Human DNA Sequencing Studies with Peddy. *Am J Hum Genet* 100, 406–
1445 413. <https://doi.org/10.1016/j.ajhg.2017.01.017>.
1446
- 1447 102. Privé, F., Arbel, J., and Vilhjálmsson, B.J. (2021). LDpred2: better, faster, stronger.
1448 *Bioinformatics* 36, 5424–5431. <https://doi.org/10.1093/bioinformatics/btaa1029>.
1449
- 1450 103. Altshuler, D.M., Gibbs, R.A., Peltonen, L., Schaffner, S.F., Yu, F., Dermitzakis, E.,
1451 Bonnen, P.E., De Bakker, P.I.W., Deloukas, P., Gabriel, S.B., et al. (2010). Integrating
1452 common and rare genetic variation in diverse human populations. *Nature* 467, 52–58.
1453 <https://doi.org/10.1038/nature09298>.
1454
- 1455 104. Wang, J., Lin, Z.J., Liu, L., Xu, H.Q., Shi, Y.W., Yi, Y.H., He, N., and Liao, W.P. (2017).
1456 Epilepsy-associated genes. *Seizure* 44, 11–20.
1457 <https://doi.org/10.1016/j.seizure.2016.11.030>.

- 1458
- 1459 105. Thormann, A., Halachev, M., McLaren, W., Moore, D.J., Svinti, V., Campbell, A., Kerr,
1460 S.M., Tischkowitz, M., Hunt, S.E., Dunlop, M.G., et al. (2019). Flexible and scalable
1461 diagnostic filtering of genomic variants using G2P with Ensembl VEP. *Nat Commun* 10,
1462 2373. <https://doi.org/10.1038/s41467-019-10016-3>.
- 1463
- 1464 106. Wu, Y., Li, X., Liu, J., Luo, X.-J., and Yao, Y.-G. (2020). SZDB2.0: an updated
1465 comprehensive resource for schizophrenia research. *Hum Genet* 139, 1285–1297.
1466 <https://doi.org/10.1007/s00439-020-02171-1>.
- 1467
- 1468 107. Mi, H., Muruganujan, A., Ebert, D., Huang, X., and Thomas, P.D. (2019). PANTHER
1469 version 14: More genomes, a new PANTHER GO-slim and improvements in enrichment
1470 analysis tools. *Nucleic Acids Res* 47, D419–D426. <https://doi.org/10.1093/nar/gky1038>.
- 1471
- 1472 108. Sayols, S. (2023). rrvgo: a Bioconductor package for interpreting lists of Gene Ontology
1473 terms. *MicroPubl Biol* 2023. <https://doi.org/10.17912/micropub.biology.000811>.
- 1474
- 1475 109. Beck, D.B., Bodian, D.L., Shah, V., Mirshahi, U.L., Kim, J., Ding, Y., Magaziner, S.J.,
1476 Strande, N.T., Cantor, A., Haley, J.S., et al. (2023). Estimated Prevalence and Clinical
1477 Manifestations of UBA1 Variants Associated With VEXAS Syndrome in a Clinical
1478 Population. *JAMA* 329, 318–324. <https://doi.org/10.1001/jama.2022.24836>.
- 1479
- 1480 110. Staples, J., Maxwell, E.K., Gosalia, N., Gonzaga-Jauregui, C., Snyder, C., Hawes, A.,
1481 Penn, J., Ulloa, R., Bai, X., Lopez, A.E., et al. (2018). Profiling and Leveraging
1482 Relatedness in a Precision Medicine Cohort of 92,455 Exomes. *Am J Hum Genet* 102,
1483 874–889. <https://doi.org/10.1016/j.ajhg.2018.03.012>.
- 1484
- 1485 111. Backman, J.D., Li, A.H., Marcketta, A., Sun, D., Mbatchou, J., Kessler, M.D., Benner, C.,
1486 Liu, D., Locke, A.E., Balasubramanian, S., et al. (2021). Exome sequencing and analysis
1487 of 454,787 UK Biobank participants. *Nature* 599, 628–634.
1488 <https://doi.org/10.1038/s41586-021-04103-z>.

- 1489
- 1490 112. McLaren, W., Gil, L., Hunt, S.E., Riat, H.S., Ritchie, G.R.S., Thormann, A., Flicek, P.,
1491 and Cunningham, F. (2016). The Ensembl Variant Effect Predictor. *Genome Biol* 17, 122.
1492 <https://doi.org/10.1186/s13059-016-0974-4>.
- 1493
- 1494 113. Liu, X., Li, C., Mou, C., Dong, Y., and Tu, Y. (2020). dbNSFP v4: a comprehensive
1495 database of transcript-specific functional predictions and annotations for human
1496 nonsynonymous and splice-site SNVs. *Genome Med* 12, 103.
1497 <https://doi.org/10.1186/s13073-020-00803-9>.
- 1498
- 1499 114. Denny, J.C., Ritchie, M.D., Basford, M.A., Pulley, J.M., Bastarache, L., Brown-Gentry,
1500 K., Wang, D., Masys, D.R., Roden, D.M., and Crawford, D.C. (2010). PheWAS:
1501 demonstrating the feasibility of a phenome-wide scan to discover gene-disease
1502 associations. *Bioinformatics* 26, 1205–1210.
1503 <https://doi.org/10.1093/bioinformatics/btq126>.
- 1504
- 1505 115. Sanders, S.J., He, X., Willsey, A.J., Ercan-Sencicek, A.G., Samocha, K.E., Cicek, A.E.,
1506 Murtha, M.T., Bal, V.H., Bishop, S.L., Dong, S., et al. (2015). Insights into Autism
1507 Spectrum Disorder Genomic Architecture and Biology from 71 Risk Loci. *Neuron* 87,
1508 1215–1233. <https://doi.org/10.1016/j.neuron.2015.09.016>.
- 1509
- 1510

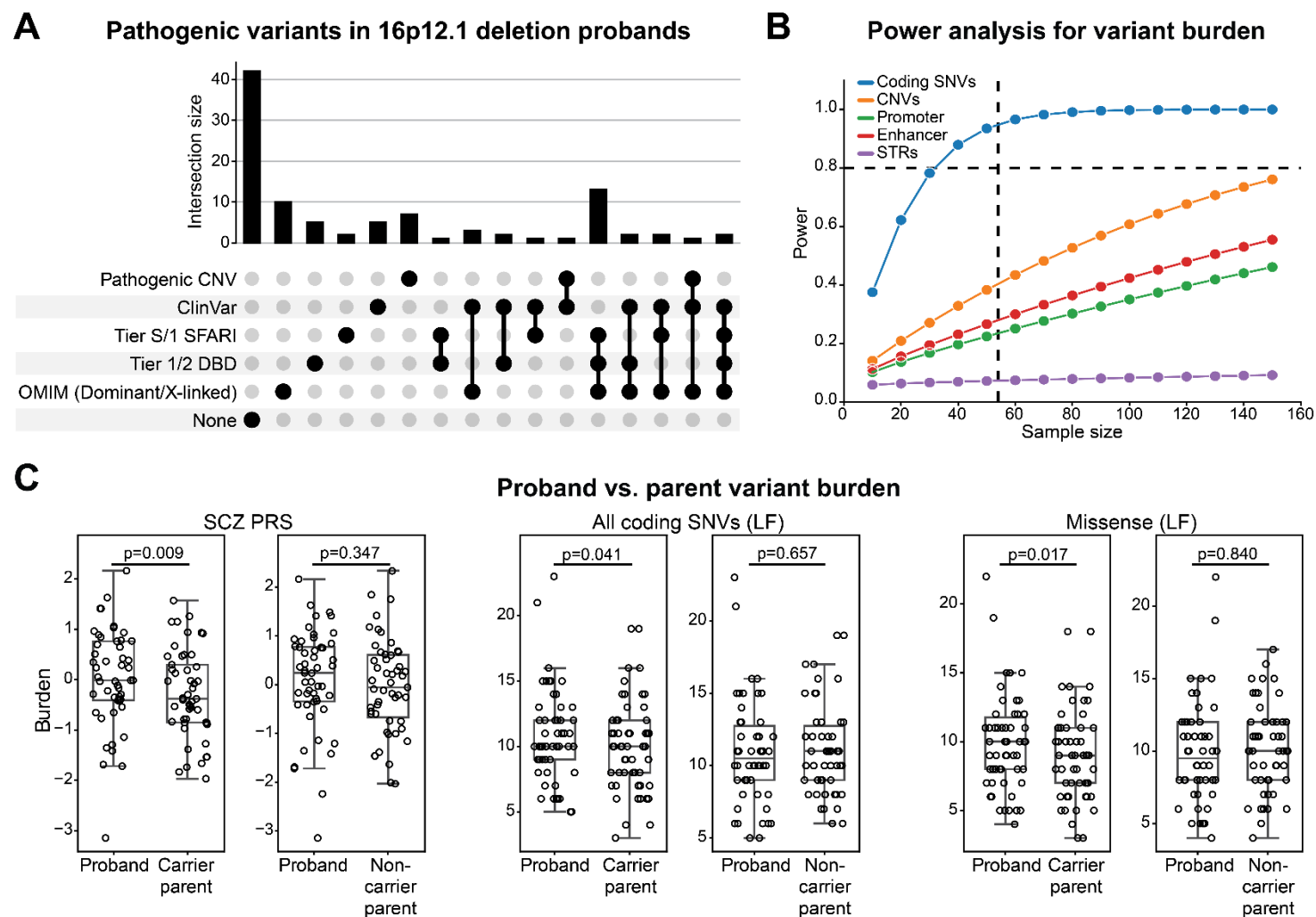


Figure S1. Secondary variant burden comparisons among 16p12.1 deletion family members (related to Figure 3). (A) UpSet plot shows the number of 16p12.1 deletion probands with secondary variants in one or more disease-associated categories, potentially indicative of multiple genetic diagnoses. (B) Power analysis for detecting changes in burden of rare variant classes

among individuals with the 16p12.1 deletion (see Methods). Dashed horizontal line indicates 80% power and dashed vertical line represents the sample size of proband-carrier parent pairs (n=54). (C) Changes in burden of SCZ PRS (left), missense (LF) variants (center), and all coding SNVs (LF) (right) between 16p12.1 deletion probands and their carrier (left, n=49-54) and noncarrier parents (right, n=50-51). P-values from one-tailed (rare variants) or two-tailed (SCZ PRS) t-tests.

1512

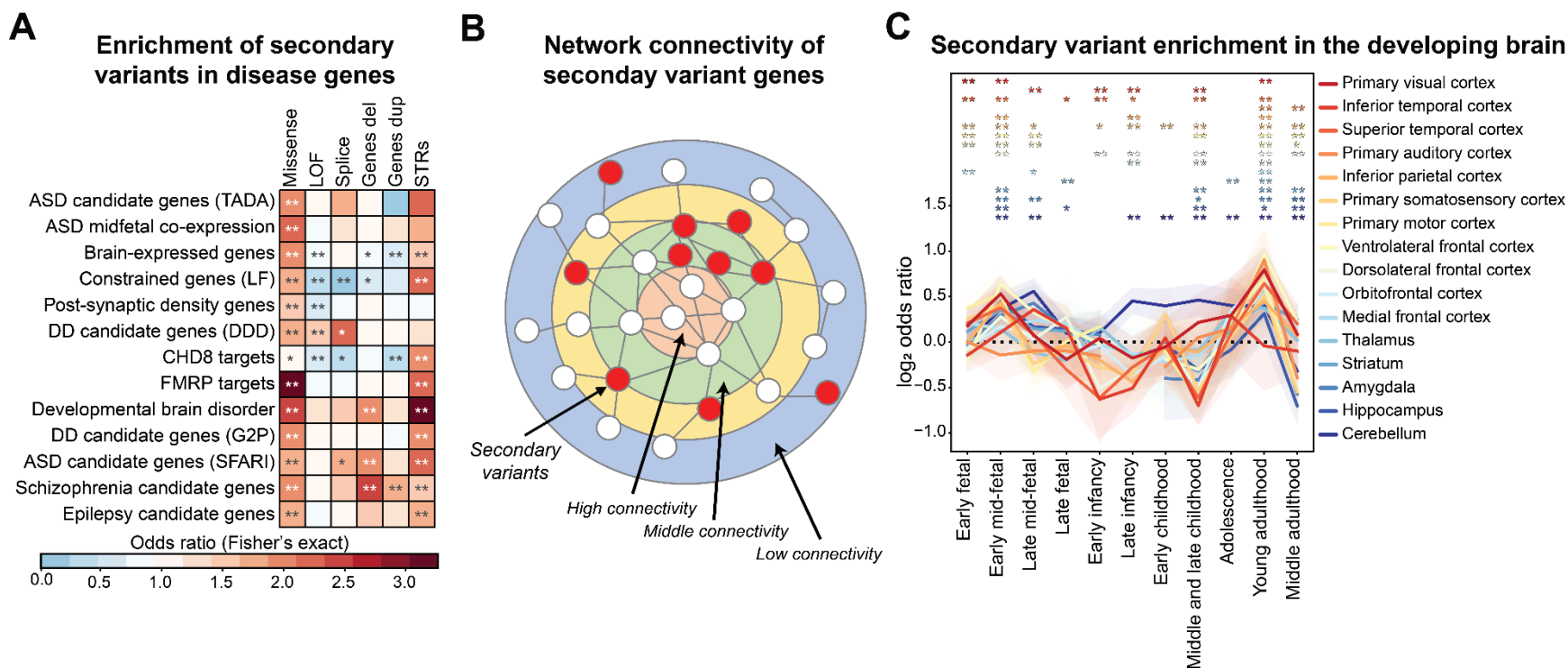


Figure S2. Functional effects of secondary variants observed in 16p12.1 deletion probands (related to Figure 3). (A)

Enrichment of secondary variant classes in 16p12.1 deletion probands for sets of genes involved with neurodevelopmental disease and related functions. Fisher's exact test, $*p \leq 0.05$, $**$ Benjamini-Hochberg $FDR \leq 0.05$. (B) Diagram illustrating the distribution of secondary variants (red nodes) in genes with varying connectivity (colored rings) in a brain-specific interaction network. Highly connected genes (light red ring) are depleted for secondary variants, while genes with intermediate connectivity (light green ring) are enriched for variants. (C) Line plot shows enrichment (log-odds ratios with 95% confidence intervals; y-axis) of secondary variants in 16p12.1 deletion probands among genes preferentially expressed in 16 brain tissues (colored lines) over 11 developmental timepoints (x-axis). Fisher's exact test, $*p \leq 0.05$, $**$ Benjamini-Hochberg $FDR \leq 0.05$.

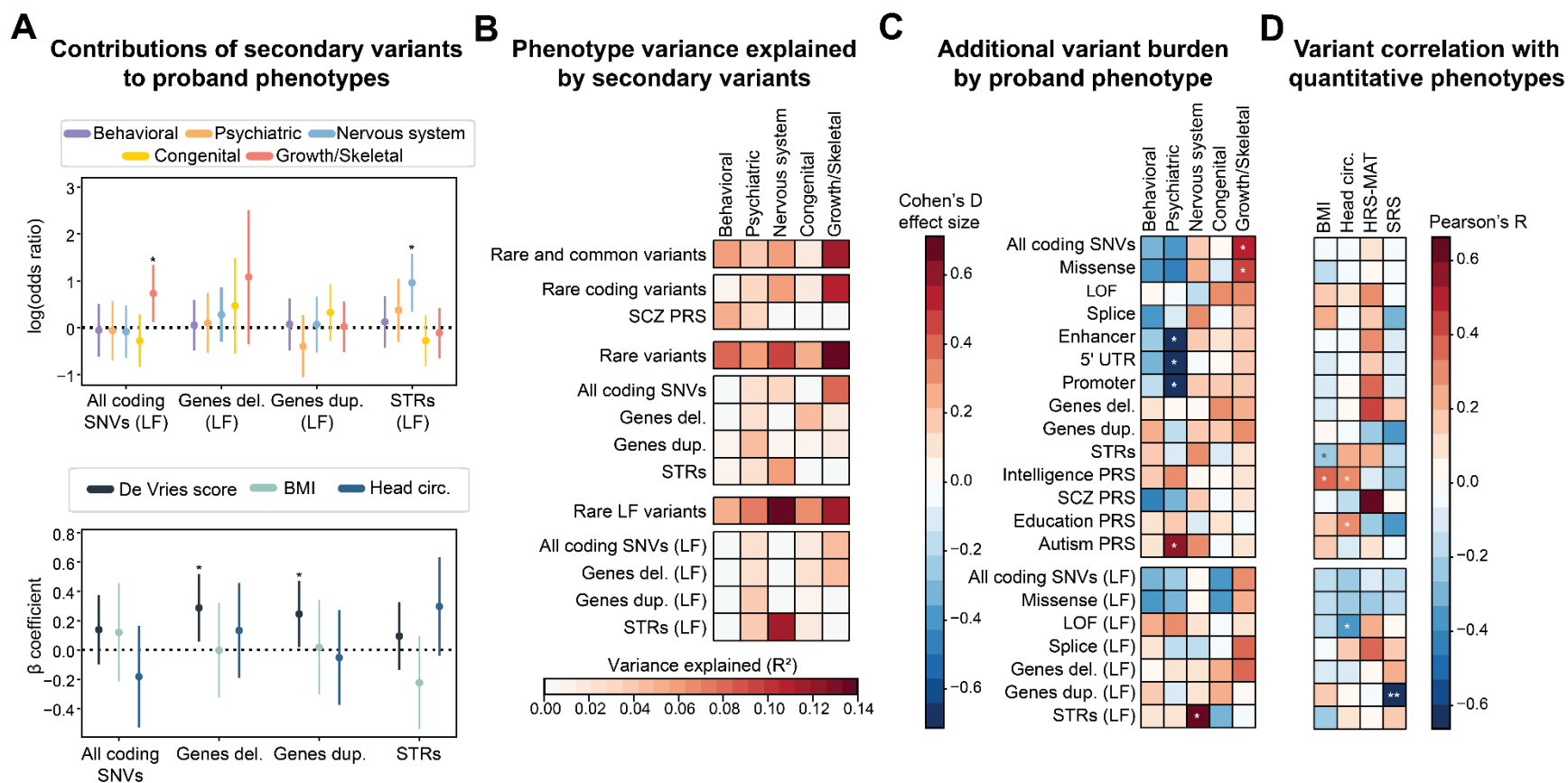


Figure S3. Secondary variant associations with 16p12.1 deletion phenotypic domains (related to Figure 4). (A) (Top) Forest plots show log-scaled odds ratios from logistic regression models for secondary variant burden in constrained genes for probands ($n=47-71$) with higher complexity scores in five phenotypic domains, compared with probands with lower complexity scores for each domain. $*p < 0.05$. (Bottom) β coefficients from linear regression models for quantitative phenotypes in probands ($n=43-76$). $*p < 0.05$. (B) Variance explained by secondary variant burden from logistic regression models (McFadden's pseudo- R^2), both for individual variant classes and joint contributions from combinations of classes. (C) Comparisons of secondary variant burden between 16p12.1 deletion probands ($n=53-84$) with higher and lower complexity scores for each phenotypic domain. Two-tailed t-test, $*p < 0.05$, **Benjamini-Hochberg FDR ≤ 0.05 . (D) Pearson correlations between quantitative phenotypes and secondary variant burden in deletion probands ($n=9-59$). $*p < 0.05$, **Benjamini-Hochberg FDR ≤ 0.05 .

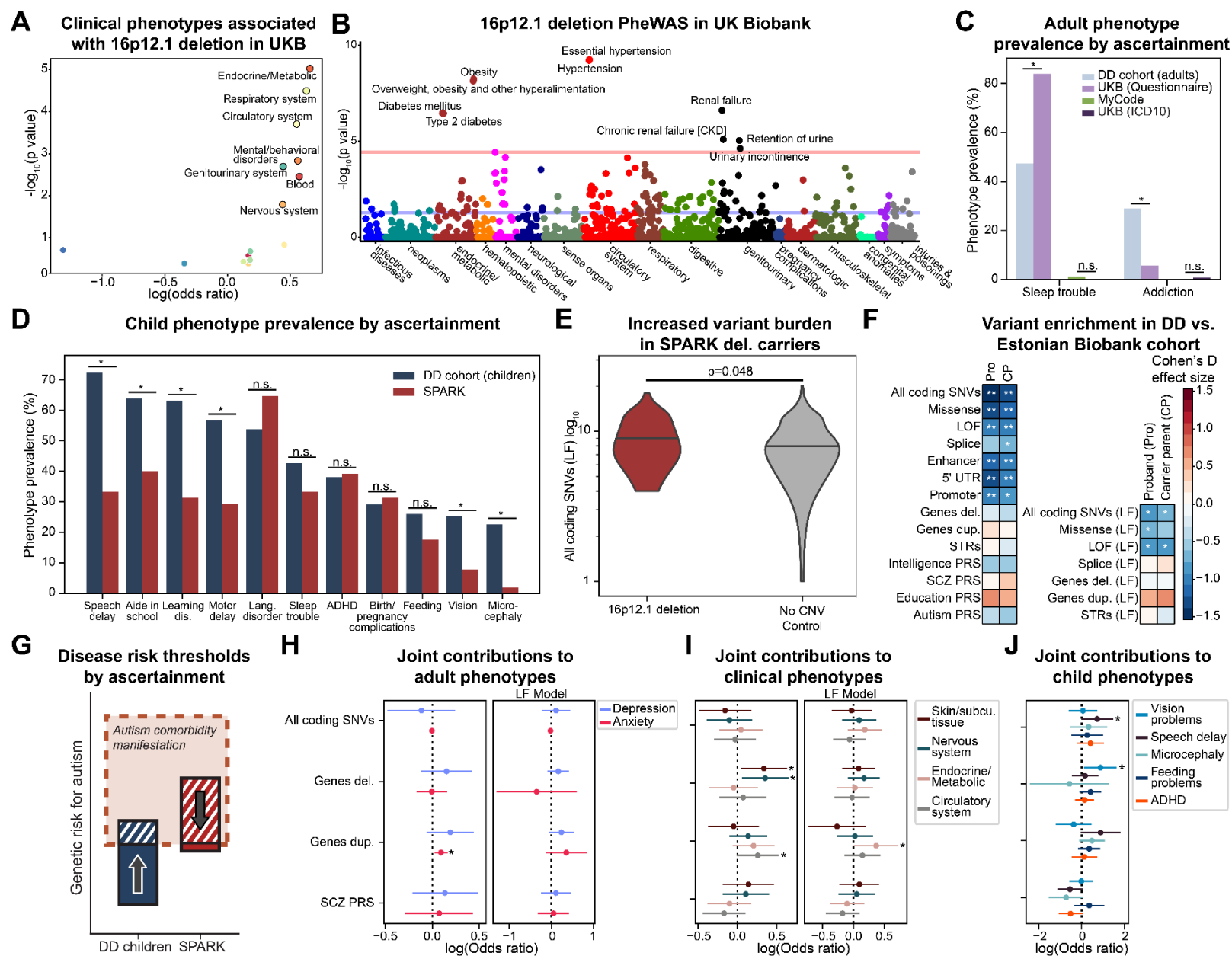


Figure S4. Effects of ascertainment on phenotypes and secondary variant associations in = 16p12.1 deletion carriers

(related to Figure 5). (A) Enrichment of select ICD10 chapters from logistic regression models in UK Biobank (UKB) 16p12.1 deletion carriers compared to controls without large rare CNVs (n=3,488). P-values from logistic regression. Labeled points indicate phenotypes with Benjamini-Hochberg $FDR \leq 0.05$. (B) PheWAS analysis for 16p12.1 deletion carriers in UKB (n=99,363-255,262). Colored circles indicate individual phenotype membership in respective ICD10 chapters. Red line indicates phenome-wide significance (Bonferroni $p=0.05$) and blue line indicates nominal significance ($p=0.05$). (C) Comparison of sleep disturbance and addiction phenotype prevalence in adults from the DD cohort (n=38) and individuals from UK Biobank (questionnaire n=35-249, ICD10 n=217) and MyCode (n=160). * $p \leq 0.05$, Fisher's Exact test. (D) Prevalence of developmental and psychiatric phenotypes in children with 16p12.1 deletion from the DD (n=80-151) and SPARK (n=40-51) cohorts. Phenotypes shown were restricted to those present in >20% of probands in either cohort. * $p \leq 0.05$, Fisher's Exact test. (E) Comparison of SNV (LF) burden in SPARK individuals with 16p12.1 deletion (n=89, left) to age and sex-matched controls without large rare (>500kb) CNVs (n=356, right). P-value from two-tailed t-test. (F) Changes in secondary variant burden between probands ("Pro", n=97-99) and carrier parents ("CP", n=54-57) in the DD cohort with 16p12.1 deletion individuals the Estonian Biobank (n=5-8). Blue indicates a depletion in secondary variant burden for Estonian Biobank deletion carriers. One-tailed t-test, * $p \leq 0.05$, **Benjamini-Hochberg $FDR \leq 0.05$. (G) Schematic outlining the proposed relationship between different genetic risk factors in individuals with 16p12.1 deletion across different ascertainments. In cohorts where a majority of participants have a particular disorder, such as autism in SPARK, established risk factors (such as autism PRS) may not show the expected correlations for comorbid features. However, these correlations would be observed in cohorts with different ascertainments (such as the DD cohort). (H-J) Forest plots show associations of secondary variants in all genes and constrained genes ("LF Model") with select phenotypes from joint logistic models in (H) DD cohort adults, UK Biobank, and MyCode individuals for psychiatric features (n=331); (I) UK Biobank and MyCode individuals for clinical phenotypes from EHR data (n=321); and (J) children from the DD cohort and SPARK (n=98-125). * $p \leq 0.05$. Full results are available in **Table S4G**.

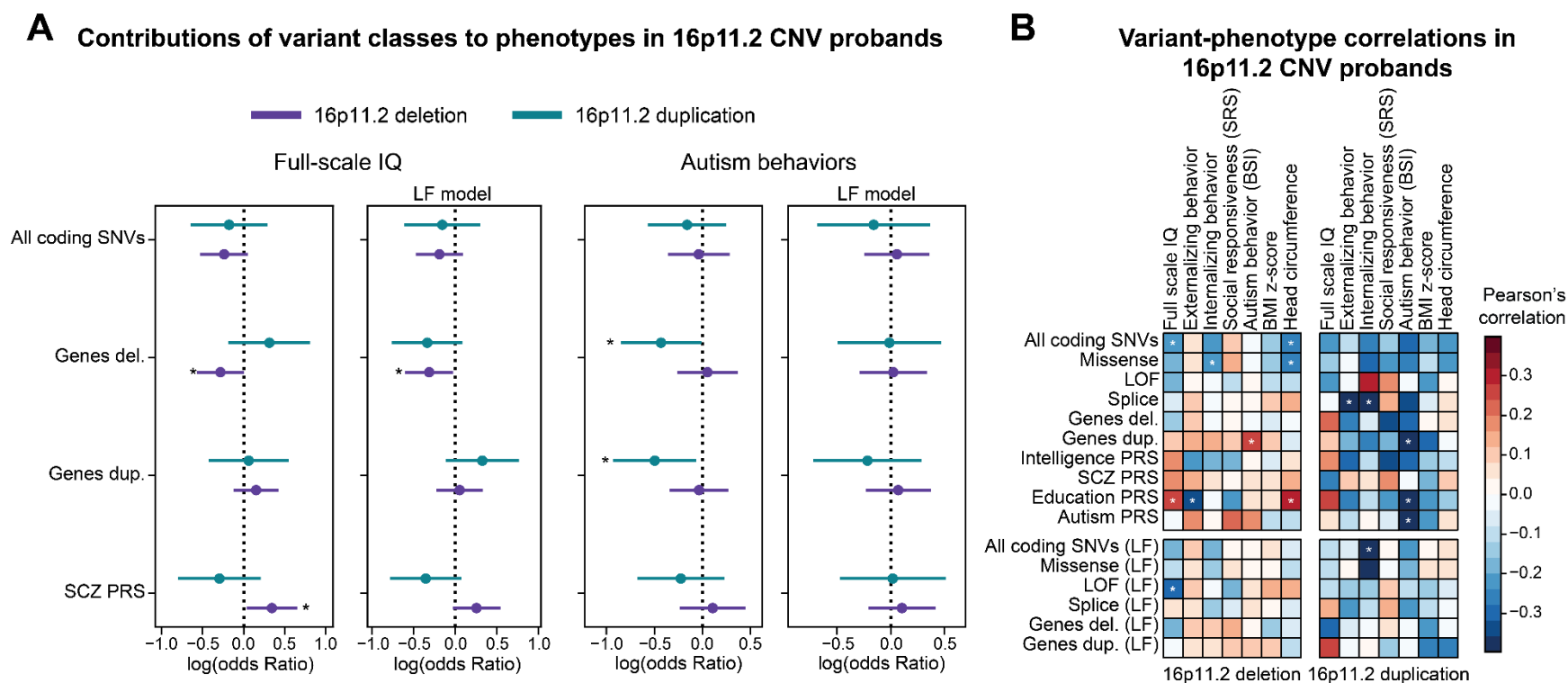
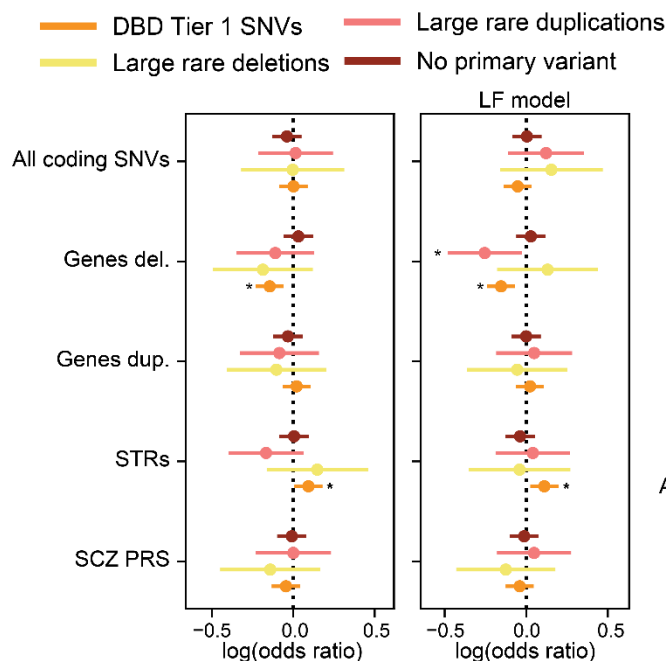


Figure S5. Associations between secondary variants and developmental features of 16p11.2 CNV probands (related to Figure 6). (A) Example forest plots show results from select linear regression models for associations between secondary variant classes and full-scale IQ (left) and autism-related behavior (BSI, right) for probands with the 16p11.2 deletion ($n=52-57$, purple) and duplication ($n=21-25$, teal). “LF model” indicates models where rare variants are selected for genes under evolutionary constraint. $*p \leq 0.05$. Full results are available in **Table S5A**. (B) Pearson’s correlations between secondary variant burden and quantitative phenotypes of probands with 16p11.2 deletion ($n=58-89$) and duplication ($n=25-37$). $*p \leq 0.05$.

A Additional variant effects on full-scale IQ of primary variant probands



B Variant-phenotype correlations across primary variants

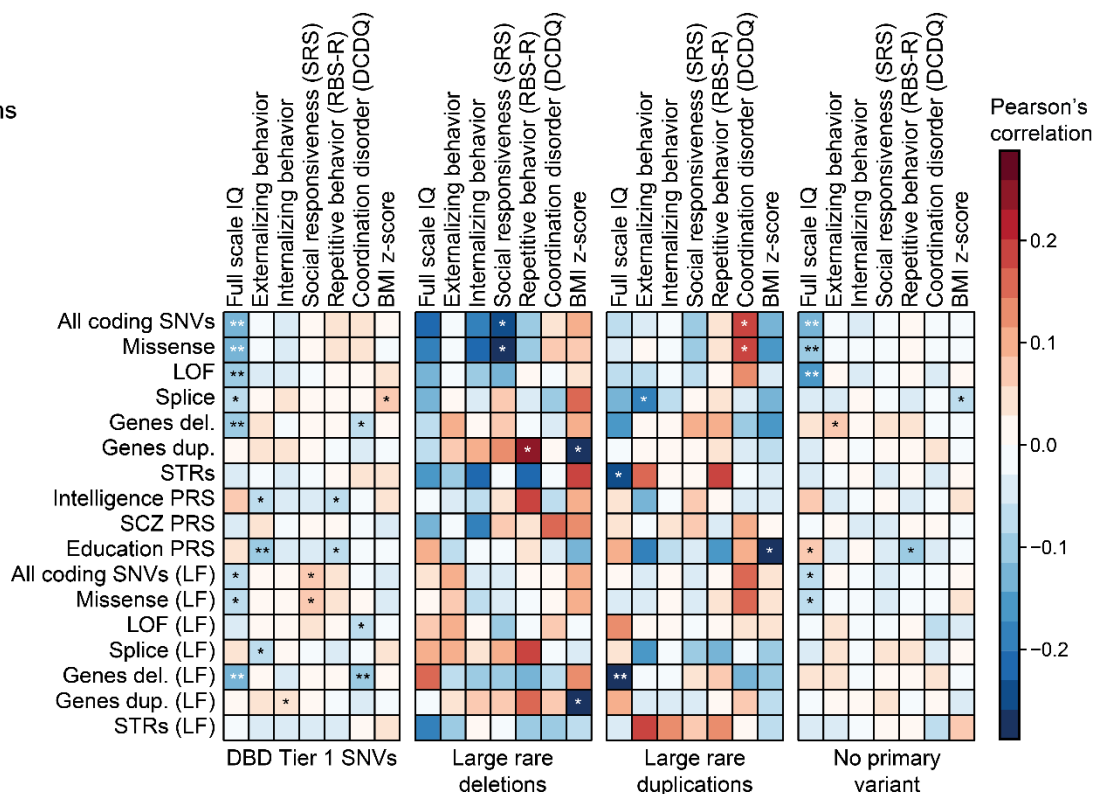


Figure S6. Associations between secondary variants and developmental features of probands with primary variants (related to Figure 6).

(A) Example forest plots show results from select linear regression models for associations between secondary variant classes and full-scale IQ for probands with pathogenic SNVs in candidate neurodevelopmental genes ($n=660$, orange), large, rare deletions ($n=51$, yellow) and duplications ($n=85$, pink), and probands without such variants ($n=632$, red) from the SSC cohort. $*p \leq 0.05$. “LF model” indicates models where rare variants are selected for genes under evolutionary constraint. Full results are available in **Table S5A**. (B) Pearson’s correlations between secondary variant burden and quantitative developmental phenotypes of SSC probands with pathogenic SNVs ($n=736-1,236$), rare deletions ($n=49-78$), rare duplications ($n=102-148$), and probands without such variants ($n=671-1,083$). $*p \leq 0.05$, $**$ Benjamini-Hochberg $FDR \leq 0.05$.

1519 **SUPPLEMENTAL TABLES**

1520 **Table S1. Description of the DD cohort and phenotypic data (Excel file; related to Figure**
1521 **1). Table S1A** lists all 452 individuals in the DD cohort, including family relationships, age,
1522 biological sex, and 16p12.1 deletion status if known. The table also lists secondary variant
1523 burden (rare variant counts and PRS), complexity scores for phenotypic domains (child and
1524 adult), quantitative measures (BMI, head circumference, IQ, and SRS), and age at developmental
1525 milestone achievement for all individuals with available data. **Table S1B** contains the scoring
1526 rubric used to calculate complexity scores for phenotypic domains in children. **Table S1C** lists
1527 minimum age thresholds used for identifying psychiatric features in pediatric family members.
1528 **Table S1D** lists pathogenic variants or deleterious variants in genes associated with disease that
1529 were identified in 16p12.1 deletion probands.

1530
1531 **Tables S2-S5. Statistical analyses (Excel files).** All statistics supplementary tables are linked to
1532 the analyses presented in specific figures, which are detailed in the first sheet of each file. The
1533 tables list sample sizes, statistic test used, effect sizes/odds ratios, confidence intervals, and p-
1534 values with and without multiple testing correction, depending on the analysis. Gene set
1535 enrichments and gene lists for specific analyses are listed under separate table headings.
1536 Additional data (i.e., GO enrichments) are also provided in some of the files, which are described
1537 below.

1538
1539 **Table S2. Statistics analysis for Figures 2, 3, S1, and S2 (Excel file).**

1540
1541 **Table S3. Statistics analysis for Figures 4 and S3 (Excel file).** **Table S3B** lists enriched GO
1542 terms for genes with secondary variants among 16p12.1 deletion probands with the five
1543 phenotypic domains.

1544
1545 **Table S4. Statistics analysis for Figures 5 and S4 (Excel file).** **Table S4A** details how
1546 psychiatric phenotypes were matched across questionnaire and EHR (ICD10) datasets across
1547 cohorts with different ascertainment.

1548

1549 **Table S5. Statistics analysis for Figures 6, S5, and S6 (Excel file). Tables S6C and S6D list**
1550 **enriched GO terms for genes with secondary variants among SSC probands with primary SNVs**
1551 **and CNVs and without primary variants (S6C), and Searchlight probands with 16p11.2 deletions**
1552 **and duplications (S6D).**
1553