

1     **Association between genes regulating neural pathways for quantitative traits of**  
2                                   **speech and language disorders**

3     Penelope Benchek, Ph.D.<sup>1\*</sup>, Robert P Igo Jr., Ph.D.<sup>1\*‡</sup>, Heather Voss-Hoynes, Ph.D.<sup>1\*</sup>,  
4     Yvonne Wren, Ph.D.<sup>5</sup>, Gabrielle Miller, M.S.<sup>2</sup>, Barbara Truitt, M.S.<sup>1</sup>, Wen Zhang,  
5     Ph.D.<sup>7</sup>, Michael Osterman, MPH<sup>1</sup>, Lisa Freebairn, M.S.<sup>2</sup>, Jessica Tag, M.A.<sup>2</sup>, H. Gerry  
6     Taylor, Ph.D.<sup>3,6</sup>, E. Ricky Chan, Ph.D.<sup>1</sup>, Panos Roussos, Ph.D.<sup>7,8</sup>, Barbara Lewis,  
7     Ph.D.<sup>2,4#</sup>, Catherine M. Stein, Ph.D.<sup>1#</sup>, Sudha K. Iyengar, Ph.D.<sup>1#</sup>

8             \* These authors contributed equally as first authors of this work  
9             # These authors contributed equally as senior authors of this work  
10            ‡ Recently deceased

11  
12     <sup>1</sup> Department of Population & Quantitative Health Sciences, and <sup>2</sup> Department of Psychological  
13     Sciences, Case Western Reserve University, Cleveland OH; <sup>3</sup> Department of Pediatrics, Case  
14     Western Reserve University, and Rainbow Babies & Children’s Hospital, University Hospital  
15     Case Medical Center, Cleveland, OH; <sup>4</sup> Cleveland Hearing and Speech Center, Cleveland, OH;  
16     <sup>5</sup> Bristol Dental School, Faculty of Health Sciences, University of Bristol, and Bristol Speech and  
17     Language Therapy Research Unit, North Bristol NHS Trust, Bristol, UK; <sup>6</sup> Nationwide Children’s  
18     Hospital Research Institute and Department of Pediatrics, The Ohio State University, Columbus,  
19     OH; <sup>7</sup> Department of Psychiatry, Friedman Brain Institute, and Department of Genetics and  
20     Genomic Science and Institute for Multiscale Biology, Icahn School of Medicine at Mount Sinai,  
21     New York, NY; <sup>8</sup> Mental Illness Research, Education, and Clinical Center (VISN 2 South),  
22     James J. Peters VA Medical Center, Bronx, NY

23  
24     **Corresponding authors:**

25     Catherine M. Stein, Ph.D. or Sudha K. Iyengar, Ph.D.  
26     Department of Population and Quantitative Health Sciences  
27     Case Western Reserve University  
28     2103 Cornell Rd  
29     Cleveland, OH 44016 USA

30  
31     **Running Title:** GWAS of communication disorder endophenotypes  
32

33 ABSTRACT

34 Speech sound disorders (SSD) manifest as difficulties in phonological memory and awareness,  
35 oral motor function, language, vocabulary, reading and spelling. Families enriched for SSD are  
36 rare, and typically display a cluster of deficits. We conducted a genome-wide association study  
37 (GWAS) in 435 children from 148 families in the Cleveland Family Speech and Reading study  
38 (CFSRS), examining 16 variables representing 6 domains. Replication was conducted using  
39 the Avon Longitudinal Study of Parents and Children (ALSPAC). We identified 18 significant  
40 loci (combined  $p < 10^{-8}$ ) that we pursued bioinformatically. We prioritized 5 novel gene regions  
41 with likely functional repercussions on neural pathways, some which colocalized with  
42 differentially methylated regions in our sample. Polygenic risk scores for receptive language,  
43 expressive vocabulary, phonological awareness, phonological memory, spelling, and reading  
44 decoding associated with increasing clinical severity. In summary, neural genetic influence on  
45 SSD is primarily multigenic and acts on genomic regulatory elements, similar to other  
46 neurodevelopmental disorders.

47

## 48 INTRODUCTION

49 Communication disorders are highly prevalent in the United States with approximately one in  
50 twelve children ages 3-17 years demonstrating a disorder<sup>1</sup>. The most common difficulties are a  
51 speech problem (5%) or language problem (3.3%). Speech Sound disorders (SSD) include  
52 both errors of articulation or phonetic structure (errors due to poor motor abilities associated  
53 with the production of speech sounds) and phonological errors (errors in applying linguistic rules  
54 to combine sounds to form words). SSD have a prevalence of approximately 16% in children 3  
55 years of age<sup>2</sup>, with an estimated 3.8% of children persisting with speech delay at 6 years of  
56 age<sup>3</sup>. More than half of these children encounter later academic difficulties in language,  
57 reading, and spelling<sup>7-11</sup>. Because of the relative rarity of persistent speech problems and their  
58 correlation with other communication domains, endophenotypes are key to the study of genetic  
59 underpinnings.

60  
61 Vocabulary is core to speech acquisition<sup>4</sup>. Children with difficulties in speech sound  
62 development often have difficulties with oral language and later reading and spelling disability<sup>2,5-</sup>  
63 <sup>8</sup>. Thus, speech, language, reading, and spelling measures are highly correlated and often  
64 have common genetic associations<sup>9,10</sup>. Moreover, speech and other communication  
65 phenotypes follow a developmental trajectory, where some speech and language disorders  
66 resolve with age, whereas others persist; genetic influences on the less easily resolved  
67 manifestations are generally stronger<sup>11,12</sup>. Because of the common genetic underpinnings and  
68 pathologic associations between speech and other communication phenotypes, it is conceivable  
69 that genetic replication interweaves with different communication measures. Of 7 known  
70 GWASs, none overlap in their top results (at  $p < 5 \times 10^{-5}$ , see Table 3<sup>13</sup>), because they only  
71 focused on a limited number of phenotypes, or these measures were assessed at different ages  
72 (either pre-school or early school-age)<sup>13-20</sup>, they only present results from one or a few

73 measures and/or a binary trait; thus, the complexity of shared genetic influences is poorly  
74 understood. Most have not focused on children with SSD, particularly measures of articulation.  
75 Our sample represents a unique set of deeply phenotyped individuals with information on 6  
76 domains that form the core of speech and language.

77  
78 SSD are likely due to deficits in both motor ability and broader neural dysfunction. While motor  
79 deficits contribute to problems in speech production, abnormalities in other neural systems likely  
80 influence formation of phonological representation, which is common to SSD as well as reading  
81 and language impairment. We hypothesize that genetic regulation of these neural pathways is  
82 associated with variation common to speech, language, reading, and spelling ability. We  
83 conducted a GWAS in the Cleveland Family Speech and Reading Study (CFSRS), a cohort  
84 ascertained through a proband with SSD. We also conducted a methylome-wide study (i.e.  
85 MWAS) to determine the functional implications of these genetic associations, and replicated  
86 findings in a population-based cohort. We utilized a family-based cohort as our discovery  
87 sample because we hypothesized it would be enriched for disease-associated variants<sup>21,22</sup>. In  
88 these analyses, we identified new candidate genes for correlated communication  
89 endophenotypes, and bioinformatic annotation of these loci revealed that regulation of neural  
90 pathways is associated with variation in these measures.

91

## 92 SUBJECTS AND METHODS

### 93 **Subject ascertainment – Cleveland Family Speech and Reading Study**

94 From the Cleveland Family Speech and Reading Study (CFSRS)<sup>23-28</sup>, we examined 435  
95 individuals from 148 families who had both DNA and endophenotype data available (Table 1).  
96 As previously described, families were ascertained through a proband with SSD identified from

97 caseloads of speech-language pathologists in the Greater Cleveland area and referred to the  
98 study; detailed inclusion criteria are provided in the Supplemental Methods. Diagnosis of CAS  
99 was confirmed by an experienced licensed speech-language pathologist upon enrollment into  
100 the study. Socioeconomic status was determined at the initial assessment based on parent  
101 education levels and occupations using the Hollingshead Four Factor Index of Social Class<sup>29</sup>.  
102 This study was approved by the Institutional Review Board of Case Medical Center and  
103 University Hospitals and all parents provided informed consent and children older than 5 years  
104 provided assent.

105

## 106 **Communication Measures in CFSRS**

107 We examined diadochokinetic rates using the *Robbins and Klee Oral Speech Motor Control*  
108 *Protocol*<sup>30</sup> or *Fletcher Time-by-Count Test of Diadochokinetic Syllable Rate*<sup>31</sup>. The merged  
109 variable is referred to as DDK. Expressive vocabulary was assessed with the *Expressive One*  
110 *Word Picture Vocabulary Test-Revised (EOWPVT)*<sup>32</sup> and receptive vocabulary with the  
111 *Peabody Picture Vocabulary Test- Third Edition (PPVT)*<sup>33</sup>, and phonological memory with the  
112 *Nonsense Word Repetition (NSW)*<sup>34</sup>, *Multisyllabic Word Repetition (MSW)*<sup>34</sup>, and *Rapid Color*  
113 *Naming*<sup>35</sup> task. In addition to examining the total number of words correct for the MSW and  
114 NSW, we also examined the percent phonemes correct for both of these tasks (NSW-PPC and  
115 MSW-PPC, respectively). Phonological awareness was assessed using the *Elision* subtest of  
116 the *Comprehensive Test of Phonological Processing – 2<sup>nd</sup> Edition*<sup>36</sup>. Reading was assessed  
117 using the *Woodcock Reading Mastery Test-Revised, Word Attack subtest (WRMT-AT)* and  
118 *Word Identification Subtest (WRMT-ID)*, the *Reading Comprehension subtest (WIAT-RC)* and  
119 *Listening Comprehension subtest (WIAT-LC)* of the *Wechsler Individual Achievement Test*<sup>37</sup>  
120 Spelling was assessed on the *Test of Written Spelling-3 (TWS)* using the total score<sup>38</sup>.  
121 Expressive and receptive language were assessed using the *Test of Language Development*

122 (*TOLD*<sup>39</sup>) and *Clinical Evaluation of Language Fundamentals-Revised (CELF*<sup>40</sup>). referred to as  
123 the CELF-E (expressive) and CELF-R (receptive), respectively. Additional details about these  
124 measures are provided in the Supplemental Methods. For each of our tests we selected the first  
125 available assessment for each individual (Supplemental Table 1).

126

## 127 **GWAS analysis**

128 Genotyping methods and quality control (QC) are described in the Supplemental Methods.  
129 Principal components (PC) obtained from principal component analysis (PCA) and the genetic  
130 relationship matrix (GRM) were generated using genotyped markers that met QC criteria. We  
131 used PC-AiR and PC-Relate from the Bioconductor package GENESIS<sup>41</sup> to generate our PCs  
132 and GRM, respectively. PC-AiR accounts for sample relatedness to provide ancestry inference  
133 that is not confounded by family structure, while PC-Relate uses the ancestry representative  
134 PCs from PC-AiR to provide relatedness estimates due only to recent family (pedigree)  
135 structure.

136

137 To examine cross-trait correlation, we used GCTA<sup>42</sup> to run a bivariate REML analysis for each  
138 pair of tests and tested for genetic correlations equal to 0. GCTA's bivariate REML analysis  
139 estimates the genetic variance of each test and the genetic covariance between the two tests  
140 that can be captured by all SNPs<sup>43</sup>. Here we included all SNPs with MAF  $\geq$  0.01. The genetic  
141 variance/covariance calculated was adjusted for sex and the first two PCs.

142

143 We used RVTtests, version 2.0<sup>44</sup> to run our GWAS. We specifically relied on RVTtest's  
144 Grammar-gamma test<sup>45</sup>, which performs a linear mixed model association test while allowing for  
145 genotype dosages and accounting for family structure using the Genetic Relationship Matrix

146 (GRM). Because each of our tests were age-normed we included only sex and the first two PCs  
147 as covariates in our regression models.

148  
149 In addition, we generated endophenotype-based polygenic risk scores (PRS) in the European  
150 subset of the CFSRS where genotype data, as well as clinical group data (no disorder, SSD  
151 only, language impairment (LI) only, SSD+LI, CAS) were available. Risk scores were derived  
152 from association statistics from our CFSRS GWASs (see GWAS methods section for details)  
153 and were constructed using PLINK 1.9<sup>46</sup> (clump and score functions). Additional details are in  
154 the Supplemental Methods. These polygenic risk scores were used to examine the hypothesis  
155 that an increase in PRS score would associate with more complex clinical phenotypes when  
156 comparing SSD only versus SSD+LI and CAS.

157

## 158 **Statistical analysis of Methylome-wide data**

### 159 *Methylome-wide association study (MWAS)*

160 Quality control analysis of methylation data is detailed in Supplemental Methods. We tested for  
161 association between CpG beta values and endophenotypes using the linear mixed model  
162 approach of GRAMMAR-Gamma<sup>45</sup> as implemented in RVtests<sup>44</sup>. Because our phenotypes  
163 were age-normed, we did not adjust for age, but rather for sex and one to four PCs. We also  
164 examined methylation-QTLs (meQTL) as described in the Supplemental Methods.

165

### 166 **Replication dataset – ALSPAC**

167 To replicate our GWAS findings, we obtained data from the Avon Longitudinal Study of Parents  
168 and Children (ALSPAC). The ALSPAC study was a prospective population-based birth cohort  
169 of babies born from > 14,000 pregnancies between April 1991-December 1992, who were  
170 followed prospectively with a wide battery of developmental tests, parental questionnaires, child-

171 completed questionnaires, and health outcomes<sup>47-49</sup>. Pregnant women resident in Avon, UK  
172 with expected dates of delivery 1st April 1991 to 31st December 1992 were invited to take part  
173 in the study. The study website contains details of all the data that is available through a fully  
174 searchable data dictionary  
175 (<http://www.bris.ac.uk/alspac/researchers/data-access/data-dictionary>). Ethical approval for the  
176 study was obtained from the ALSPAC Ethics and Law Committee and Institutional Review  
177 Board of Case Medical Center and University Hospitals. Because this was a birth cohort, all  
178 children were included, regardless of diagnosis. We obtained both parental report data on  
179 speech development in the children, and also communication measures similar to those that we  
180 analyzed (see Communication Measures above and Supplemental Table 3). As this was a  
181 longitudinal study, different measures were given at different ages, and when the same domain  
182 was tested at two different ages, the identical measure was not used. At some ages, only  
183 random subsets were selected, so the sample size available from each age is not the same. In  
184 Supplemental Table 4, we list the measures given in the CFSRS battery along with the most  
185 similar measure given in ALSPAC.

186

#### 187 *GWAS in ALSPAC data*

188 QC analyses of ALSPAC data are described in Supplemental Methods. Because of the format  
189 of data that were provided, we used slightly different methods for statistical analyses. Genetic  
190 association testing was performed using linear regression in Hail 0.1. Covariates adjustments  
191 included sex and the first two PCs. Age was not a consideration as ALSPAC is a longitudinal  
192 birth cohort study and age differences were negligible for any given measure.

193

#### 194 **Functional annotation and results integration**

195 In this analysis, we considered CFSRS the discovery sample, since families were ascertained  
196 through a child with SSD, and used ALSPAC as the replication sample. We identified  
197 associated loci with SNPs significant at  $p < 10^{-5}$  in CFSRS and  $p < 0.05$  in ALSPAC, with effects in  
198 the same direction.

199

#### 200 *Functional annotation*

201 Because the majority of our findings are intergenic and/or fall in noncoding regions, we relied on  
202 annotation tools FUMA and HaploReg to characterize which genes our variants might affect, as  
203 well as variants' functionality. We utilized FUMA<sup>50</sup> for mapping genes to our variants based on  
204 genomic proximity, eQTL evidence and chromatin interactions evidence. Default settings in  
205 FUMA were used, with the exception of tissue specificity. We hypothesized that gene  
206 expression and regulation would be most relevant in brain and neural tissues, as well as  
207 muscles related to speech. In FUMA we focused on eQTL and chromatin interaction evidence  
208 in our target tissues (brain, muscle and esophagus). HaploReg v.4.1 was used to examine the  
209 chromatin state evidence predicting whether the variant fell in a promoter or enhancer region. In  
210 HaploReg we focused on chromatin state evidence in our target tissues (brain and muscle).

211

#### 212 *Locus prioritization*

213 In order to further prioritize and synthesize our findings, we annotated associated loci as  
214 described above, including annotation of associated effects of these loci in the literature, and  
215 incorporated supportive findings from our MWAS. We summarize findings in Table 2, and  
216 generated a simple locus priority score as the number of times a locus included an enhancer  
217 and/or promoter, included an eQTL, was previously associated with a communication disorder  
218 and/or neuropsychiatric disorder, showed eQTL or chromatin state evidence specific to brain

219 and/or neural tissues, mapped to a gene that was a *FOXP2* target in brain tissue<sup>51-53</sup>, and a  
220 meQTL in that region (at  $p < 5 \times 10^{-5}$ ) with an associated methylation site (at  $p < 0.05$ ) with the  
221 same phenotype as the associated GWAS loci. We applied the EpiXcan pipeline<sup>54</sup> to identify  
222 eQTLs with our associated SNPs that are differentially expressed in the dorsolateral prefrontal  
223 cortex (DLPFC)<sup>55</sup> (Supplemental Methods).

224

## 225 RESULTS

226 The CFSRS sample included 435 subjects from 148 families (Table 1). Of these, 27% had SSD  
227 only, 4% had LI only, 16% had SSD+LI without CAS, and 11% had CAS (Table 1). Of the  
228 subjects in the ALSPAC sample, the prevalence of speech problems by parental report varied  
229 from 4%-6% (Supplemental Table 3).

230

### 231 **Genetic correlation analysis reveals new relationships among endophenotypes**

232 Genetic correlation analysis revealed that while many of the patterns of correlation were  
233 consistent with phenotypic correlations we have previously reported<sup>10</sup>, polygenic correlations  
234 enable a deeper understanding of these measures, which will inform examination of replication  
235 of association effects both within the CFSRS data set and with measures from ALSPAC (Figure  
236 1). For example, while previous studies have demonstrated a strong genetic correlation  
237 between reading and spelling measures, polygenic correlation analysis additionally reveals  
238 correlations between those skills and Elision. Not surprisingly, expressive and receptive  
239 language, as measured on the CELF, are strongly correlated with vocabulary (EOWPVT and  
240 PPVT) in addition to reading (WRMT-AT and WRMT-ID). Vocabulary is also strongly correlated  
241 with listening comprehension (WIAT-LC).

242

### 243 **Most significant findings from GWAS reveal 5 new candidate genes**

244 The majority of associated SNPs ( $p < 10^{-5}$ ) were intergenic, with a lesser number of intronic  
245 SNPs (Supplemental Figure 2). Noncoding regions harboring a significant proportion of risk  
246 alleles is consistent with previous findings related to neuropsychiatric disease and behavioral  
247 traits<sup>56</sup>. We focused on SNPs that had a  $p$ -value  $< 1 \times 10^{-5}$  in CFSRS with replication with a related  
248 trait in ALSPAC ( $p < 0.05$ ), or Fisher combined  $p$ -values  $< 1 \times 10^{-7}$ , that had functional relevance  
249 based on our gene priority score (Table 2).

250

251 Among the 5 prominent loci, all had enhancers or promoters for muscle, brain, and/or neuronal  
252 progenitor cells, 4 out of 5 had significant methylation and meQTL effects, and 3 out of 5 had  
253 eQTLs for brain and/or skeletal-muscle tissue (Figure 2, Supplemental Table 5). EpiXcan  
254 analysis suggested that the SNP in the chromosome 1 *IFI6* region is associated with expression  
255 in the DLPF cortex (Elision  $p = 0.018$ , TWS  $p = 0.008$ ; Supplemental Tables 6 and 7). The first  
256 region on chromosome 14, including *NFKBIA* and *PPP2R3C*, shows significant chromatin  
257 interaction mapping in adult cortex tissue. *NFKBIA*, which codes for a component of the NF- $\kappa$ B  
258 pathway, is associated with neurogenesis, neuritogenesis, synaptic plasticity, learning and  
259 memory<sup>57</sup>. The second region on chromosome 14 includes *PP2R3C*, which is within the  
260 topologically associating domain (TAD) boundary of the *NFKBIA* locus in Hippocampus and  
261 DLPFC. EpiXcan analysis showed *NFKBIZ*, a gene in the same pathway as *NFKIBA*, is also  
262 associated with expression in the DLPFC (Elision  $p = 0.000452$ , TWS  $p = 0.004939$ ; Supplemental  
263 Tables 6 and 7).

264

### 265 **Replication of previous communication disorder loci**

266 *ATP2C2* was associated with WRMT-ID ( $p=7.6 \times 10^{-8}$ ), WRMT-AT ( $p=4.6 \times 10^{-5}$ ), and Elision  
267 ( $p=4.6 \times 10^{-5}$ ), consistent with prior literature<sup>58</sup> (Supplemental Figures 3 and 4). Similarly,  
268 *CYP19A1* was associated with WRMT-AT ( $p=2.8 \times 10^{-5}$ ), Elision ( $p=3.3 \times 10^{-4}$ ), and WRMT-ID  
269 ( $p=5.0 \times 10^{-4}$ ), validating a previous association<sup>59</sup>. *CNTNAP2* was associated with CELF-R  
270 ( $p=5.2 \times 10^{-6}$ ), and DDK ( $p=2.9 \times 10^{-5}$ ), replicating a previous association<sup>58</sup>. While SNPs within  
271 *ROBO1* and *ROBO2* were not significantly associated with our measures, SNPs in the  
272 intergenic region were associated with WRMT-ID ( $p=3.6 \times 10^{-6}$ ); *ROBO1* was originally  
273 associated with dyslexia while *ROBO2* was originally associated with expressive vocabulary<sup>20,60</sup>.  
274 Finally, SNPs within the *DCDC2-KIAA0319-TTRAP* and in *FOXP2* regions were associated with  
275 various traits at  $p < 0.01$ . Within the ALSPAC cohort, a different pattern of replication emerged  
276 (Supplemental Figure 5), with sometimes different SNPs and/or different phenotypes than those  
277 associated with CFSRS.

278  
279 In addition, we examined loci (genes and/or SNPs) associated in recently published GWAS  
280 studies of language and reading<sup>13-20</sup> (Supplemental Table 8); we restricted our examination to  
281 the CFSRS data, since the ALSPAC data were included in some of the published studies. In  
282 these analyses, we often observed cross-trait replication, with most genes originally associated  
283 with dyslexia, and associated with other traits in our sample. These included *ZNF385D*<sup>14</sup>, which  
284 was associated with all CFSRS traits at  $p < 0.005$ , *CDH13*<sup>19</sup>, associated with all CFSRS traits at  
285  $p < 0.005$ , *GRIN2B*<sup>15</sup>, associated with TWS, EOWPVT, and Elision at  $P < 0.0005$  and all CFSRS  
286 traits at  $P < 0.05$ , *NKAIN*<sup>15</sup>, associated with CELF-R at  $9.7 \times 10^{-5}$  ( $rs16928927$   $p=1 \times 10^{-4}$ ) and  
287 WIAT-RC ( $p=4 \times 10^{-4}$ ), and *MACROD2*<sup>17</sup> associated with all CFSRS traits at  $p < 0.005$ .

288

289 **Polygenic risk scores are associated with increasing clinical severity**

290 In Figure 3, we illustrate polygenic risk scores (PRS) for 6 endophenotypes representing the  
291 major domains (receptive language, expressive vocabulary, phonological awareness,  
292 phonological memory, spelling, and reading decoding), by quintile, across the clinical subgroups  
293 (all endophenotypes are illustrated in Supplemental Figure 6). Generally, we found that  
294 polygenic load, indicated by increasing risk scores, was associated with clinical severity  
295 ( $p < 1 \times 10^{-8}$  by ANOVA), with typical children having the lowest scores, followed by children with  
296 SSD-only, and children with SSD+LI and CAS having the greatest scores. The exception to this  
297 trend is receptive language, where the genetic load is greatest for children with LI, for whom  
298 receptive language is a focal deficit. Thus, in general, an increase in PRS score is associated  
299 with greater clinical severity.

300

## 301 DISCUSSION

302 Communication disorders are genetically complex, manifested by a variety of deficiencies in  
303 articulation, vocabulary, receptive and expressive language, phonological awareness, reading  
304 decoding and comprehension, and spelling. This GWAS ascertained children through an  
305 earlier-presenting clinical disorder and examined several key communication measures, and is  
306 thus one of the first studies of its kind. This study is also novel in that it is the first GWAS to  
307 include a measure of phonological awareness, as well as a motor speech measure. By  
308 analyzing several endophenotypes together, we can draw conclusions about the common  
309 genetic basis across these seemingly dissimilar skills. Here, we have identified five new  
310 candidate regions, some containing multiple genes, that have connections to neurological  
311 function and regulation of neurological pathways. We also found that increased polygenic load  
312 is associated with more severe communication disorders. Finally, by examining genetic  
313 correlations among these traits, we conclude that different domains of communication have

314 some common genetic influences. All of these aspects together add new clarity regarding the  
315 genetic underpinnings of speech and language skills.

316

317 First, the novel candidate genes that we have identified all have roles in neurological function as  
318 evidenced by expression levels of those genes in brain and/or neural tissue, and associations  
319 with other communication and/or psychiatric phenotypes. This commonality between  
320 communication traits and brain and neural pathways was also demonstrated by a mouse study  
321 of vocalization<sup>61</sup>, and pleiotropy between brain, learning, and psychiatric phenotypes was  
322 recently demonstrated by a large GWAS of brain phenotypes<sup>62</sup>. Existence of enhancers,  
323 promoters, and methylation effects in the associated regions further emphasizes the importance  
324 of regulatory effects on these traits. Deletions spanning *SETD3* and *CCNK* have been  
325 associated with syndromic neurodevelopmental disorders<sup>63</sup> and variants in *SETX*, within this  
326 same family of genes, have been associated with CAS<sup>64</sup>. In addition, *CCNK* is in the *FOXP2*  
327 pathway in brain tissue<sup>51-53</sup>. *NFKBIA* is involved in regulation of the NF- $\kappa$ B pathway, which is  
328 involved a number of brain-related processes including neurogenesis, neuritogenesis, synaptic  
329 plasticity, learning, and memory<sup>65</sup>. *PPP2R3C* has been associated with schizophrenia<sup>66</sup>. *IFI6*  
330 expression has been associated with autism<sup>67</sup> and overexpression of *IFI6* in the brain is present  
331 in chronic neurodegeneration<sup>68</sup>. Finally, *DACT1* may be involved in excitatory synapse  
332 organization and dendrite formation during neuronal differentiation<sup>69</sup> and is mainly expressed  
333 within the first two trimesters of pregnancy, just before the first evidence of speech processing is  
334 observed in preterm neonates<sup>70</sup>. Interestingly, *SETD3*, *NFKBIA*, and *IFI6* are all also tied to the  
335 immune system, and a recent study identified an excess of T cells in brains of individuals with  
336 autism<sup>71</sup>.

337

338 Second, understanding the genetic architecture across these endophenotypes is essential for  
339 understanding how loci are associated with different measures in different study cohorts or  
340 across the developmental trajectory. Strong genetic correlations are observed between  
341 spelling, reading comprehension and decoding, expressive and receptive language, vocabulary,  
342 and phonological awareness. The strongest replications were for a variety of measures  
343 collected in CFSRS with ALSPAC from older youth. Consistent with these findings, we  
344 previously demonstrated that spelling at later ages has a higher estimated heritability than  
345 spelling at school-age<sup>11</sup>. Measures administered in older youth may also be more sensitive to  
346 variations in clinical manifestation of SSD. Examination of the ALSPAC measures suggests that  
347 many of those administered at younger ages may have tapped different domains than intended,  
348 or may have been less sensitive to later emerging reading and spelling skills. Methods of cohort  
349 ascertainment may also be important in comparing our findings to those of other studies. Our  
350 families were ascertained through a child with SSD whereas other studies ascertained subjects  
351 through LI or dyslexia. These different ascertainment schemes affect both the available  
352 measures, as well as the distribution of scores and power to detect association. Since both LI  
353 and dyslexia emerge later than SSD, longitudinal studies that ascertain through a proband with  
354 SSD will be able to capture variants associated with all three disorders, as there is high  
355 comorbidity. In addition to the plethora of studies ascertaining children at a variety of ages,  
356 which has an impact on the heritability of traits<sup>10</sup>, these studies use a wide variety of measures,  
357 even for the same endophenotype. Moreover, these studies have been conducted in  
358 populations that speak different languages of varying orthographic transparency, which makes  
359 them difficult to compare. As noted by Carrion-Castillo et al.<sup>13</sup>, most of the novel loci identified  
360 through GWAS have been unique to each study, and these aforementioned issues may explain  
361 that lack of replication. Thus, examination of the genetic correlation matrix is essential for  
362 interpretation of results across studies, as it is nearly impossible to analyze the same exact  
363 traits, as we have demonstrated with our replication study cohort (ALSPAC).

364

365 Third, we replicated candidate genes that had been previously primarily associated with reading  
366 and/or language impairment: *CNTNAP2*, *ATP2C2*, and *CYP19A1*. These analyses extend  
367 previous findings to show that these genes are associated with articulation (*CNTNAP2*) and  
368 phonological awareness (*ATP2C2* and *CYP19A1*). This further illustrates the pleiotropic nature  
369 of these genes. While we did not observe association with SNPs within the coding regions of  
370 *ROBO1* and *ROBO2*, we did observe significant associations with SNPs between these two  
371 genes, which may have regulatory influences on *ROBO1/ROBO2*. We also replicated ( $p < 5 \times 10^{-3}$ )  
372 loci identified in recent GWAS of reading and/or language traits. Similar to another  
373 association study between *FOXP2* variants and language<sup>72</sup>, we did not observe statistically  
374 significant association between *FOXP2* and measures in CFSRS, though there was replication  
375 of some traits at a less stringent ( $p < 0.01$ ) level<sup>72</sup>.

376

377 Finally, our analysis of polygenic risk scores shows strong associations between these risk  
378 scores and clinical outcomes of increasing severity. Because of the strong significance of these  
379 findings, this suggests that the genetic architecture of communication disorders maybe largely  
380 polygenic, which may additionally explain the lack of replication and/or genome-wide  
381 significance. While other studies have examined polygenic risk scores associated with  
382 language<sup>15,73</sup>, ours is the first to examine polygenic risk associated with other communication  
383 endophenotypes. It is noteworthy that our associated SNPs fell outside of gene coding regions  
384 but resided in regulatory regions, even having potential regulatory effects themselves. This  
385 further illustrates the genetic complexity of communication disorders; perhaps the search for  
386 single gene dysfunction is misplaced, and rather regulatory functions are more relevant.

387

388 This study has several limitations. The sample size of the CFSRS cohort was modest,  
389 potentially reducing power. There was not clear correspondence between measures obtained  
390 in ALSPAC with those in CFSRS, necessitating consideration of cross-trait replication. We  
391 restricted analyses in both cohorts to individuals of European descent because of low sample  
392 size in other ethnic groups, reducing generalizability.

393  
394 In summary, this first GWAS of communication measures ascertained through families with SSD  
395 identified five new candidate genes, all with potential relevance in central nervous system  
396 function. Polygenic risk is strongly associated with more severe speech and language  
397 outcomes. Careful consideration of genetic correlation among domains of verbal and written  
398 language shows that these loci have general effects on communication, not specific to any  
399 single domain, suggesting a common genetic architecture. Further research is needed to more  
400 closely examine the impact of regulatory variants on these outcomes.

401

## 402 ACKNOWLEDGMENTS

403 We would like to thank the families who have so generously participated in this study for many  
404 years. This research was supported by the Genomics Core Facility of the CWRU School of  
405 Medicine's Genetics and Genome Sciences Department. This work made use of the High  
406 Performance Computing Resource in the Core Facility for Advanced Research Computing at  
407 Case Western Reserve University. This work was supported by NIH grant R01DC000528  
408 awarded to Dr. Lewis and R01DC012380 awarded to Dr. Iyengar. We are extremely grateful to  
409 all the families who took part in the ALSPAC study, the midwives for their help in recruiting  
410 them, and the whole ALSPAC team, which includes interviewers, computer and laboratory  
411 technicians, clerical workers, research scientists, volunteers, managers, receptionists and

412 nurses. The UK Medical Research Council and Wellcome (Grant ref: 217065/Z/19/Z) and the  
413 University of Bristol provide core support for ALSPAC. This publication is the work of the  
414 authors and Dr. Sudha Iyengar will serve as guarantor for the contents of this paper. GWAS  
415 data for ALSPAC was generated at the Genotyping Facilities at Wellcome Sanger Institute.

416

#### 417 DATA AVAILABILITY

418 Data from the Cleveland Family Speech and Reading study are not available for broad genetic  
419 data sharing because of IRB restrictions. Please contact the corresponding author, Dr. Sudha  
420 Iyengar, to request data, which will require an IRB application.

421

#### 422 CONFLICT OF INTEREST

423 The authors have no conflicts of interest to report.

424

#### 425 MATERIALS AND CORRESPONDANCE

426 Please contact Dr. Sudha Iyengar, [ski@case.edu](mailto:ski@case.edu), regarding access to summary statistics.

427 TABLES

428 **Table 1. Characteristic table for CFSRS GWAS sample**

---

N*	435
Number of Families	148
Age range	[2.5, 64]
Female N (%)	194 (45%)
Speech Disorder Subgroup N (%)	
CAS	47 (11%)
SSD + LI (no CAS)	70 (16%)
SSD only	119 (27%)
Lang only	17 (4%)
No CAS/SSD/Lang	177 (41%)
Missing	5 (1%)
Hollingshead SES	
1 (lowest)	3 (1%)
2	30 (7%)
3	67 (15%)
4	167 (38%)
5 (highest)	147 (34%)
Missing	21 (5%)

---

\*Sample considered is the union of all samples across the 16 tests. Specific test sample sizes and age ranges are shown in supplemental Table 1.

429

430

431 **Table 2. Annotation of most significant loci with replication in CFSRS and ALSPAC**

Locus (Chr location)	Gene(s)	# Associated SNPs	# independently associated SNPs (after conditional analysis)	Gene priority score	Expression in brain / neural tissue	Associated with Communication and/or psych phenotype	Associated with multiple CFSRS traits	Promoter (esophagus, muscle, brain, neural)	Enhancer (esophagus, muscle, brain, neural)	eQTL (esophagus, muscle, brain, neural)	Target of FOXP2 (brain)	Methylation / meQTL
1:30732871	LINC01648;MATN1	1	1	3	1	1	0	1	0	0		
1:55494735*	BSND;PCSK9	5	1	4	1	0	1	1	0	0		1
1:146988760	LINC00624	1	1	5	1	1	0	0	1	1		1
1:159028378	IFI16, AIM2	23	1	5	0	1	0	0	1	1		2
2:143378805	LRP1B;KYNLU	4	1	4	1	1	1	1	0	0		
2:169280713	STK39;CERS6	1	1	3	1	1	0	0	1	0		
3:1942898	CNTN6;CNTN4	1	1	3	1	1	0	0	1	0		
3:39743136	MOBP;MYRIP	1	1	2	0	1	0	0	0	0		1
4:27297733	LINC02261;MIR4275	9	1	2	1	0	0	0	1	0		
4:73572756	ADAMTS3;COX18	7	1	4	1	1	1	0	1	0		
4:77531588	SHROOM3	1	1	3	1	0	0	1	1	0		
5:72144005	TNPO1	1	1	4	1	1	0	1	1	0		
5:132043351	KIF3A	1	1	4	1	1	0	0	1	1		
5:170102906	KCNIP1	2	1	1	0	0	0	0	0	0		1
5:172924967	MIR8056;LOC285593	15	1	4	1	0	0	1	1	0		1
7:123604182	SPAM1	10	1	4	1	1	0	1	1	0		
7:154706515	DPP6;PAXIP1-AS2	1	1	5	1	1	0	1	1	0	1	
9:114335864	PTGR1;ZNF483	0	1	6	1	1	0	1	1	1		1
10:46027420	MARCH8	2	1	4	1	1	0	0	1	1		
12:21002703	SLCO1B3	2	1	1	0	0	0	0	0	0		1
12:103677691**	LOC101929058; C12orf4	1	1	1	0	1	0	0	0	0		
12:131389783	RAN; ADGRD1	1	1	0	0	0	0	0	0	0		
13:28329109	POLR1D; GSX1	18	1	1	0	0	0	1	0	0		
13:79839523	LINC00331; RBM26	10	1	2	0	1	0	0	0	1		
14:35837476	PSMA6; NFKBIA	26	1	7	1	1	0	1	1	1		2
14:59210646	DACT1; LINC01500	7	1	5	1	1	0	1	1	0		1
14:93195374	LGMN	1	1	4	1	1	0	1	1	0		
14:94993936*	SERPINA12; SERPINA4	5	1	5	0	1	1	1	1	0		1
14:99858970	BCL11B; SETD3	1	1	8	1	1	0	1	1	1	1	2
16:77231207	MON1B	1	1	7	1	1	0	1	1	1		2
18:4023876	DLGAP1	1	1	3	1	1	0	0	1	0		
18:40822793	RIT2; SYT4	1	1	1	0	1	0	0	0	0		
18:56462735	MALT1; LINC01926	1	1	3	0	1	0	0	1	0		1

# Associated SNPs include those associated in CFSRS (p<10<sup>-5</sup>) and Alspac (p<0.05) or fisher combined (p< 10<sup>-7</sup>)

\*Alspac led locus. No CFSRS SNPs showed association at P < 10<sup>-5</sup>.

\*\* CFSRS P = 1.3 \* 10<sup>-5</sup> and Alspac P = 5.8 \* 10<sup>-5</sup> (fisher P = 1.6 \* 10<sup>-8</sup>).

432

433 **FIGURES (attached separately)**

434 **Figure 1. Genetic correlation matrix across traits in CFSRS.** Figure 1 shows cross-trait  
435 correlation results for each pair of tests using GCTA's bivariate REML analysis. Cross-trait  
436 correlation was tested under the null hypothesis of 0 correlation. Circles shown are for results  
437 significant at  $P < 0.05$ , with increasing diameter/color corresponding with increasing correlation  
438 (circles omitted otherwise).

439

440 **Figure 2. Locus zoom plots for most significant findings.** Figure 2 shows association  
441 results for the top loci. P-values displayed are for CFSRS and are for the test for which the top  
442 SNP was observed. Circles show P-values for SNP associations and triangles show P-values  
443 for methylation associations (specifically those for which the top SNP is an meQTL). The larger  
444 plot shows the top SNP for each region +/- 200 kb. The window highlights the region that spans  
445 significant association results ( $P \leq 1 \times 10^{-5}$  in CFSRS. **A.** IFI16 region (window spans  
446 chr1:159001292-159028378) rs855865 was associated with NSW in CFSRS ( $p = 7 \times 10^{-6}$ ) and  
447 with vocabulary (WISC-V) in ALSPAC ( $p = 0.01$ ). This region also includes an meQTL  
448 (rs12124059,  $p = 4 \times 10^{-8}$ ) for methylation marker cg07196514, and this methylation marker was  
449 also associated with NSW ( $p = 0.018$ ). **B.** NFKBIA region (window spans chr14:35770806-  
450 35846092). rs57645874 was associated with Elision in CFSRS ( $p = 1 \times 10^{-6}$ ) and with reading  
451 accuracy (NARA-A) in ALSPAC ( $p = 0.02$ ). This region also contains an meQTL, rs4981288, for  
452 cg07166546 ( $p = 2 \times 10^{-50}$ ), and this methylation marker was associated with Elision ( $p = 3 \times 10^{-5}$ ),  
453 TWS ( $p = 0.0005$ ) and WRMT-ID ( $p = 0.002$ ). **C.** DACT1 region (window spans chr14:59210335-  
454 59221002). rs856379 was associated with MSW in CFSRS ( $p = 3 \times 10^{-6}$ ) and with nonword  
455 reading (ALSPACread) in ALSPAC ( $p = 0.036$ ). This SNP is an meQTL for methylation marker  
456 cg13972423 ( $p = 3 \times 10^{-5}$ ), **D.** SETD3 region (window spans chr14:99858970-99942692).  
457 rs1257267 was associated with WRMT-AT in CFSRS ( $p = 6.58 \times 10^{-6}$ ) and with nonsense word

458 repetition (CNrep5) in ALSPAC ( $p=0.05$ ). While only 1 SNP replicated between CFSRS and  
459 ALSPAC, 14 additional SNPs showed association in CFSRS at  $p<10^{-5}$ . This SNP is an meQTL  
460 for cg18949721 ( $p=4\times 10^{-12}$ ), which was also associated with WRMT-AT ( $p=0.003$ ). **E. MON1B**  
461 region (window spans chr16:77231207-77248555). rs4888606 was associated with MSW in  
462 CFSRS ( $p=9 \times 10^{-6}$ ) and with nonword reading (ALSPACread) in ALSPAC ( $p=0.046$ ). While  
463 only 1 SNP replicated between CFSRS and ALSPAC, 18 additional SNPs showed association  
464 in CFSRS at  $p<10^{-5}$ . This SNP falls in an intron of *MON1B* and is an meQTL for cg06128999  
465 ( $p=4\times 10^{-23}$ ) and cg05007098 ( $p=1\times 10^{-15}$ ), which were also associated with MSW ( $p=0.045$  and  
466  $p=0.12$ , respectively). Functional annotation is in Supplemental Figure 2.

467

468 **Figure 3. Polygenic risk scores across major domains.** We constructed polygenic risk  
469 scores for 587 individuals who were both genotyped and had clinical subgroup information  
470 available. Polygenic risk scores are displayed by quantile across the clinical subgroups for six  
471 endophenotypes representing the major domains (A Receptive language; B Expressive  
472 vocabulary; C Phonological awareness; D Phonological memory; E Spelling; F Reading  
473 decoding).

474

475

#### 476 DESCRIPTION OF SUPPLEMENTAL DATA

477 Supplemental Methods: Describes behavioral phenotypes in detail and detailed methods for  
478 genetic methylation analysis

479 Supplemental Tables and Figures:

480 Supplemental Table 1. Descriptive statistics for CFSRS measures

481 Supplemental Table 2. Results of methylation analysis of candidate gene regions

482 Supplemental Table 3. Descriptive statistics for ALSPAC sample

483 Supplemental Table 4. Correspondence between CFSRS and ALSPAC measures

484 Supplemental Table 5. Annotation of functional implications of most significant loci from GWAS

485 Supplemental Table 6. PsychEncode EpiXcan method using Meta-analysis results of Elision

486 GWAS

487 Supplemental Table 7. PsychEncode EpiXcan method using Meta-analysis results of TWS

488 GWAS

489 Supplemental Table 6 – Association results from regions identified from published GWAS of  
490 reading and language phenotypes

491 Supplemental Figure 1. Distribution of associated SNPs

492 Supplemental Figure 2. Functional annotation corresponding to Figure 3.

493 Supplemental Figure 3. Clustering of Significant Variants ( $P < 0.01$ ) among Known Speech  
494 Genes across CFSRS Tests

495 Supplemental Figure 4. LocusZoom plots of candidate genes where at least one trait had a SNP  
496 significant at  $p < 10^{-4}$

497 Supplemental Figure 5. Clustering of Significant Variants ( $P < 0.01$ ) among Known Speech  
498 Genes across ALSPAC Tests

499 Supplemental Figure 6. Polygenic Risk score across all individual measures

500

## 501 REFERENCES

502

- 503 1 *Almost 8 percent of US children have a communication or swallowing disorder*, 2015).
- 504 2 Catts, H. W., Adlof, S. M., Hogan, T. P. & Weismer, S. E. Are specific language impairment and  
505 dyslexia distinct disorders? *Journal of speech, language, and hearing research : JSLHR* **48**, 1378-  
506 1396, doi:10.1044/1092-4388(2005/096) (2005).
- 507 3 Shriberg, L., Tomblin, J. & McSweeney, J. Prevalence of speech delay in 6-year-old children and  
508 comorbidity with language impairment. *Journal of Speech, Language, and Hearing Research* **42**,  
509 1461-1481 (1999).
- 510 4 McLeod, S. B., E. in *Children's speech: an Evidence-based approach to assessment and*  
511 *intervention* (ed S. McLeod, Baker, E) 181-184 (Pearson Education, 2017).
- 512 5 Lemons, C. J. & Fuchs, D. Phonological awareness of children with Down syndrome: its role in  
513 learning to read and the effectiveness of related interventions. *Research in developmental*  
514 *disabilities* **31**, 316-330, doi:10.1016/j.ridd.2009.11.002 (2010).
- 515 6 Al Otaiba, S., Puranik, C., Zilkowski, R. & Curran, T. Effectiveness of Early Phonological Awareness  
516 Interventions for Students with Speech or Language Impairments. *The Journal of special*  
517 *education* **43**, 107-128, doi:10.1177/0022466908314869 (2009).

- 518 7 Larivee, L. C., HW. Early reading achievement in children with expressive phonological disorders.  
519 *Am J Speech Lang Pathol* **8**, 119-128 (1999).
- 520 8 Scarborough, H. in *Specific Reading Disabilities: A view of the spectrum* (ed BK; Accardo Shapiro,  
521 PJ; Capute, AJ) 75-119 (York Press, 1990).
- 522 9 Lewis, B. A. *et al.* The Genetic Bases of Speech Sound Disorders: Evidence From Spoken and  
523 Written Language. *J Speech Lang Hear. Res* **49**, 1294-1312 (2006).
- 524 10 Stein, C. M. *et al.* Pleiotropic effects of a chromosome 3 locus on speech-sound disorder and  
525 reading. *Am J Hum Genet* **74**, 283-297 (2004).
- 526 11 Lewis, B. A. *et al.* Heritability and longitudinal outcomes of spelling skills in individuals with  
527 histories of early speech and language disorders. *Learning and individual differences* **65**, 1-11,  
528 doi:10.1016/j.lindif.2018.05.001 (2018).
- 529 12 Stevenson, J., Graham, P., Fredman, G. & McLoughlin, V. A twin study of genetic influences on  
530 reading and spelling ability and disability. *Journal of child psychology and psychiatry, and allied*  
531 *disciplines* **28**, 229-247, doi:10.1111/j.1469-7610.1987.tb00207.x (1987).
- 532 13 Carrion-Castillo, A. *et al.* Evaluation of results from genome-wide studies of language and  
533 reading in a novel independent dataset. *Genes Brain Behav* **15**, 531-541, doi:10.1111/gbb.12299  
534 (2016).
- 535 14 Eicher, J. D. *et al.* Genome-wide association study of shared components of reading disability  
536 and language impairment. *Genes Brain Behav* **12**, 792-801, doi:10.1111/gbb.12085 (2013).
- 537 15 Gialluisi, A. *et al.* Genome-wide association scan identifies new variants associated with a  
538 cognitive predictor of dyslexia. *Translational psychiatry* **9**, 77, doi:10.1038/s41398-019-0402-0  
539 (2019).
- 540 16 Gialluisi, A. *et al.* Genome-wide screening for DNA variants associated with reading and  
541 language traits. *Genes Brain Behav* **13**, 686-701, doi:10.1111/gbb.12158 (2014).
- 542 17 Harlaar, N. *et al.* Genome-wide association study of receptive language ability of 12-year-olds. *J*  
543 *Speech Lang Hear Res* **57**, 96-105, doi:10.1044/1092-4388(2013/12-0303) (2014).
- 544 18 Kornilov, S. A. *et al.* Genome-Wide Association and Exome Sequencing Study of Language  
545 Disorder in an Isolated Population. *Pediatrics* **137**, doi:10.1542/peds.2015-2469 (2016).
- 546 19 Luciano, M. *et al.* A genome-wide association study for reading and language abilities in two  
547 population cohorts. *Genes Brain Behav* **12**, 645-652, doi:10.1111/gbb.12053 (2013).
- 548 20 St Pourcain, B. *et al.* Common variation near ROBO2 is associated with expressive vocabulary in  
549 infancy. *Nature communications* **5**, 4831, doi:10.1038/ncomms5831 (2014).
- 550 21 Morris, N., Elston, R. C., Barnholtz-Sloan, J. S. & Sun, X. Novel approaches to the analysis of  
551 family data in genetic epidemiology. *Front Genet* **6**, 27, doi:10.3389/fgene.2015.00027 (2015).
- 552 22 Ott, J., Kamatani, Y. & Lathrop, M. Family-based designs for genome-wide association studies.  
553 *Nat Rev Genet* **12**, 465-474, doi:10.1038/nrg2989 (2011).
- 554 23 Lewis, B. & Freebairn, L. Speech production skills of nuclear family members of children with  
555 phonology disorders. *Speech and Language* **41**, 45-61 (1998).
- 556 24 Lewis, B., Freebairn, L. & Taylor, H. Follow-up of children with early expressive phonology  
557 disorders. *Journal of Learning Disabilities* **33**, 433-444 (2000).
- 558 25 Lewis, B. A. *et al.* Literacy outcomes of children with early childhood speech sound disorders:  
559 impact of endophenotypes. *J Speech Lang Hear. Res* **54**, 1628-1643 (2011).
- 560 26 Lewis, B. A. *et al.* Family pedigrees of children with suspected childhood apraxia of speech.  
561 *Journal of Communication Disorders* **37**, 157-175 (2004).
- 562 27 Lewis, B. A., Freebairn, L. A., Hansen, A. J., Iyengar, S. K. & Taylor, H. G. School-age follow-up of  
563 children with childhood apraxia of speech. *Language, speech, and hearing services in schools* **35**,  
564 122-140 (2004).

- 565 28 Lewis, B. A. *et al.* Speech and language skills of parents of children with speech sound disorders.  
566 *Am J Speech Lang Pathol* **16**, 108-118 (2007).
- 567 29 Hollingshead, A. (Department of Sociology, Yale University, New Haven, CT. 06520, 1975).
- 568 30 Robbins, J. & Klee, T. Clinical assessment of oropharyngeal motor development in young  
569 children. *Journal of Speech and Hearing Research* **52**, 271-277 (1987).
- 570 31 Fletcher, D. (C.C. Publications, Inc., Tigard, OR, 1977).
- 571 32 Gardner, M. (Academic Therapy Publications, Novato, CA, 1990).
- 572 33 Dunn, L. & Dunn, L. (American Guidance Service, Inc, Circle Pines, MN, 1997).
- 573 34 Catts, H. Speech production/phonological deficits in reading disordered children. *Journal of*  
574 *Learning Disabilities* **19**, 504-508 (1986).
- 575 35 Denkla, M. & Rudel, R. Rapid 'automatized' naming (R.A.N.): dyslexia differentiated from other  
576 learning disabilities. *Neuropsychologia*, 471-479 (1976).
- 577 36 Wagner, R. T., J; Rashotte, C; Pearson, NA. (Pearson, London, England, 2013).
- 578 37 Wechsler, D. (The Psychological Corporation, San Antonio, TX, 1991).
- 579 38 Larsen, S. H., D. (The Psychological Corporation, San Antonio, TX, 1994).
- 580 39 Newcomer, P. & Hammill, D. *Test of language development - Primary, Second Edition.* (Pro-Ed.,  
581 1988).
- 582 40 E, S., Wiig, E. & Secord, W. *Clinical evaluation of language fundamentals-Revised.* (The  
583 Psychological Corporation, 1987).
- 584 41 GENESIS: GENetic ESTimation and Inference in Structured samples (GENESIS): Statistical methods  
585 for analyzing genetic data from samples with population structure and/or relatedness. R  
586 package version (2019).
- 587 42 Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide complex trait  
588 analysis. *Am. J. Hum Genet* **88**, 76-82 (2011).
- 589 43 Lee, S. H., Yang, J., Goddard, M. E., Visscher, P. M. & Wray, N. R. Estimation of pleiotropy  
590 between complex diseases using single-nucleotide polymorphism-derived genomic relationships  
591 and restricted maximum likelihood. *Bioinformatics (Oxford, England)* **28**, 2540-2542,  
592 doi:10.1093/bioinformatics/bts474 (2012).
- 593 44 Zhan, X., Hu, Y., Li, B., Abecasis, G. R. & Liu, D. J. RVTESTS: an efficient and comprehensive tool  
594 for rare variant association analysis using sequence data. *Bioinformatics (Oxford, England)* **32**,  
595 1423-1426, doi:10.1093/bioinformatics/btw079 (2016).
- 596 45 Svishcheva, G. R., Axenovich, T. I., Belonogova, N. M., van Duijn, C. M. & Aulchenko, Y. S. Rapid  
597 variance components-based method for whole-genome association analysis. *Nature genetics* **44**,  
598 1166-1170, doi:10.1038/ng.2410 (2012).
- 599 46 Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based linkage  
600 analyses. *Am J Hum Genet* **81**, 559-575 (2007).
- 601 47 Fraser, A. *et al.* Cohort Profile: the Avon Longitudinal Study of Parents and Children: ALSPAC  
602 mothers cohort. *International journal of epidemiology* **42**, 97-110, doi:10.1093/ije/dys066  
603 (2013).
- 604 48 Golding, J., Pembrey, M. & Jones, R. ALSPAC--the Avon Longitudinal Study of Parents and  
605 Children. I. Study methodology. *Paediatric and perinatal epidemiology* **15**, 74-87,  
606 doi:10.1046/j.1365-3016.2001.00325.x (2001).
- 607 49 Boyd, A. *et al.* Cohort Profile: the 'children of the 90s'--the index offspring of the Avon  
608 Longitudinal Study of Parents and Children. *International journal of epidemiology* **42**, 111-127,  
609 doi:10.1093/ije/dys064 (2013).
- 610 50 Watanabe, K., Taskesen, E., van Bochoven, A. & Posthuma, D. Functional mapping and  
611 annotation of genetic associations with FUMA. *Nature communications* **8**, 1826-1826,  
612 doi:10.1038/s41467-017-01261-5 (2017).

- 613 51 MacDermot, K. *et al.* Identification of FOXP2 truncation as a novel cause of developmental  
614 speech and language deficits. *Am J Hum Genet* **76**, 1074-1080 (2005).
- 615 52 Spiteri, E. *et al.* Identification of the transcriptional targets of FOXP2, a gene linked to speech  
616 and language, in developing human brain. *American journal of human genetics* **81**, 1144-1157,  
617 doi:10.1086/522237 (2007).
- 618 53 Vernes, S. C. *et al.* High-throughput analysis of promoter occupancy reveals direct neural targets  
619 of FOXP2, a gene mutated in speech and language disorders. *American journal of human*  
620 *genetics* **81**, 1232-1250, doi:10.1086/522238 (2007).
- 621 54 Zhang, W. *et al.* Integrative transcriptome imputation reveals tissue-specific and shared  
622 biological mechanisms mediating susceptibility to complex traits. *Nature communications* **10**,  
623 3834, doi:10.1038/s41467-019-11874-7 (2019).
- 624 55 Wang, D. *et al.* Comprehensive functional genomic resource and integrative model for the  
625 human brain. *Science* **362**, doi:10.1126/science.aat8464 (2018).
- 626 56 Goriounova, N. A. & Mansvelder, H. D. Genes, Cells and Brain Areas of Intelligence. *Front Hum*  
627 *Neurosci* **13**, 44-44, doi:10.3389/fnhum.2019.00044 (2019).
- 628 57 Zhang, Y. & Hu, W. NF $\kappa$ B signaling regulates embryonic and adult neurogenesis. *Front Biol*  
629 *(Beijing)* **7**, 10.1007/s11515-11012-11233-z, doi:10.1007/s11515-012-1233-z (2012).
- 630 58 Newbury, D. F. & Monaco, A. P. Genetic advances in the study of speech and language disorders.  
631 *Neuron* **68**, 309-320 (2010).
- 632 59 Anthoni, H. *et al.* The aromatase gene CYP19A1: several genetic and functional lines of evidence  
633 supporting a role in reading, speech and language. *Behav Genet* **42**, 509-527,  
634 doi:10.1007/s10519-012-9532-3 (2012).
- 635 60 Hannula-Jouppi, K. *et al.* The axon guidance receptor gene ROBO1 is a candidate gene for  
636 developmental dyslexia. *PLoS Genet* **1**, e50 (2005).
- 637 61 Ashbrook, D. G. *et al.* Born to Cry: A Genetic Dissection of Infant Vocalization. *Front Behav*  
638 *Neurosci* **12**, 250-250, doi:10.3389/fnbeh.2018.00250 (2018).
- 639 62 Zhao, B. *et al.* Genome-wide association analysis of 19,629 individuals identifies variants  
640 influencing regional brain volumes and refines their genetic co-architecture with cognitive and  
641 mental health traits. *Nature genetics* **51**, 1637-1644, doi:10.1038/s41588-019-0516-6 (2019).
- 642 63 Fan, Y. *et al.* De Novo Mutations of CCNK Cause a Syndromic Neurodevelopmental Disorder with  
643 Distinctive Facial Dysmorphism. *American journal of human genetics* **103**, 448-455,  
644 doi:10.1016/j.ajhg.2018.07.019 (2018).
- 645 64 Worthey, E. A. *et al.* Whole-exome sequencing supports genetic heterogeneity in childhood  
646 apraxia of speech. *Journal of neurodevelopmental disorders* **5**, 29, doi:10.1186/1866-1955-5-29  
647 (2013).
- 648 65 Lanzillotta, A. *et al.* NF- $\kappa$ B in Innate Neuroprotection and Age-Related Neurodegenerative  
649 Diseases. *Front Neurol* **6**, 98-98, doi:10.3389/fneur.2015.00098 (2015).
- 650 66 Gusev, A. *et al.* Transcriptome-wide association study of schizophrenia and chromatin activity  
651 yields mechanistic disease insights. *Nature genetics* **50**, 538-548, doi:10.1038/s41588-018-0092-  
652 1 (2018).
- 653 67 El-Ansary, A. & Al-Ayadhi, L. GABAergic/glutamatergic imbalance relative to excessive  
654 neuroinflammation in autism spectrum disorders. *J Neuroinflammation* **11**, 189-189,  
655 doi:10.1186/s12974-014-0189-0 (2014).
- 656 68 Nazmi, A. *et al.* Chronic neurodegeneration induces type I interferon synthesis via STING,  
657 shaping microglial phenotype and accelerating disease progression. *Glia* **67**, 1254-1276,  
658 doi:10.1002/glia.23592 (2019).

659 69 Okerlund, N. D. *et al.* Dact1 is a postsynaptic protein required for dendrite, spine, and excitatory  
660 synapse development in the mouse forebrain. *J Neurosci* **30**, 4362-4368,  
661 doi:10.1523/JNEUROSCI.0354-10.2010 (2010).

662 70 Le Guen, Y. *et al.* A DACT1 enhancer modulates brain asymmetric temporal regions involved in  
663 language processing. *bioRxiv*, 539189, doi:10.1101/539189 (2019).

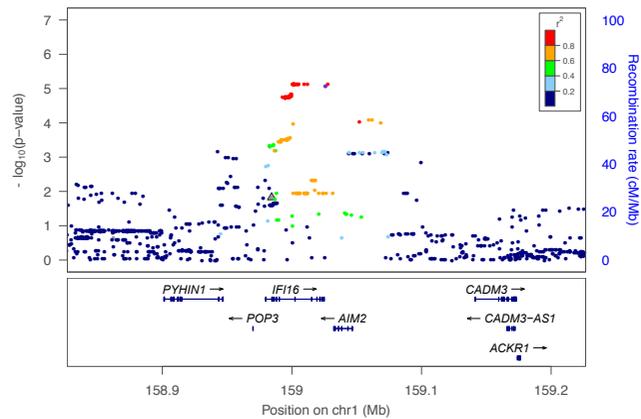
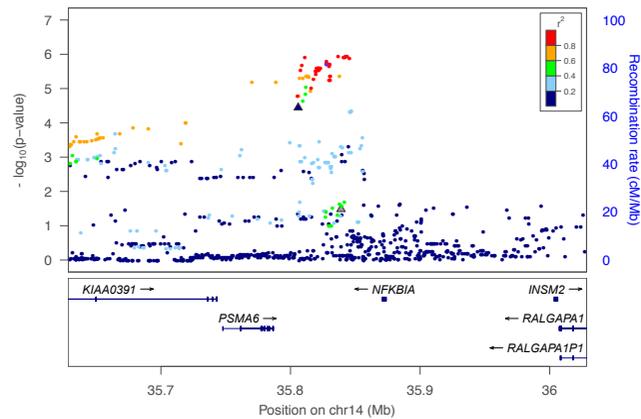
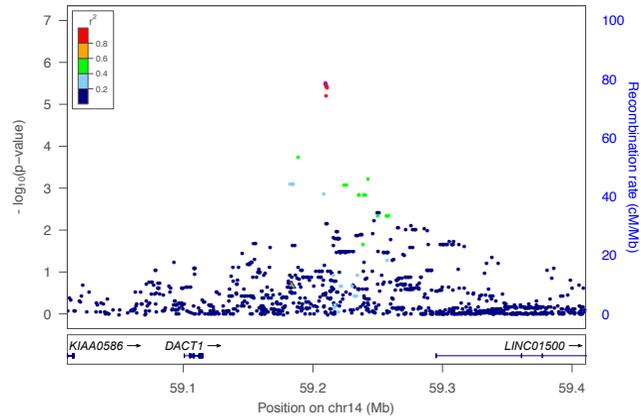
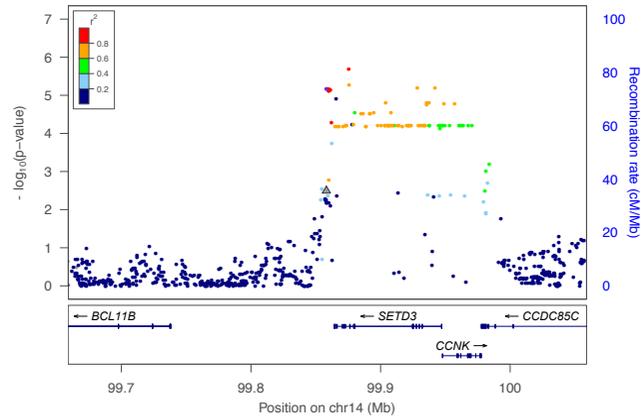
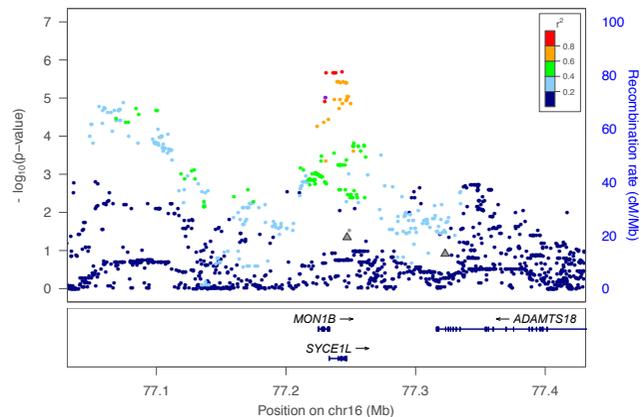
664 71 DiStasio, M. M., Nagakura, I., Nadler, M. J. & Anderson, M. P. T lymphocytes and cytotoxic  
665 astrocyte blebs correlate across autism brains. *Ann Neurol* **86**, 885-898, doi:10.1002/ana.25610  
666 (2019).

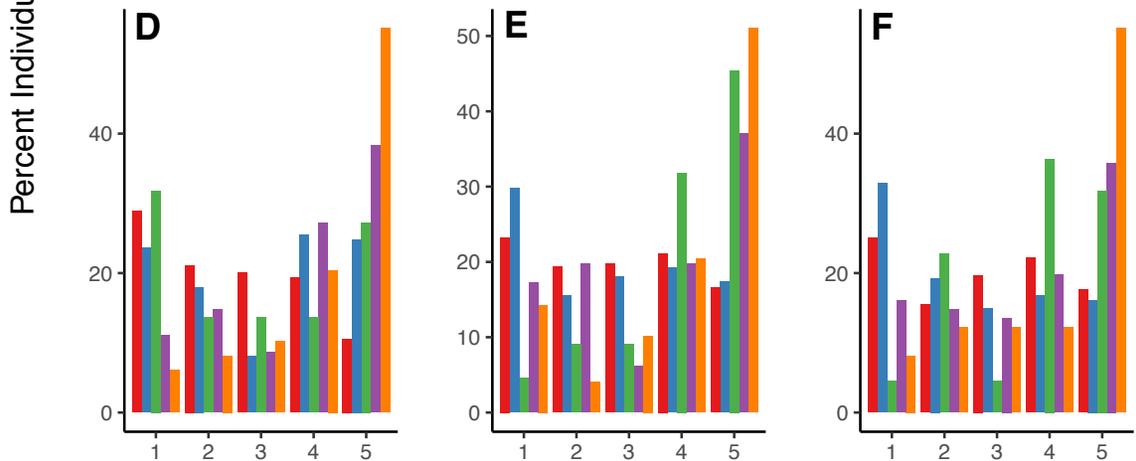
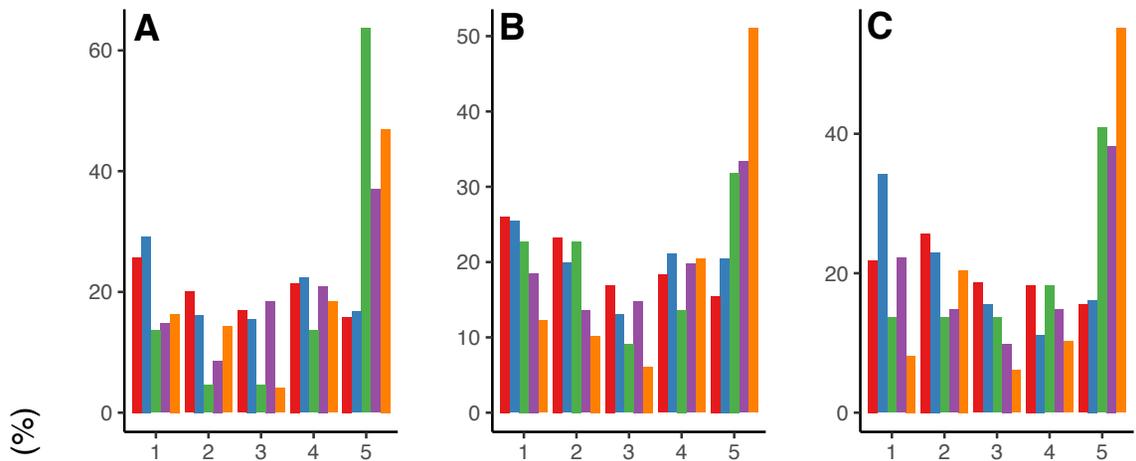
667 72 Mueller, K. L. *et al.* Common Genetic Variants in FOXP2 Are Not Associated with Individual  
668 Differences in Language Development. *PLoS One* **11**, e0152576,  
669 doi:10.1371/journal.pone.0152576 (2016).

670 73 Nudel, R. *et al.* Language deficits in specific language impairment, attention deficit/hyperactivity  
671 disorder, and autism spectrum disorder: An analysis of polygenic risk. *Autism research : official  
672 journal of the International Society for Autism Research*, doi:10.1002/aur.2211 (2019).

673



**A****B****C****D****E**



Subgroup ■ No SSD/LI/CAS ■ SSD only ■ LI only ■ SSD & LI ■ CAS

Genetic Risk Score (quantile)

