

# 1 Integrating *de novo* and inherited variants in over 42,607 autism cases 2 identifies mutations in new moderate risk genes

3 Xueya Zhou<sup>1,8,11</sup>, Pamela Feliciano<sup>2,11</sup>, Tianyun Wang<sup>3,11</sup>, Irina Astrovskaia<sup>2,11</sup>, Chang Shu<sup>1,8,11</sup>,  
4 Jacob B. Hall<sup>2</sup>, Joseph U. Obiajulu<sup>1,8</sup>, Jessica Wright<sup>2</sup>, Shwetha Murali<sup>3</sup>, Simon Xuming Xu<sup>2</sup>, Leo  
5 Brueggeman<sup>4</sup>, Taylor R. Thomas<sup>4</sup>, Olena Marchenko<sup>2</sup>, Christopher Fleisch<sup>2</sup>, Sarah D. Barns<sup>2</sup>,  
6 LeeAnne Green Snyder<sup>2</sup>, Bing Han<sup>2</sup>, Timothy S. Chang<sup>5</sup>, Tychele N. Turner<sup>6</sup>, William Harvey<sup>3</sup>,  
7 Andrew Nishida<sup>7</sup>, Brian J. O'Roak<sup>7</sup>, Daniel H. Geschwind<sup>5</sup>, The SPARK Consortium, Jacob J.  
8 Michaelson<sup>4</sup>, Natalia Volfovsky<sup>2</sup>, Evan E. Eichler<sup>3</sup>, Yufeng Shen<sup>8,9,12</sup> and Wendy K. Chung<sup>1,2,10,12</sup>

9 <sup>1</sup>Department of Pediatrics, Columbia University Medical Center, New York, NY, <sup>2</sup>Simons  
10 Foundation, New York, New York, <sup>3</sup>Department of Genome Sciences, University of Washington  
11 School of Medicine, Seattle, WA, <sup>4</sup>Department of Psychiatry, University of Iowa Carver College  
12 of Medicine, Iowa City, IA, <sup>5</sup>Program in Neurogenetics, Department of Neurology, David Geffen  
13 School of Medicine, University of California, Los Angeles, Los Angeles, CA, <sup>6</sup>Department of  
14 Genetics, Washington University, St. Louis, MO, <sup>7</sup>Department of Molecular & Medical Genetics,  
15 Oregon Health & Science University, Portland, OR, <sup>8</sup>Department of Systems Biology, Columbia  
16 University Medical Center, New York, NY, <sup>9</sup> Department of Biomedical Informatics, Columbia  
17 University Medical Center, New York, NY, <sup>10</sup> Department of Medicine, Columbia University  
18 Medical Center, New York, NY, <sup>11</sup>These authors contributed equally: Xueya Zhou, Pamela  
19 Feliciano, Tianyun Wang, Irina Astrovskaia and Chang Shu. <sup>12</sup>These authors jointly supervised this  
20 work: Yufeng Shen, Wendy K. Chung. A full list of the SPARK Consortium members appears at the end of this paper.  
21

## 22 Abstract

23 Despite the known heritable nature of autism spectrum disorder (ASD), studies have primarily  
24 identified risk genes with *de novo* variants (DNVs). To capture the full spectrum of ASD genetic  
25 risk, we performed a two-stage analysis of rare *de novo* and inherited coding variants in 42,607  
26 ASD cases, including 35,130 new cases recruited online by SPARK. In the first stage, we analyzed  
27 19,843 cases with one or both biological parents and found that known ASD or  
28 neurodevelopmental disorder (NDD) risk genes explain nearly 70% of the genetic burden  
29 conferred by DNVs. In contrast, less than 20% of genetic risk conferred by rare inherited loss-of-  
30 function (LoF) variants are explained by known ASD/NDD genes. We selected 404 genes based  
31 on the first stage of analysis and performed a meta-analysis with an additional 22,764 cases and  
32 236,000 population controls. We identified 60 genes with exome-wide significance ( $p < 2.5e-6$ ),  
33 including five new risk genes (*NAV3*, *ITSN1*, *MARK2*, *SCAF1*, and *HNRNPUL2*). The association of  
34 *NAV3* with ASD risk is entirely driven by rare inherited LoFs variants, with an average relative  
35 risk of 4, consistent with moderate effect. ASD individuals with LoF variants in the four  
36 moderate risk genes (*NAV3*, *ITSN1*, *SCAF1*, and *HNRNPUL2*,  $n = 95$ ) have less cognitive  
37 impairment compared to 129 ASD individuals with LoF variants in well-established, highly  
38 penetrant ASD risk genes (*CHD8*, *SCN2A*, *ADNP*, *FOXP1*, *SHANK3*) (59% vs. 88%,  $p = 1.9e-06$ ).  
39 These findings will guide future gene discovery efforts and suggest that much larger numbers of  
40 ASD cases and controls are needed to identify additional genes that confer moderate risk of  
41 ASD through rare, inherited variants.

42

### 43 **Introduction**

44 Autism spectrum disorder (ASD) is a neurodevelopmental condition characterized by impaired  
45 social communication and repetitive behaviors<sup>1</sup>. Previous studies in ASD utilized family-based  
46 designs to focus on *de novo* variants (DNVs) identified from parent-offspring trios<sup>2-8</sup>. Over one-  
47 hundred high confidence ASD genes enriched with likely deleterious DNVs have been  
48 identified<sup>8</sup>, most of which are also enriched for DNVs in other neurodevelopment disorders  
49 (NDDs)<sup>9-11</sup>. Statistical modeling suggests there are ~1000 genes with DNV variants in ASD<sup>12,13</sup>.  
50 However, despite the large effect size of individual pathogenic DNVs, all DNVs together only  
51 explain ~ 2% of variance in liability for ASD<sup>8,14</sup>.

52 On the other hand, ASD is highly heritable (estimated heritability over 0.5)<sup>14-16</sup>. Previous  
53 studies estimated that common variants explain up to half of the heritability<sup>14</sup>, although only  
54 five genome-wide significant loci have been identified<sup>17</sup>. The role of inherited coding variants  
55 has been evaluated using familial segregation of loss-of-function (LoF) variants (stop-gain, splice  
56 site and frameshift variants) carried by parents without ASD diagnoses or intellectual disability.  
57 Rare LoF variants only in genes intolerant of variation<sup>9,18</sup> are over-transmitted to probands  
58 compared with siblings without ASD<sup>7,8,19-22</sup>. However, identification of the individual risk genes  
59 enriched by such inherited variants has remained elusive.

60 We have created a large longitudinal research cohort, SPARK (SPARKForAutism.org<sup>23</sup>) to  
61 advance research on the genetic, behavioral, and clinical features associated with ASD. SPARK  
62 represents the largest ASD cohort in the world, with over 100,000 individuals with ASD  
63 enrolled.

64 Rare, LoF variants are enriched in developmental disorders including ASD<sup>22,24</sup>, but LoF variants  
65 in the general population are also enriched for sequencing and annotation artefacts<sup>25</sup>, which  
66 present technical challenges in large sequencing studies. Methods to distinguish between high  
67 and low confidence LoF variants<sup>18,26,27</sup> have been used to quantify gene level LoF  
68 intolerance<sup>18,26,28,29</sup> and to refine the role of *de novo* LoF variants in NDDs<sup>20</sup>.

69 Here we present an integrated analysis of *de novo* and inherited coding variants in over 42,607  
70 ASD cases, including cases from previously published ASD cohorts and 35,130 new cases from  
71 SPARK. To our knowledge, this analysis is the largest sequencing study of ASD to date. In our  
72 two-stage design, we first characterized the contribution of DNVs and rare inherited LoF  
73 variants to ASD risk. Results from the first stage informed the second stage, in which we  
74 conducted a meta-analysis of 404 genes. By combining evidence from DNVs, over-transmission,  
75 and case-control comparison, we identified 60 ASD risk genes with exome-wide significance,  
76 including five new genes not previously implicated in neurodevelopmental conditions. Finally,  
77 we estimated the effect sizes of known and newly significant genes and used them for power  
78 calculations to inform the design of future studies.

79

## 80 Results

### 81 *Overview of data and workflow*

82 We aggregated exome or whole genome sequencing (WGS) data of 35,130 new cases from the  
83 SPARK study and 7,665 cases from published ASD studies (ASC<sup>3,8</sup>, MSSNG<sup>6</sup>, and SSC<sup>2,30</sup>)  
84 (**Supplementary Table S1**) and performed a two-stage analysis (**Figure 1**). In the first stage, we  
85 analyzed *de novo* coding variants (DNVs) in 16,877 ASD trios and assessed transmission of rare  
86 LoF variants in 20,491 parents without ASD diagnoses or intellectual disability to offspring with  
87 ASD (including 9,504 trios and 2,966 single-parent-proband duos). For DNVs, we characterized  
88 the enrichment pattern in known and candidate risk genes, mutation intolerance (ExAC pLI<sup>18</sup>  
89 and gnomAD metrics<sup>26</sup>) and performed gene-based burden tests of *de novo* LoF and missense  
90 variants by DeNovoWest<sup>11</sup>. For rare inherited LoFs, we estimated the over-transmission from  
91 parents without an ASD diagnosis to ASD cases in all genes and gene sets predefined by  
92 functional genomic data or results from DNV analysis. Based on DNV enrichment and over-  
93 transmission patterns in gene sets, we selected 404 genes for meta-analysis in stage 2 utilizing  
94 22,764 new cases with exome or WGS data. In stage 2, we applied DeNovoWEST on DNVs,  
95 conducted transmission-disequilibrium tests on inherited LoFs in trios or duos, performed  
96 burden tests on rare LoFs in cases compared with population controls (104,068 subjects from  
97 gnomAD exome, non-neuro subset v2.1.1 and 132,345 TOPMed subjects), and combined the p-  
98 values to estimate a final p-value for each of the 404 genes. Finally, we performed a mega-  
99 analysis of rare LoFs in all cases and controls to estimate the effect sizes of known or new  
100 candidate ASD genes to inform future studies.

### 101 *Known ASD or NDD risk genes explain two-thirds of population attributable risk of de novo* 102 *coding variants in ASD*

103 In the first stage, we combined data from four large-scale ASD cohorts, resulting in 16,877  
104 unique ASD trios and 5,764 unaffected trios (**Supplementary Table S1**). The four cohorts show  
105 similar exome-wide burden of DNVs in simplex families. The burden of *de novo* LoF variants in  
106 cases with a family history of ASD is significantly lower than those without a reported family  
107 history ( $p=1.1e-4$  by Poisson test), whereas the burden of predicted *de novo* damaging  
108 missense (D-mis, defined by REVEL score<sup>31</sup> $\geq 0.5$ ) and synonymous variants are similar  
109 (**Supplementary Figure S1**).

110 Compared to unaffected offspring, the excess of damaging DNVs (*de novo* LoF and D-mis  
111 variants) in individuals with ASD is concentrated in LoF-intolerant genes, defined as genes with  
112 a probability of being LoF intolerant ( $pLI$ )<sup>18</sup>  $\geq 0.5$  in the Exome Aggregation Consortium (ExAC).  
113 Using LoF observed/expected upper-bound fraction (LOEUF), a recently developed gene  
114 constraint metric<sup>26</sup>, the burden of damaging DNVs is highest among genes ranked in the top  
115 20% of LOEUF scores (**Figure 2A**). Overall, the population attributable risk (PAR) from damaging  
116 DNVs is about 10%. We assembled 618 previously established dominant (“known”) ASD or NDD  
117 risk genes (**Supplementary Table S2**). These genes explained about 2/3 of the PAR from  
118 damaging DNVs. Excluding these genes, the fold enrichment of damaging DNVs was greatly  
119 attenuated (**Figure 2A**).

120 To assess the evidence of DNVs in individual genes, we applied DeNovoWEST<sup>11</sup>, which  
121 integrates DNV enrichment with clustering of missense variants in each gene. The initial  
122 DeNovoWEST scan of DNVs in 16,877 ASD trios identified 159 genes with  $p < 0.001$   
123 (**Supplementary Table S3**).

#### 124 *Rare inherited LoF variants contribute to ASD risk mostly through unknown risk genes*

125 To analyze the contribution of rare inherited LoF variants to ASD risk, we evaluated  
126 transmission disequilibrium in ultra-rare (allele frequency  $< 1e-5$ ) high-confidence (by LOFTEE<sup>26</sup>  
127 and pext<sup>27</sup>; see Methods and Supplementary Note) LoF variants from parents without ASD  
128 diagnoses or intellectual disability to affected offspring with ASD in 9,504 trios and 2,966 duos  
129 from the first stage (**Supplementary Table S4**). For a given set of genes, we quantified  
130 transmission disequilibrium using the number of over-transmitted (excess in transmission over  
131 non-transmission) LoF variants per trio; parent-offspring duos were considered half-trios.

132 Among autosomal genes, the overall transmission disequilibrium signal of ultra-rare LoF  
133 variants is enriched in LoF intolerant genes (ExAC  $pLI \geq 0.5$ ) and in genes within the top 20% of  
134 LOEUF scores (**Figure 2B**), similar to the burden of damaging DNVs. We observed both over-  
135 transmission to affected and under-transmission to unaffected offspring, especially in genes  
136 within the top 10% of LOEUF scores. However, known ASD/NDD genes only explain ~20% of  
137 over-transmission of LoF variants to affected offspring (**Figure 2B**). On the X chromosome, we  
138 only considered transmission from mothers without ASD diagnoses to 9,883 affected sons and  
139 2,571 affected daughters (**Supplementary Table S4**). Rare LoF variants in mothers without ASD  
140 diagnoses only show significant over-transmission to affected sons but not affected daughters  
141 and remain significant after removing known ASD/NDD genes (**Supplementary Figure S2**).  
142 Together, these data suggest that most genes conferring inherited ASD risk are yet to be  
143 identified. Autosomal rare D-mis variants also show evidence of transmission disequilibrium to  
144 affected offspring, although the signal is much weaker and dependent on gene set, D-mis  
145 prediction method, pExt and allele frequency filters (**Supplementary Figure S3**).

146 To characterize the properties of genes contributing to ASD risk through rare inherited variants,  
147 we defined 25 gene sets from five categories representing both functional and genetic evidence  
148 relevant to ASD (**Supplementary Table S5 and Supplementary Figure S4**). We limited the genes  
149 to 5,754 autosomal constrained genes (ExAC  $pLI \geq 0.5$  or top 20% of LOEUF scores) and  
150 performed TDT (**Supplementary Table S6**). For each gene set, we tested if high-confidence rare  
151 LoF variants show a higher frequency of transmission to ASD offspring than the remaining  
152 genes in the overall constrained gene set. As a comparison with DNVs, we also tested if the  
153 same set of genes are more frequently disrupted by damaging DNVs than the rest of the genes  
154 in ASD trios using the framework of dnEnrich<sup>32</sup>.

155 We first considered functional gene sets derived from the neuronal transcriptome, proteome,  
156 or regulome. We confirmed significant enrichment in damaging DNVs ( $p < 0.005$  by simulation)  
157 in the gene sets that were previously suggested to be enriched for ASD risk genes including  
158 expression module M2/3<sup>33</sup>, RBFOX1/3 targets<sup>34</sup>, FMRP targets<sup>35</sup>, and CHD8 targets<sup>36</sup>. However,  
159 this enrichment can be largely explained by known ASD/NDD genes (**Supplementary Figure S5**).  
160 For ultra-rare inherited LoF variants, we found the proportion of transmission to ASD  
161 individuals in most functional gene sets is close to all genes in the background; only RBFOX

162 targets show a weak enrichment but can be largely explained by known genes (**Figure 3**). We  
163 also applied two recently developed machine learning methods to prioritize ASD risk genes:  
164 forecASD<sup>37</sup> that integrates brain expression, gene network, and other gene level metrics, and A-  
165 risk<sup>38</sup> that uses cell-type specific expression signatures in developing brain. Although  
166 enrichment of DNVs in genes predicted by these methods are mainly explained by known  
167 genes, genes prioritized by A-risk are significantly enriched with inherited LoFs that cannot be  
168 explained by known genes. Using A-risk $\geq 0.4$  (recommended threshold), 30% of constrained  
169 genes (n=1,464) were prioritized and explain 64% of the over-transmission of LoF variants to  
170 ASD offspring ( $p=2.6e-5$  by chi-squared test). The enrichment is even higher than genes  
171 prioritized by the LOEUF score: 33% of genes (N=1,777) in the top decile of LOEUF account for  
172 55% over-transmission ( $P=3.5e-4$  by chi-squared test) (**Figure 3**).

173 We also considered gene sets that have evidence of genetic association with DNVs. Genes  
174 nominally enriched by DNVs ( $P<0.01$  by DeNovoWEST; N=300) in ASD from the current study  
175 have a significantly higher over-transmission rate than other constrained genes (Odds  
176 ratio=1.39,  $p=3.0e-5$  by chi-squared test) (**Figure 3**), although these genes only account for 21%  
177 of the over-transmission. Genes nominally enriched by DNVs in other NDDs<sup>11</sup> are also  
178 significantly enriched by DNVs in ASD and weakly enriched by inherited LoFs in ASD; however,  
179 both can be largely explained by known genes (**Figure 3**). This suggests that a subset of ASD  
180 genes increase risk by both *de novo* and inherited variants, and new genes can be identified by  
181 integrating evidence from DNV enrichment and TDT.

#### 182 *DNVs and a subset of rare inherited LoFs are associated with cognitive impairment*

183 To evaluate the association of genotypes with phenotype in ASD, we used self-reported  
184 cognitive impairment in SPARK, a Vineland score of  $<70$  in the SSC or the presence of  
185 intellectual disability in ASC. Damaging DNVs in genes ranked within the top 10% of LOEUF  
186 scores show a higher burden ( $p=1.1e-24$ , by chi-squared test) in ASD cases with evidence of  
187 cognitive impairment than other cases, consistent with previous results<sup>2,8</sup> (**Figure 4A**). Once  
188 known ASD/NDD genes were excluded, the residual burden of damaging DNVs in genes at the  
189 top 10% LOEUF is greatly reduced and not significantly associated with cognitive phenotype in  
190 ASD (**Figure 4A**). Over-transmission of rare LOFs in genes within the top 10% of LOEUF genes to  
191 ASD cases with cognitive impairment is about 2.7 times higher than to the cases without  
192 cognitive impairment ( $p=4.6e-3$  by chi-squared test) and is still 2x higher ( $p=0.04$  by chi-squared  
193 test) once known ASD/NDD genes were excluded (**Figure 4B**). However, rare LoFs in genes  
194 prioritized by A-risk, in which there is significant over-transmission to all cases overall, are not  
195 associated with cognitive impairment (**Supplementary Figure S6**). Taken together, these results  
196 suggest that rare variants in the top 10% of LOEUF genes—most of which are already known to  
197 be ASD/NDD risk genes—are associated with cognitive impairment. However, a subset of rare,  
198 inherited variants, particularly those prioritized by A-risk, are not associated with cognitive  
199 impairment.

200

201 *Meta-analysis of de novo and rare inherited LoF variants identifies 5 new risk genes with exome-*  
202 *wide significance*

203 Based on results from the first stage of analysis, 404 genes showed plausible evidence of  
204 contributing to ASD risk, including: 1) 260 genes with evidence of TDT (TDT statistic<sup>39</sup> $\geq 1$ ) and in  
205 gene sets enriched with rare inherited LoFs (top 10% LOEUF or within top 20% LOEUF and A-  
206 risk $\geq 0.4$ ) (**Supplementary Table S6**) and 2) 159 genes with  $p < 0.001$  from the DeNovoWEST  
207 analysis of DNVs (with 15 genes by both) (**Supplementary Table S3**). We performed a meta-  
208 analysis on the 367 autosomal genes with all data from Stage 1 and Stage 2, which includes  
209 6,174 new ASD trios, 1,942 new duos, 15,780 unrelated cases (see Methods), and 236,000  
210 population controls.

211 In the meta-analysis, we used Fisher's method<sup>40</sup> to combine 3 p-values that estimate  
212 independent evidence of DNVs, TDT, and case-control comparison: (1) DeNovoWEST with DNVs  
213 from both Stage 1 and 2 ( $n=23,039$  trios, **Supplementary Table S1**) using the parameters  
214 estimated in Stage 1, (2) TDT with rare LoF variants in parents without ASD diagnoses or  
215 intellectual disability with affected offspring in 15,586 trios and 4,907 duos (**Supplementary**  
216 **Table S4**), and (3) unrelated cases (**Supplementary Table S7**) compared to population controls  
217 using a binomial test. We used two sets of controls: gnomAD exome v2.1.1 non-neuro subset  
218 ( $n=104,068$ ) and TOPMed WGS (freeze 8,  $n=132,345$ ). We performed a case-control burden test  
219 using the two sets separately and input the larger p-value for the Fisher's method. This  
220 approach avoids any sample overlap and provides sensitivity analysis to ensure that significant  
221 genes are not dependent on the choice of population reference. Although population reference  
222 data were processed by different bioinformatics pipelines, the cumulative allele frequencies  
223 (CAFs) of high-confidence (HC, see Methods) LoF variants are similar between internal pseudo-  
224 controls (see Methods) and the two population references after applying the same LoF filters  
225 (**Supplementary Figure S7**). Previous population genetic simulations predict that for genes  
226 under moderate to strong selection (selection coefficient $>0.001$ ), deleterious variants are  
227 expected to arise within 1,000 generations and population demographic histories do not  
228 confound the CAFs of deleterious alleles in these genes<sup>41</sup>. For 367 selected autosomal genes,  
229 the point estimates of selection coefficient under mutation-selection balance model<sup>42</sup> are all  
230 greater than 0.01 (**Supplementary Figure S8**). Consistent with the theoretical predictions, most  
231 HC LoF variants in these genes are ultra-rare (**Supplementary Figure S9**) and the CAFs of HC LoF  
232 variants in European and non-European population samples are highly correlated  
233 (**Supplementary Figure S10**). Thus, we included population samples across all ancestries as  
234 controls. To make use of all genetic data collected, we also included rare variants of unknown  
235 inheritance from autism cases that were analyzed in the first stage. These variants come from  
236 cases that are part of parent-autism duos; such variants were either inherited from the parent  
237 not participating in the study or occurred *de novo*. Therefore, these data represent data  
238 independent of the transmission disequilibrium testing, even though the same cases were  
239 included in TDT.

240 We identified 60 genes with exome-wide significance ( $p < 2.5e-6$ ). Figure 5 summarizes the  
241 distribution of LoF variants (with different modes of inheritance) in genes that reached  
242 experimental-wide significance by DNV enrichment (**Figure 5A**) and other significant genes by  
243 meta-analysis (**Figure 5B, Supplementary Figure S11**). Genes that are significant only in meta-

244 analysis tend to harbor more inherited LoF variants than *de novo* variants, consistent with their  
245 lower penetrance for ASD or NDD.

246 Although most significant genes were previously known, we identified five new genes that are  
247 exome-wide significant regardless of the choice of population reference: *NAV3*, *MARK2*, *ITSN1*,  
248 *SCAF1*, and *HNRNPUL2* (**Table 1**). As expected, most supporting variants are ultra-rare, and  
249 results are robust to the allele frequency filter. These five new genes together explain 0.27%  
250 population attributable risk ratio (PAR) (**Supplementary Table S8**). *NAV3* has a similar PAR as  
251 *CHD8* and *SCN2A* (~0.095%). *ITSN1* is similar to *PTEN* (~0.065%).

252 The association of *NAV3* with ASD risk is entirely driven by rare inherited variants (**Table 1**).  
253 *NAV3* harbors a single HC *de novo* LoF variant in an unaffected sibling in the SSC and was  
254 previously included in the negative training set by A-risk<sup>38</sup>. Despite this, *NAV3* still has a high A-  
255 risk score, suggesting *NAV3*'s expression pattern is highly similar to known ASD genes  
256 (**Supplementary Data 1**)<sup>7,43</sup>. *NAV3* has high expression in inner cortical plate of developing  
257 cortex<sup>33</sup>, and in pyramidal neurons (hippocampus CA1 and somatosensory cortex) and cortical  
258 interneurons, consistent with the signatures of known ASD genes<sup>44</sup> (**Supplementary Figure**  
259 **S12**).

260 The association of *MARK2* with ASD risk is primarily driven by DNVs. *MARK2* is also associated  
261 with other NDDs<sup>11</sup> ( $P=2.7e-5$  by DeNovoWEST) including Tourette syndrome<sup>45</sup> and epilepsy<sup>46</sup>.  
262 We find that 3/8 of autistic offspring with variants in *MARK2* report epilepsy, 2/8 report  
263 Tourette syndrome and 7/8 have evidence of cognitive impairment (**Supplementary Table S9**).

264 The remaining three novel genes have support from both DNVs and rare LoFs. Two genes have  
265 suggestive evidence from other NDD studies. *ITSN1* and *SCAF1* shows nominal significance of  
266 DNV enrichment in 31,058 NDD trios<sup>11</sup> ( $P<0.05$  by DeNovoWEST). *SCAF1* was among the top 50  
267 genes from gene-based burden test in a recent schizophrenia case-control study ( $P=0.0027$  by  
268 burden test)<sup>47</sup>. Both *ITSN1* and *NAV3* have moderate effect sizes (point estimate of relative risk  
269 3~6, **Supplementary Table S8**). *ITSN1* has been highlighted in our previous study with evidence  
270 of enriched inherited LoFs<sup>7</sup>. *ITSN1* and *NAV3* also show increased CAF of LoF variants in a recent  
271 study by ASC<sup>8</sup> although the association was not significant. We also assessed deletions in these  
272 new genes. For both *ITSN1* and *NAV3*, we identified four partial or whole gene deletions in  
273 33,083 parents without ASD diagnoses or intellectual disability that also show transmission  
274 disequilibrium to affected offspring (**Supplementary Figure S13**).

275  
276 While both *de novo* and rare inherited LoFs in the most constrained genes are strongly  
277 associated with intellectual disability (ID) in ASD (**Figure 4**), the association of such variants in  
278 individual genes is heterogenous, as suggested by the lack of association of rare inherited  
279 variants in genes with high A-risk (**Supplementary Figure S5**). We calculated the burden of  
280 cognitive impairment (see **Methods**) in 87 ASD individuals with HC LoF variants in the four novel  
281 moderate risk genes and compared it to 129 individuals with HC LoF in the well-established ASD  
282 risk genes *CHD8*, *SCN2A*, *SHANK3*, *ADNP* and *FOXP1* as well as 8,731 individuals with ASD in  
283 SPARK (**Supplementary Figure S14**). Although most individuals with variants in well-established  
284 ASD risk genes have some evidence of cognitive impairment (88%,) individuals with LoF variants  
285 in the moderate risk genes had significantly lower burden (56%,  $p=4.5e-7$  by chi-squared test).  
286 Individuals with HC LOFs in the moderate risk genes did not have a significantly different

287 burden of cognitive impairment than 8,731 individuals with ASD in SPARK (56% vs. 50%,  $p =$   
288 n.s.). Individuals with LoF variants in the moderate risk genes also had a similar male: female  
289 (4:1) ratio compared to the larger cohort whereas individuals with variants in the well-  
290 established ASD risk genes showed significantly less male bias (1.6: 1,  $p = 0.009$  by chi-squared  
291 test) (**Supplementary Figure S14**), as previously reported<sup>2</sup>. We also predicted full-scale IQ on all  
292 participants based on parent-reported data using a machine learning method<sup>48</sup>. Carriers of rare  
293 LoFs in three (*NAV3*, *SCAF1*, and *HNRNPUL2*) of the four new genes with substantial  
294 contribution from rare inherited variants have similar IQ distribution as the overall SPARK  
295 cohort (**Figure 6A**), which is substantially higher than heterozygotes with rare LoFs in well-  
296 established, highly-penetrant genes that contribute to ASD primarily through *de novo* variants  
297 (“DN genes”), such as *CHD8*, *SHANK3*, and *SCN2A*. In fact, both novel and established genes  
298 with significant contribution from rare inherited LoFs are less associated with ID than DN genes  
299 (**Figure 6B**). Across these genes, there is a significant negative correlation ( $r = 0.78$ ,  $p = 0.001$ ) of  
300 estimated relative risk of rare LoFs with average predicted IQ of the individuals with these  
301 variants (**Figure 6C**). These genes could be associated with other neurobehavioral phenotypes.

302

303 Most known ASD/NDD genes that are enriched by *de novo* LoF variant harbor more *de novo*  
304 than inherited HC LoF variants in ~16,000 unrelated ASD trios (**Figure 5A and Supplementary**  
305 **Figure S15**), consistent with their high penetrance for ASD/NDD phenotypes and strong  
306 negative selection. Using population exome or WGS data, we calculated a point estimate of  
307 selection coefficient ( $\hat{s}$ )<sup>49</sup> of LoFs in each gene (**Supplementary Table S8**) and found that the  
308 fraction of *de novo* LoFs in ASD genes is higher in genes with large  $\hat{s}$ , and smaller in genes with  
309 small  $\hat{s}$  (**Supplementary Figure S7B**), consistent with population genetic theory<sup>50</sup>. We also  
310 estimated average effect size of rare LoFs in ASD genes by comparing cumulative allele  
311 frequency (CAF) in 31,976 unrelated cases and population exome or WGS data. As expected,  
312 known and newly significant ASD genes with higher risk to ASD are under stronger selection  
313 (larger  $\hat{s}$ ) (**Supplementary Figure S16**).

#### 314 *Functional similarity of new genes with known ASD genes*

315 To better appreciate the probable functional implications of the new exome-wide significant genes that  
316 confer inherited risk for ASD, we integrated mechanistic (STRING<sup>102</sup>) and phenotypic (HPO<sup>103</sup>) data  
317 into a single embedding space (six dimensions, one for each archetype coefficient) using a  
318 combination of canonical correlation analysis and archetypal analysis. This embedding space  
319 serves as an interpretive framework for putative ASD risk genes ( $N = 1,776$ ). Six  
320 functional/phenotypic archetypes were identified (**Figure 7**) that represent pathways that are  
321 well-understood to play a role in ASD: neurotransmission (archetype 1 or A1), chromatin  
322 modification (archetype 2 or A2), RNA processing (archetype 3 or A3), membrane trafficking  
323 and protein transport (archetype 4 or A4), extracellular matrix, motility, and response to signal  
324 (archetype 5 or A5), and KRAB domain and leucine-rich region proteins (archetype 6 or A6), also  
325 enriched for intermediate filaments. These archetypes organize risk genes in a way that jointly  
326 maximizes their association with mechanisms (STRING clusters) and phenotypes (HPO terms).  
327 For instance, A1 genes (neurotransmission) are enriched for the STRING cluster CL:8435 (ion  
328 channel and neuronal system) and are also associated with seizure and epileptic phenotypes. A2  
329 genes (chromatin modifiers) are enriched for nuclear factors and genes linked to growth and

330 morphological phenotypes (**Supplementary Table S10**). We call genes that strongly map to an  
331 archetype (i.e., > 2x the next highest-ranking archetype) “archetypal” and “mixed” if this  
332 criterion is not met (see methods). Archetypal genes are generally less functionally ambiguous  
333 than “mixed” genes. Of the five novel inherited risk genes, two are archetypal (suggesting  
334 function within known risk mechanisms): *NAV3* (A6: KRAB domain & LRR) and *ITSN1* (A4:  
335 membrane trafficking and protein transport). *SCAF1*, *MARK2*, and *HNRNPUL2* are mixtures of  
336 the identified archetypes, largely A4 and A5. That these new genes did not resolve clearly into  
337 archetypes (that were defined by known and suspected autism risk genes) suggests that they  
338 may operate in potentially novel or under-appreciated mechanisms. To elucidate these  
339 possibilities, we constructed an *ad hoc* “archetype,” defined by the centroid between *SCAF1*,  
340 *MARK2*, and *HNRNPUL2* (see Figure 7C). Cell-cell junction (CL:6549) was the STRING cluster  
341 most associated with this centroid ( $p = 4.12 \times 10^{-14}$  by the K-S test, Fig. 7D), which fits with its  
342 location between A4 (membrane trafficking) and A5 (ECM).

#### 343 *Power analysis*

344 The power of identifying risk genes with rare or *de novo* variants monotonically increases with  
345 increasing effect size or expected CAF under the null. New ASD genes to be discovered are likely  
346 to have smaller effect size than known ASD genes, as suggested by our results. Additionally,  
347 known ASD genes are biased toward longer genes with higher background mutation rate of  
348 damaging variants (“long genes”) (**Supplementary Figure S17**). Even though longer genes are  
349 more likely to be expressed in brain and relevant to ASD/NDD<sup>51</sup>, among most constrained  
350 genes, long genes (LoF mutation rate<sup>52,53</sup> above 80% quantile) and short genes (below 80%)  
351 have similar enrichment of damaging *de novo* variants and rare inherited LoFs (**Supplementary**  
352 **Figure S18**). Notably, for small genes, known genes have virtually no contribution to over-  
353 transmitted HC LoFs to affected offspring (**Supplementary Figure S18B**). It suggests that many  
354 smaller genes contributing to ASD risk remain to be identified. We focus on the power of  
355 detecting new ASD genes with a moderate effect size and the full range of background  
356 mutation rate.

357 We use a published framework<sup>41</sup> to analyze power based on case-control association of rare  
358 variants. For rare variants in genes under strong selection, CAF is largely determined by  
359 mutation rate and selection coefficient<sup>41</sup>. We therefore modeled power of discovering risk  
360 genes as a function of relative risk and selection coefficient. With about 5,500 constrained  
361 genes, the power of the current study was calculated for 31,976 unrelated cases and  
362 experiment-wise error rate of  $9e-6$  (**Supplementary Figure S19**).

363 We inversed the power calculation to determine required sample size to achieve 90% power  
364 under the same assumptions (**Supplementary Figure S20**). For genes at median LoF mutation  
365 rate across all genes, we estimated that it requires about 96,000 cases (three times the current  
366 sample size) to identify genes with similar effect size as *NAV3* (RR=4.5) and *ITSN1* (RR=5), about  
367 64,000 (twice the current sample size) to find genes with similar effect sizes as *SCAF1* (RR=8)  
368 and *HNRNPUL2* (RR=9). We note that it requires 10 and 5 times the current sample size to  
369 detect these types of genes by *de novo* variants alone.

## 370 Discussion

371 In this study, we assembled the largest sequencing data set of individuals with ASD to date,  
372 including 35,130 ASD cases and their family members collected by SPARK. We characterized the  
373 contribution of rare inherited variants to ASD risk and identified five new ASD risk genes by  
374 both *de novo* and rare inherited coding variants. We identified rare LoF variants in new ASD risk  
375 genes with modest effect size that are not strongly associated with ID. This finding represents a  
376 difference in phenotypic association with ID compared with other well-established, highly  
377 penetrant ASD genes. To find new risk genes with relative risks of 2-5 (comparable to the low  
378 relative risk genes from this study: *NAV3* and *ITSN1*) in the 50-percentile for gene-wide LoF  
379 mutation rate ( $2e-6$ ) and the 50-percentile for selection among known risk genes (0.2), our  
380 power analysis suggests that 52,000, 73,000, 116,000 or 227,000 total ASD cases are necessary,  
381 respectively (cf. eq 1 from power calculation in Supplementary material). Larger ASD cohorts  
382 with phenotypic data will be necessary to identify new ASD risk genes and may help to  
383 understand the biology of core symptoms of ASD in individuals without ID.

384 Our results suggest that identification of new risk genes with rare inherited variants can  
385 substantially improve genetic diagnostic yield. We found that rare inherited LoF variants  
386 account for 6% of PAR, similar to *de novo* LoF variants. Over two thirds of the PAR from *de novo*  
387 coding variants are explained by known ASD or NDD genes. In contrast, less than 20% of PAR  
388 from rare inherited LoFs variants is explained by known genes, suggesting most genes  
389 contributing to ASD risk through rare inherited variants are yet to be discovered. These  
390 unknown risk genes are still largely constrained to LoFs in the general population and/or have  
391 similar expression profiles in developing brains to known ASD risk genes. Combining evidence  
392 from both *de novo* and rare inherited variants, we identified 60 genes associated with ASD with  
393 exome-wide significance, including five novel genes. Rare LoFs in these five new genes account  
394 for a PAR of 0.27%, about half of the PAR of the 5 most common highly penetrant ASD genes  
395 (*KDM5B*, *GIGYF1*, *CHD8*, *SCN2A*, *SHANK3*).

396 *NAV3*, to our knowledge, is the first autosomal ASD risk gene discovered by association of solely  
397 rare inherited variants. Carriers of rare LoFs in *NAV3* have an average predicted IQ of 81,  
398 slightly above the SPARK cohort average (79). The prevalence of ID among *NAV3* heterozygotes  
399 is similar to the SPARK cohort average. This is distinctly different from established ASD risk  
400 genes (e.g., *CHD8*, *SHANK3*, *SCN2A*), nearly all identified by highly penetrant *de novo* variants,  
401 associated with ID in ASD cohorts<sup>2</sup>. The absence of ID is also observed in other genes (e.g.,  
402 *SCAF1*, *HNRNPUL2*, *GIGYF1*, *KDM5B*, *KMT2C*) with substantial contribution from rare inherited  
403 variants and modest effect size. Nevertheless, the data show that variants in these new ASD  
404 genes have effects on core symptoms of ASD, cognition, and other behaviors including  
405 schizophrenia, Tourette syndrome, ADHD and other behavioral conditions. Detailed  
406 phenotyping of individuals carrying these rare inherited variants is needed to understand the  
407 phenotypic effects of each gene. Such strategies should include a genetic and phenotypic  
408 assessment of family members who also carry the rare variant but may not have an ASD  
409 diagnosis. Since all individuals consented in SPARK are re-contactable, such studies will enable a  
410 more complete picture of the broad phenotypic effects of these variants without the bias of  
411 clinical ascertainment. Overall, these risk genes with modest effect size may represent a

412 different class of ASD genes that are more directly associated with core symptoms of ASD  
413 and/or neuropsychiatric conditions rather than global brain developmental and ID.

414 The approaches employed in this study made full use of rare variation, and this analytical  
415 method is generalizable to many conditions. In particular, the multiple methods used to reduce  
416 noise in LoF alleles present in control samples were particularly effective in assessing the signal  
417 within the novel genes of moderate effect. We also leveraged gene expression profiles  
418 informed by machine learning methods to help prioritize genes for the meta-analysis stage of  
419 our analysis<sup>38</sup>. Future studies that leverage additional multi-omic data such as dGTEX may  
420 further improve signal to noise.

421 Our archetypal analysis provides some clues as to the potential risk mechanisms of the five  
422 newly identified risk genes. *ITSN1* was unambiguously mapped to A4: membrane trafficking and  
423 protein transport and has a role in coordinating endocytic membrane traffic with the actin  
424 cytoskeleton<sup>53,54</sup> *NAV3* (A6: KRAB domain and LRR), is associated with both axon guidance<sup>55</sup> and  
425 malignant growth and invasion<sup>56</sup> and is thought to regulate cytoskeletal dynamics. Indeed, A6 is  
426 enriched for processes related to intermediate filaments (**Supplementary Table S10**) a known  
427 determinant of cell motility and polarity<sup>57</sup>. Although *MARK2*, *SCAF1*, and *HNRNPUL2* were not  
428 identified as archetypal (potentially suggesting divergence from well-known autism risk  
429 mechanisms) a search for functional enrichment of this interstitial region between A4 and A5  
430 found that their roles in developmental risk may be most relevant at the cell-cell junction,  
431 particularly as it relates to migration (see **Figure 7D**).

432 Taken together, our results suggest that a continued focus on *de novo* variants for ASD gene-  
433 discovery may yield diminishing returns. By contrast, studies designed to identify genomic risk  
434 from rare and common inherited variants will not only yield new mechanistic insight but help  
435 explain the high heritability of ASD. SPARK is designed to recruit individuals across the autism  
436 spectrum, without relying on ascertainment at medical centers. As a result, SPARK may be  
437 better suited to identify genes with transmitted variants that have lower penetrance and to  
438 identify the genetic contributions to the full spectrum of autism. The strategies employed by  
439 SPARK — to recruit and assess large numbers of individuals with autism across the spectrum  
440 and their available family members without costly, in-depth clinical phenotyping — is necessary  
441 to achieve the required sample size to fully elucidate genetic contributions to ASD. SPARK's  
442 ability to recontact and follow all participants will also be critical to deeply assess the  
443 phenotypes associated with the newly discovered genes and to develop and test novel  
444 treatments.

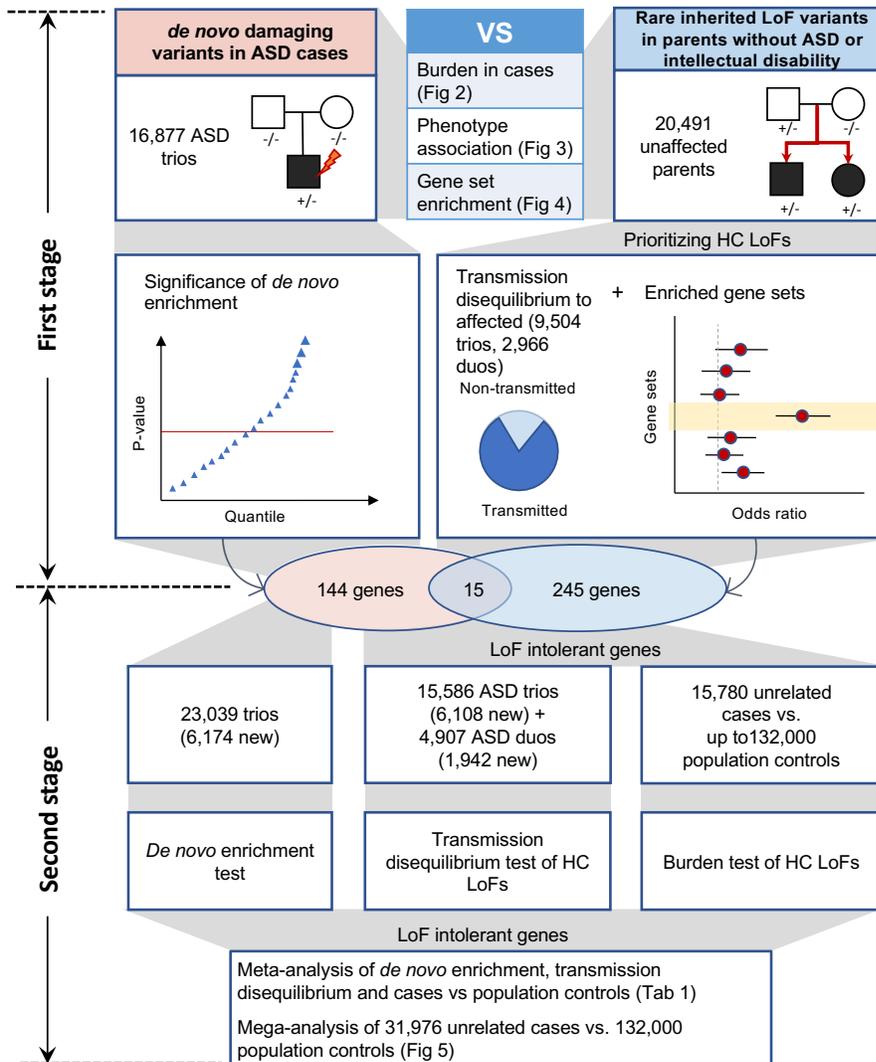
445

Gene	Prioritization	Enrichment of <i>de novo</i> damaging variants					Transmission disequilibrium of HC LoFs			case-control comparison of HC LoF rate			$P_{Meta}$
		dnLoF	$\mu_{LoF}$	dnDmis	$\mu_{Dmis}$	$P_{DNV}$	Count	Trans: Non-Trans to affected	$P_{TDT}$	Number (rate) of LoFs in cases	Rate of LoFs in controls: gnomAD exome, TOPMed	$P_{CC}$	
<i>NAV3</i>	TDT	1	1.1e-5	1	1.1e-5	0.23	17	17:2	3.6e-4	22 (1.4e-3)	3e-4, 2.6e-4	4.4e-7, 2.1e-8	1.2e-8
<i>MARK2</i>	<i>De novo</i>	5	4.4e-6	3	4.8e-6	8.9e-9	3	3:1	0.31	4 (2.5e-4)	2e-5, 6e-5	4.5e-3, 0.03	2.3e-8
<i>SCAF1</i>	TDT	2	4.8e-6	0	1.7e-7	1.3e-3	4	3:1	0.31	13 (8.2e-4)	3e-5, 7e-5	2.1e-6, 1.4e-6	2.1e-7
<i>ITSN1</i>	TDT	3	1.2e-5	2	1.3e-5	2.6e-3	18	17:2	3.6e-4	10 (6.3e-4)	1.6e-4, 2e-4	2e-3, 4e-3	4.3e-7
<i>HNRNPUL2</i>	<i>De novo</i>	3	5.8e-6	0	3.8e-6	1.8e-3	2	2:0	0.25	10 (6.3e-4)	4e-5, 5e-5	2.6e-6, 8.2e-7	2.7e-7

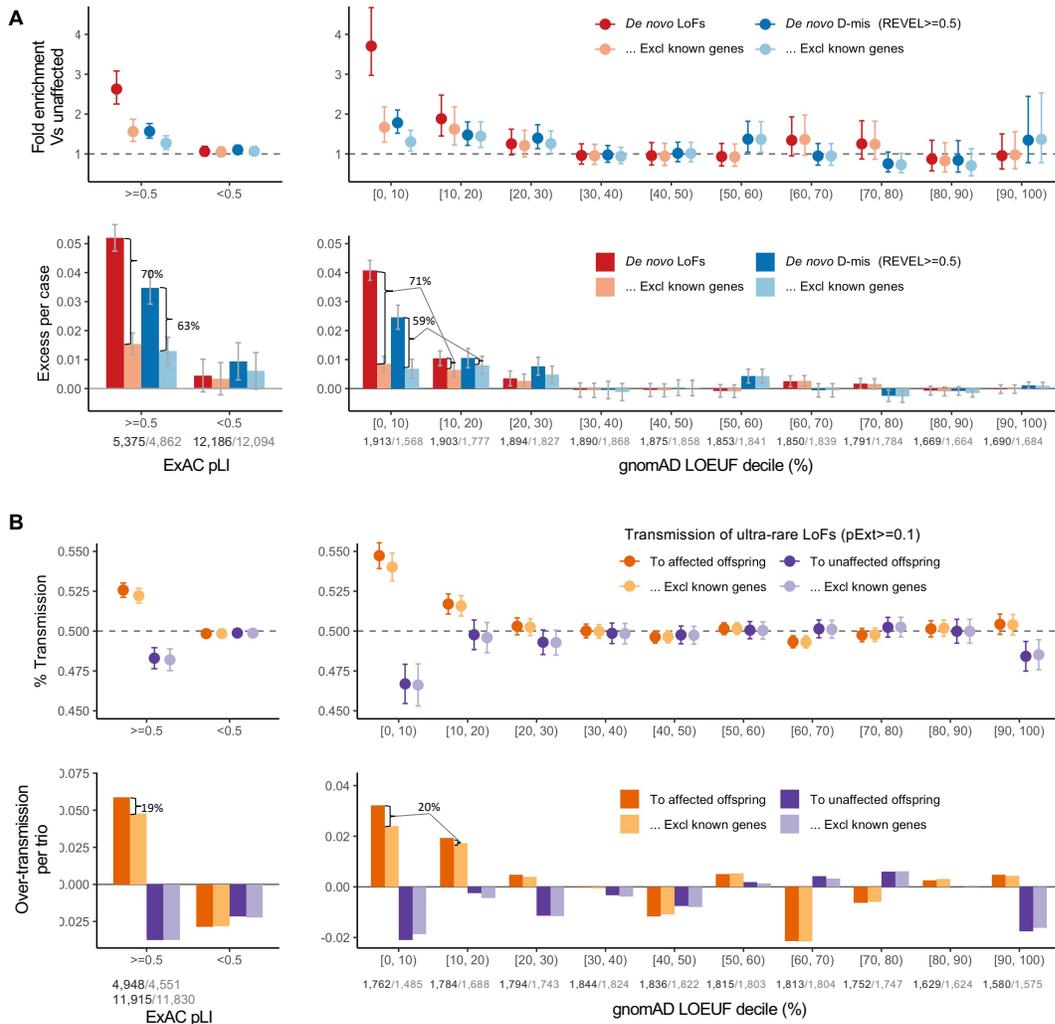
446

447 **Table 1: Statistical evidence for the five novel exome-wide significant ASD risk genes identified in this study.**

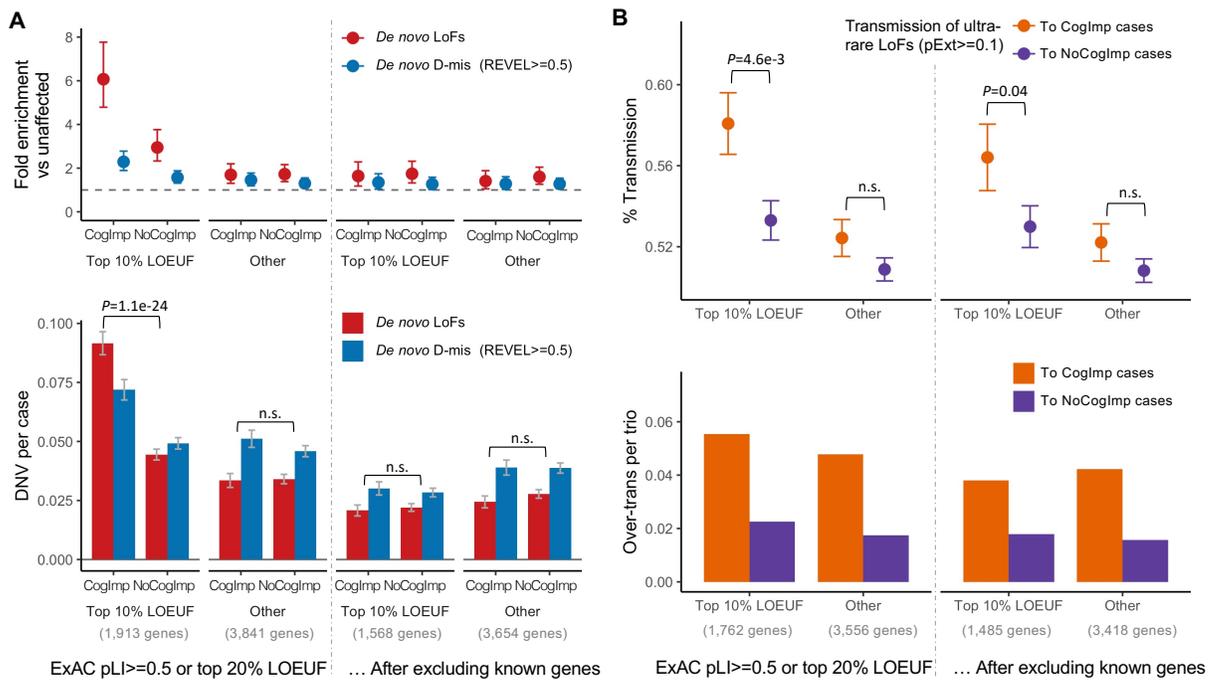
448 Control HC LoF rates are estimated from two population-based reference panels: gnomAD exome (v2.1.1, non-  
449 neuro subset, 104,068 individuals), and TopMed (freeze 8, 132,345 individuals). Meta-analysis is done by  
450 combining p-values from *de novo*, TDT and pseudo case-control analysis using Fisher's method. For pseudo case-  
451 control, we conservatively took the largest p-value for meta-analysis.  $P_{DNV}$ : One-sided p-value for enrichment of all  
452 DNVs in 23,053 ASD trios,  $P_{TDT}$ : One-sided p-value of over-transmission of HC LoFs to affected offspring in 28,556  
453 trios and 4,526 duos,  $P_{CC}$ : One-side p-value for increased HC LoF rate in 15,811 unrelated cases compared with  
454 population controls (showing two p-values from comparison with gnomAD exome and TOPMed data respectively).



455  
 456 **Figure 1. Analysis workflow.** In the discovery stage, we identified *de novo* variants in 16,877 ASD trios and rare LoF  
 457 variants in 20,491 parents without ASD diagnoses and intellectual disability. We compared properties of *de novo*  
 458 and rare variants to identify rare LoFs that contribute to genetic risk in individuals with ASD. We also evaluated  
 459 their associations with cognitive impairment and enriched gene sets. We performed an initial exome-wide scan of  
 460 genes enriched by *de novo* variants or showing transmission disequilibrium (TD) of rare LoFs to affected offspring  
 461 and selected a total of 404 genes for further replication, including 159 *de novo* enriched genes and 260 prioritized  
 462 TD genes from enriched gene sets (15 genes were in both). In the meta-analysis stage, we first evaluated evidence  
 463 from *de novo* enrichment and TD of rare, inherited LoFs in an expanded set of family-based samples including over  
 464 6,000 additional ASD trios and around 2000 additional duos. The *de novo* variants in ASD were combined with  
 465 those from additional 31,565 NDD trios to refine the filters of high confidence (HC) LoFs in *de novo* LoF enriched  
 466 genes. We also constructed an independent dataset of LoF variants of unknown inheritance from 15,780 cases that  
 467 were not used in *de novo* or transmission analysis. We compared LoF rates in cases with two population-based sets  
 468 of controls ( $n \sim 104,000$  and  $\sim 132,000$ , respectively). For 367 LoF intolerant genes on autosomes, the final gene  
 469 level evidence was obtained by meta-analyzing p-values of *de novo* enrichment, TD of HC rare, inherited LoFs, and  
 470 comparison of HC LoFs from cases and controls not used in the *de novo* or transmission analysis. We also  
 471 performed a mega-analysis that analyzed HC LoFs identified in all 31,976 unrelated ASD cases and compared their  
 472 rates with population-based controls.

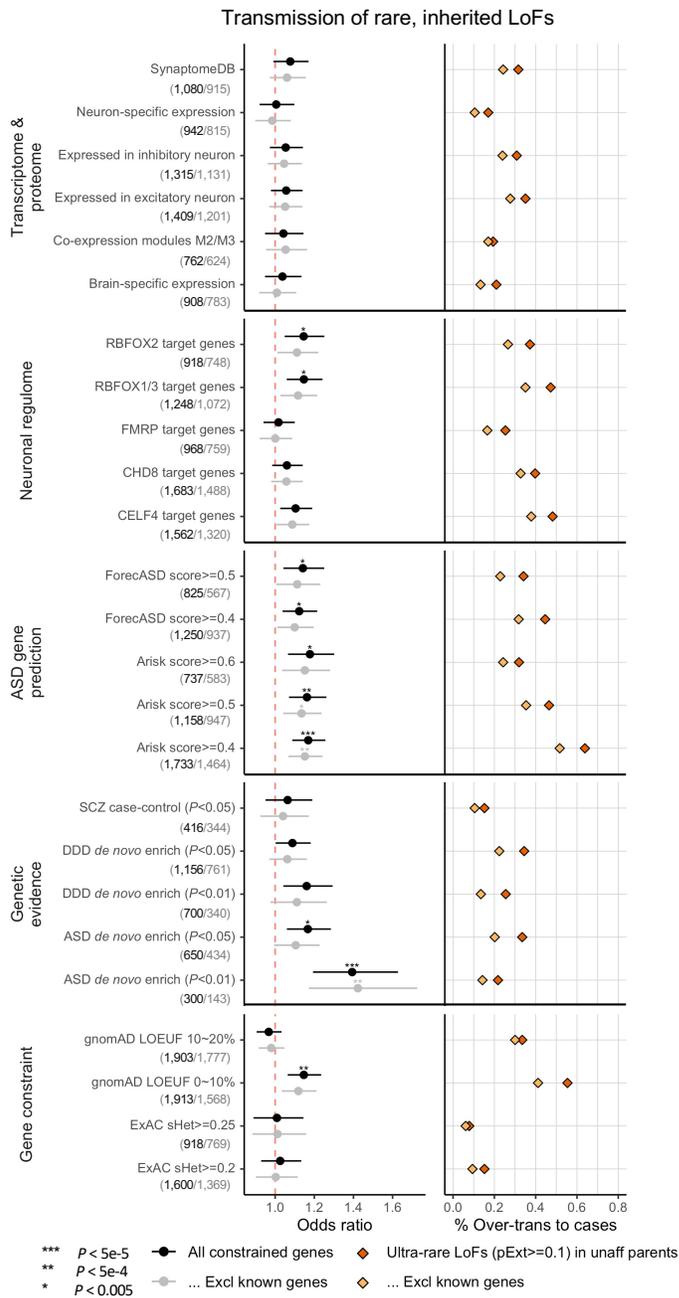


474  
 475 **Figure 2. Comparison of burden between *de novo* damaging variants and rare, inherited LoFs in ASD.** (A) The  
 476 burden of *de novo* variants was evaluated by the rate ratio and rate difference between 16,877 ASD and 5,764  
 477 unaffected trios. The exome-wide burden of *de novo* LoF and Dmis (REVEL>=0.5) variants are concentrated in  
 478 constrained genes (ExAC pLI>=0.5) and in genes with the highest levels of LoF-intolerance in the population—  
 479 defined by the top two deciles of gnomAD LOEUF scores. Burden analysis was repeated after removing known  
 480 ASD/NDD genes. The number of genes before and after removing known genes in each constraint bin is shown  
 481 below the axis label. Among constrained genes (ExAC pLI>=0.5 or the top 20% of gnomAD LOEUF scores), close to  
 482 two thirds of case-control rate differences of *de novo* LoF and Dmis variants can be explained by known genes. (B)  
 483 The burden of inherited LoFs was evaluated by looking at the proportion of rare LoFs in 20,491 parents without  
 484 ASD diagnoses or intellectual disability that are transmitted to affected offspring in 9,504 trios and 2,966 duos and  
 485 show evidence of over-transmission of LoFs per ASD trio. As a comparison, we also show the transmission  
 486 disequilibrium pattern to unaffected offspring in 5,110 trios and 129 duos. Using ultra-rare LoFs with pExt>=0.1,  
 487 exome-wide signals of transmission disequilibrium of rare, inherited LoF variants also concentrate in constrained  
 488 genes (ExAC pLI>=0.5) and in genes within the top two deciles of gnomAD LOEUF scores. Analysis was restricted to  
 489 autosomal genes and repeated after removing known ASD/NDD genes (number of genes in each constrained bin  
 490 before and after removing known genes is shown below the axis label). Among all constrained genes, only one-fifth  
 491 of over-transmission of LoFs to ASD trios can be explained by known ASD/NDD genes.



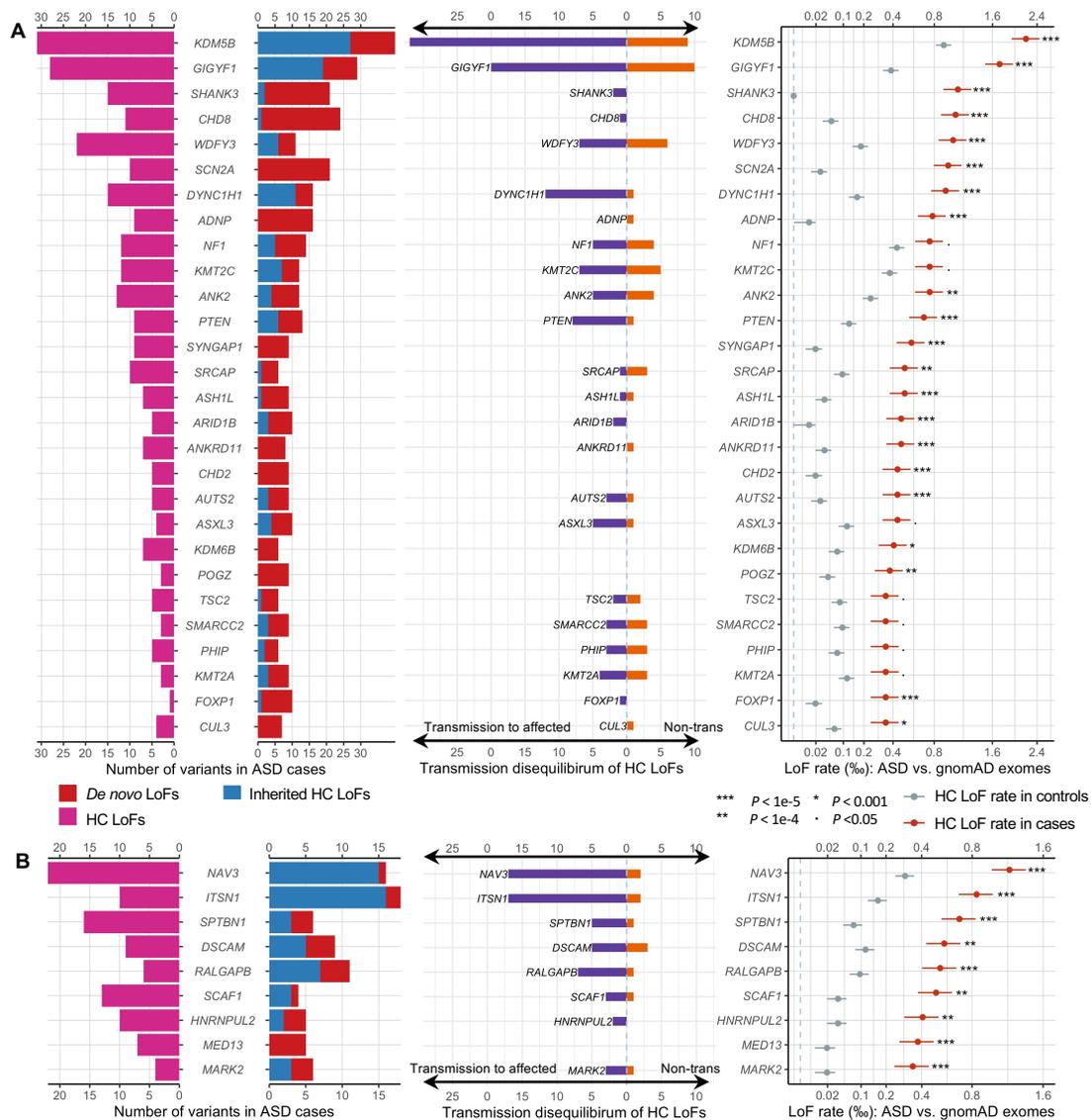
492  
493  
494  
495  
496  
497  
498  
499

**Figure 3. Association of rare, inherited LoFs with cognitive impairment in ASD cases.** Ultra-rare inherited LoFs with pExt $\geq$ 0.1 in genes with the top 10% gnomAD LOEUF scores also show a higher proportion of transmission and a higher over-transmission rate to ASD offspring with cognitive impairment than those without. Rare LoFs in other constrained genes are not significantly associated with phenotypic severity. The increased burden of inherited LoFs in cases with cognitive impairment remains significant after removing known ASD/NDD genes.

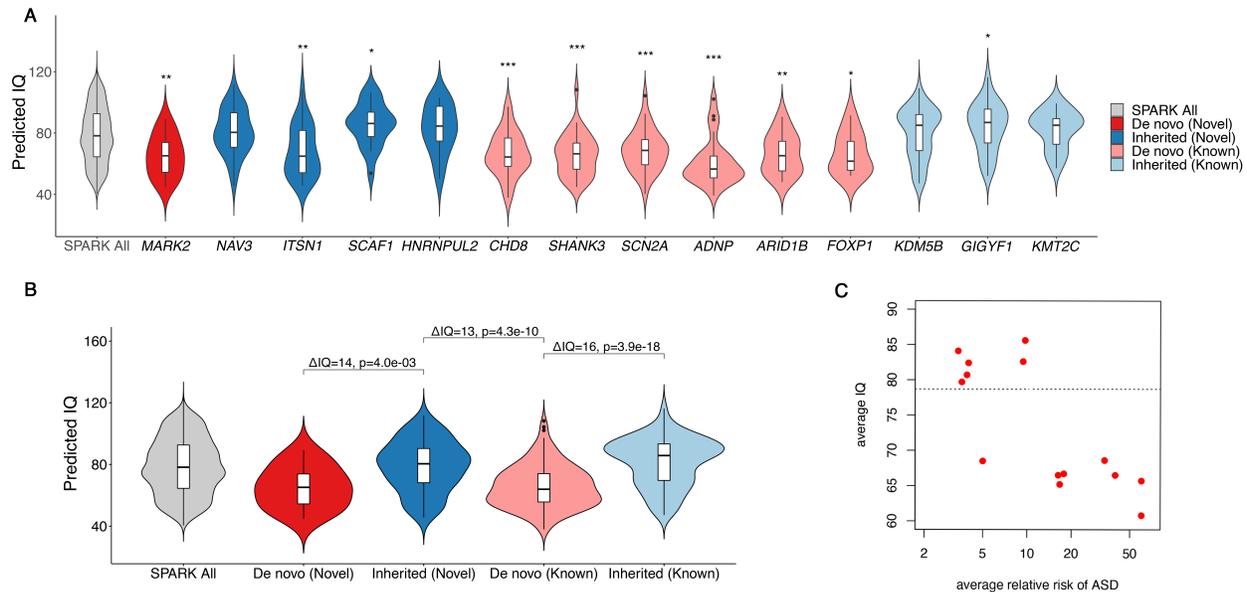


500

501 **Figure 4. Enrichment of rare LoF variants in ASD cases across gene sets.** Gene sets were defined and grouped by  
 502 transcriptome proteome, neuronal regulome, ASD gene prediction scores, genetic evidence from neuropsychiatric  
 503 diseases, and gene level constraint. Analyses were repeated after removing known ASD/NDD genes. (Number of  
 504 genes in each set before and after removing known genes are shown in bracket below gene set.) Dots represent  
 505 fold enrichment of DNVs or odds ratios for over-transmission of LoFs in each set. Horizontal bars indicate the 95%  
 506 confidence interval. For each gene set, we show the percentage of over-transmission of rare LoFs to cases.  
 507 Enrichment of rare, inherited LoFs was evaluated by comparing the transmission and non-transmission of ultra-  
 508 rare LoFs with  $p_{Ext} \geq 0.1$  in the gene set versus those in all other constrained genes using a 2-by-2 table.  $P$ -values  
 509 were given using the chi-squared test.

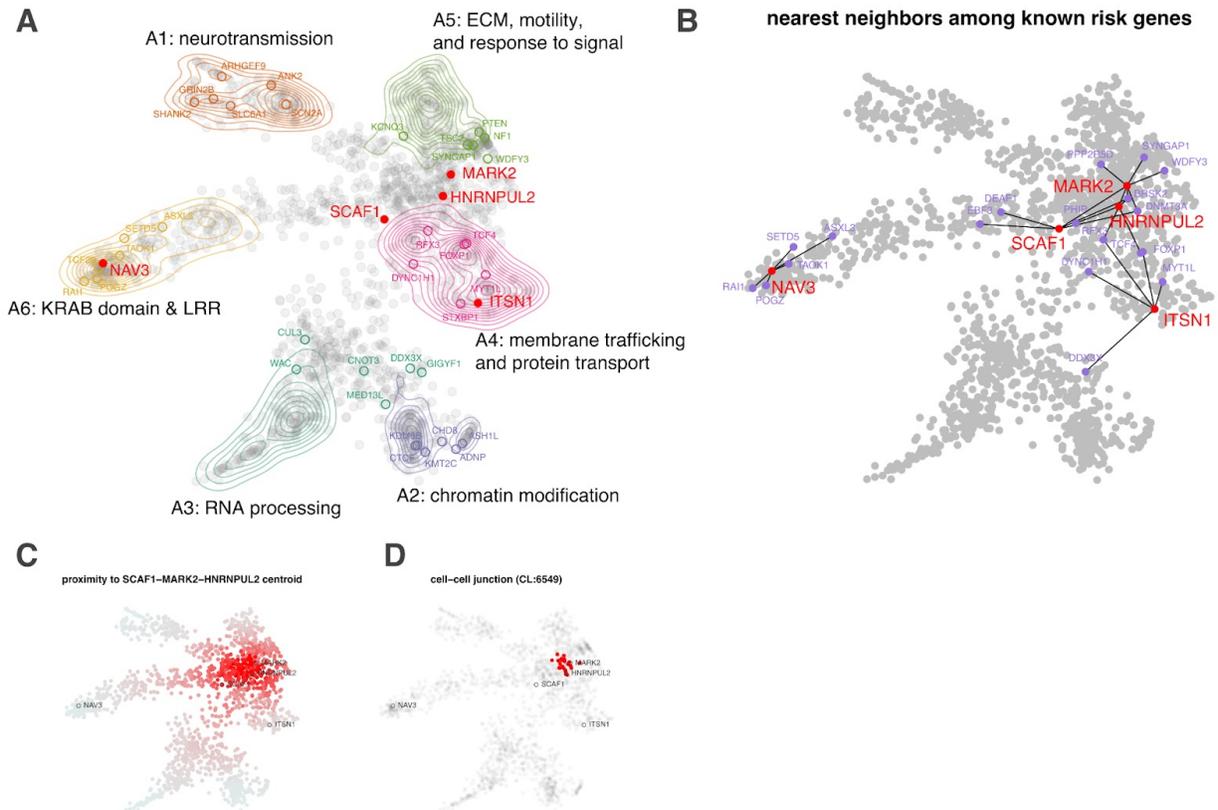


510  
 511 **Figure 5. Distribution of *de novo* and inherited LoF variants in known and novel ASD genes in cases and**  
 512 **population controls.** From left to right: pyramid plots summarizing the number of *de novo* LoF variants in 15,857  
 513 ASD trios, inherited HC LoFs in 18,720 unrelated offspring included in transmission analysis, and HC LoFs in 15,780  
 514 unrelated cases; bar plot of transmission vs. non-transmission for rare HC LoFs identified in parents without ASD  
 515 diagnoses or intellectual disability; three plots comparing the HC LoF rate in 31,976 unrelated ASD cases with  
 516 gnomAD exomes (non-neuro subset, 104,068 individuals). Horizontal bars indicate standard errors. (A) The upper  
 517 panel shows 28 known ASD/NDD genes in which LOEUF scores are in the top 30% of gnomAD, have a p-value for  
 518 enrichment among all DNVs ( $p < 9e-6$ ) in 23,039 ASD trios, and have more than 10 LoFs. (B) The lower panel shows  
 519 9 additional ASD risk genes that achieved a p-value of  $< 9e-6$  in Stage 2 of this analysis. The majority of genes in the  
 520 lower panels harbor more inherited LoFs than *de novo* variants. All five novel genes (**Error! Reference source not**  
 521 **found.**) are shown in the lower panel. Note that the x-axes of LoF rates are in the squared root scale.



522

523 **Figure 6. Predicted full-scale IQ (FSIQ) in individuals with pathogenic variants in inherited or *de novo* genes in**  
 524 **SPARK.** We examined the distribution of predicted IQ by a machine learning method<sup>48</sup> for individuals with ASD  
 525 with a LoF mutation in one of the five novel exome-wide significant genes (*MARK2*, *NAV3*, *ITSN1*, *SCAF1*,  
 526 *HNRNPUL2*) and nine known ASD genes (*CHD8*, *SHANK3*, *SCN2A*, *ADNP*, *ARID1B*, *FOXP1*, *KDM5B*, *GIGYF1*, *KMT2C*),  
 527 compared with 2,545 SPARK participants with ASD and known IQ scores. We denote the genes contributing to ASD  
 528 primarily through *de novo* LoF variants in our analysis as “*De novo*” (in red), and the genes primarily through  
 529 inherited LoF variants as “*Inherited*” (in blue). (A) Distribution of predicted IQ between individuals with ASD with  
 530 LoF mutations in the five novel genes, 9 known genes and all participants with ASD and known IQ scores in SPARK  
 531 ( $n = 2,545$ ). We compared the mean predicted IQ between participants with LoF mutations in ASD genes and all  
 532 participants by two-sample t-test. Significance level is denoted by the star sign above each violin plot (\*:  $0.01 \leq$   
 533  $p < 0.05$ , \*\*:  $0.001 \leq p < 0.01$ , \*\*\*:  $p < 0.001$ ). Individuals with pathogenic variants in *de novo* risk genes have  
 534 significantly lower predicted IQ than overall SPARK participants with ASD and known IQ scores, while individuals  
 535 with LoF variants in moderate risk, inherited genes with show similar predicted IQ as the overall SPARK  
 536 participants, with the exception of *ITSN1*. (B) Distribution of predicted IQ between individuals with ASD gene  
 537 grouped by both inheritance status (“*De novo*” or “*Inherited*”) and whether the ASD genes are novel (“*Novel*” or  
 538 “*Known*”). We compared the mean predicted IQ between individuals with pathogenic variants in “*De novo*” genes  
 539 and “*Inherited*” genes among our five novel genes and nine known genes. Overall, people with LoF mutations in  
 540 “*De novo*” genes have an average of 13-16 points lower predicted IQ than individuals with LoF mutations in  
 541 “*Inherited*” genes, regardless of whether the ASD genes are novel or known. (C) Average relative risk of ASD and  
 542 average predicted IQ among different groups. Each dot shows the average of individuals with rare LoFs of a gene  
 543 selected in panel A. The relative risk is estimated from mega analysis and capped at 60. Pearson correlation  
 544 between average IQ and log relative risk is  $-0.78$  ( $p=0.001$ ). The horizontal line represents the average IQ (IQ=79)  
 545 of all SPARK individuals with predicted IQs. *ITSN1* is an outlier at the bottom left corner.



546

547 **Figure 7. Functional/phenotypic embedding of ASD risk genes.** Using a combination of archetypal analysis and  
 548 canonical correlation analysis, putative autism risk genes were organized into  $k=6$  archetypes that represent  
 549 distinct mechanistic (STRING) and phenotypic (HPO) categorizations (A; neurotransmission, chromatin  
 550 modification, RNA processing, transport, extracellular matrix, motility and response to signal, and leucine-rich  
 551 repeat/KRAB domain containing genes). Genes implicated by our meta-analysis are indicated by their label, with  
 552 novel genes indicated in red. For each of the five novel genes, we identified the five nearest neighbors in the  
 553 embedding space among the 62 meta-analysis genes (B). *SCAF1*, *MARK2*, and *HNRNPUL2* were identified as  
 554 “mixed” rather than “archetypal” in their probable risk mechanisms. To gain further insight into possible risk  
 555 mechanisms, we calculated the embedding distance to the centroid of these three genes (C), which was then used  
 556 as an index variable to perform gene set enrichment analysis. A STRING cluster (CL:6549) containing genes related  
 557 to cell-cell junctions and the gap junction was identified as being highly localized in this region of the embedding  
 558 space ( $p = 4.12 \times 10^{-14}$  by the KS test) (D). This may suggest that these genes confer autism risk through  
 559 dysregulation of processes related to cell adhesion and migration.

560

## 561 **Methods**

562 We performed an integrated analysis of coding variants in over 35,130 new ASD cases in SPARK  
563 and additional cases from previously published autism cohorts (ASC<sup>3,8</sup>, MSSNG<sup>6</sup>, and SSC<sup>2,30</sup>),  
564 using a two-stage analysis workflow (**0 1**). In the first stage, we analyzed over 10,000 ASD cases  
565 from family-based samples and systematically compared damaging DNVs and rare, inherited  
566 LoF variants. Then we performed an exome-wide scan of genes enriched by DNVs in ASD cases  
567 and prioritized genes with suggestive evidence of DNV enrichment. We filtered for high-  
568 confidence (HC) LoF variants and searched for genes enriched by inherited HC LoFs using a  
569 transmission disequilibrium test (TDT)<sup>54</sup>. In the second stage, we added 22,764 ASD cases and  
570 used meta-analysis to further assess the prioritized genes for enrichment of DNVs and TDT of  
571 HC LoFs. For LoF intolerant genes, we compared frequency of HC LoF variants in unrelated  
572 cases, population controls, and pseudo-controls in ASD families. Finally, we performed a case-  
573 control analysis of ASD cases vs population controls to estimate effect sizes for known and  
574 newly significant genes and used them for power calculations to estimate sample sizes needed  
575 for future studies.

576

### 577 *ASD Cohorts*

#### 578 *SPARK*

579 We established SPARK (Simons Foundation Powering Autism Research for Knowledge) cohort to  
580 facilitate genotype driven research of ASD at scale<sup>23</sup>. Eligibility criteria for SPARK study is  
581 residence in the United States and a professional diagnosis of ASD or a family member of a  
582 proband in SPARK. SPARK has recruited over 50,000 re-contactable families with ASD cases at  
583 31 different clinical centers across the United States as well as through social and digital media.  
584 Individuals with known genetic diagnoses and individuals with and without a family history of  
585 autism are included. Whenever possible, parents and family members with or without autism  
586 were enrolled and included in the genetic analysis.

587 Saliva was collected using the OGD-500 kit (DNA Genotek) and DNA was extracted at  
588 PreventionGenetics (Marshfield, WI). The samples were processed with custom NEB/Kapa  
589 reagents, captured with the IDT xGen capture platform, and sequenced on the Illumina NovaSeq  
590 6000 system using S2/S4 flow cells. Samples were sequenced to a minimum standard of >85% of  
591 targets covered at 20X. 97% of samples have at least 20x coverage in >95% of region (99% of  
592 samples — in 89% of regions). Pending sample availability, any sample with 20X coverage below  
593 88% was re-processed and the sequencing events were merged to achieve sufficient coverage.  
594 The Illumina Infinium Global Screening Array v1.0 (654,027 SNPs) was used for genotyping. The  
595 average call rate is 98.5%. Less than 1% of samples have a call rate below 90%.

596

597 In the first stage of analysis, we included 28,649 SPARK individuals including 10,242 ASD cases  
598 from over 9,000 families with exome sequencing data that passed QC (**Error! Reference source  
599 not found.**). A subset of 1,379 individuals was part of the previously published pilot study<sup>7</sup>. To  
600 replicate prioritized genes from the discovery stage, we performed a second stage analysis that  
601 included an additional 39,926 individuals with 16,970 ASD cases from over 20,000 families with

602 exome or whole genome sequencing (WGS) data available after of the analysis in discovery  
603 cohort was completed. For new samples in this study, exome sequences were captured by IDT  
604 xGEN research panel and sequenced on the Illumina NovaSeq system. DNA samples were also  
605 genotyped for over 600K SNPs by Infinium Global Screening Array.

606 We used KING<sup>55</sup> to calculate statistics for pairwise sample relatedness from genotypes of  
607 known biallelic SNPs, and validated participant-reported familial relationships (**Supplementary**  
608 **Figure S21A-B**). The relatedness analysis also identified cryptically related families that are  
609 connected by unreported parent-offspring or full sibling pairs. Pedigrees were reconstructed  
610 manually from inferred pairwise relationships and validated by PRIMUS<sup>56</sup> and we used inferred  
611 pedigree for all analyses. Sample sex was validated by normalized sequencing depths or array  
612 signal intensities of X and Y chromosomes which also identified X and Y chromosome  
613 aneuploidies (**Supplementary Figure S21C-D**). To infer genetic ancestry, we first performed  
614 principal component (PC) analysis on SNP genotypes of non-admixed reference population  
615 samples from 1000 Genomes Projects<sup>57</sup> (Africans, Europeans, East Asians and South Asians) and  
616 Human Genome Diversity Project<sup>58,59</sup> (Native Americans), then projected SPARK samples onto  
617 PC axes defined by the five reference populations using EIGNSOFT<sup>60</sup> (**Supplementary Figure**  
618 **S22**). The projected coordinates on first four PC axes were transformed into probabilities of five  
619 population ancestries using the method of SNPweights<sup>61</sup>. The inferred ancestral probabilities  
620 show general concordance with self-reported ethnicities (**Supplementary Figure S22B**).  
621 Samples were predicted from a reference population if the predicted probability was  $\geq 0.85$ .

622 The phenotypes of participants are based on self- or parent-report provided at enrollment and  
623 in a series of questionnaires from the Simons Foundation Autism Research Initiative database,  
624 SFARI Base. We used SFARI Base Version 4 for the discovery cohort and Version 5 for the  
625 replication cohort. In the discovery cohort, information about self-reported cognitive  
626 impairment (or intellectual disability/developmental delay) was available for 99.2% of ASD  
627 cases and 83.5% of other family members at recruitment or from the Basic Medical Screening  
628 Questionnaire available on SFARIbase. For phenotype-genotype analyses in individuals with  
629 variants in specific ASD risk genes, we defined an individual as having cognitive impairment if 1)  
630 there was self- or parent-report of cognitive impairment at registration or in the Basic Medical  
631 Screening Questionnaire, 2) the participant was at or over the age of 6 at registration and was  
632 reported to speak with less than full sentences or the participant was at or above age 4 at  
633 registration and reported as non-verbal at that time, 3) the parent reported that cognitive  
634 abilities were significantly below age level, 4) the reported IQ or the estimated cognitive age  
635 ratio (ratio IQ<sup>62,63</sup>) was  $< 80$  or 5) the parent reported unresolved regression in early childhood  
636 without language returning and the participant does not speak in full sentences. The  
637 continuous full-scale IQ was imputed based on a subset of 521 samples with full scale IQ and  
638 phenotypic features by the elastic net machine learning model<sup>48</sup>. In a subset of cases for which  
639 full-scale IQ data or standardized Vineland adaptive behavior scores (version 3) was available,  
640 we found self-reported cognitive impairment shows higher correlation with Vineland score than  
641 full-scale IQ (**Supplementary Figure S23**). ASD cases with self-reported cognitive impairment  
642 were defined as Cognitively Impaired cases, and other cases as Not Cognitively Impaired cases.  
643 Other non-ASD family members were considered as unaffected if they were also not indicated  
644 to have cognitive impairment. In total of 18.5% families, proband has at least one first-degree

645 relative with ASD who was recruited in the study and/or reported by a family member. Those  
646 families were referred to as multiplex, and other families with only a single ASD individual as  
647 simplex. The majority (>85%) of affected relative pairs in multiplex families were siblings.  
648 Multiplex families have slightly lower male-to-female ratio and lower proportion of cognitive  
649 impairment among affected offspring (**Supplementary Figure S24A-B**). In comparison, only 1%  
650 of parents in the discovery cohort are affected of which two thirds are females and less than 3%  
651 have cognitive impairment (**Supplementary Figure S24A-B**). In addition, non-ASD family  
652 members in multiplex families show significantly higher frequency of self-reported cognitive  
653 impairment, learning/language disorders, other neuropsychiatric conditions, and other types of  
654 structural congenital anomalies (**Supplementary Figure S24C**). Non-ASD parents in multiplex  
655 families also have lower educational attainment (**Supplementary Figure S24D**).

#### 656 SSC

657 SSC (Simon Simplex Collection) collected over 2,500 families with only one clinically confirmed  
658 ASD cases who have no other affected first or second degree relatives as an effort to identify *de*  
659 *novo* genetic risk variants for ASD<sup>64</sup>. SSC data have been published before<sup>2,19,30,65</sup>. Here we  
660 included 10,032 individuals including 2,633 cases with exome or WGS data available and passed  
661 QC (**Error! Reference source not found.**). The data were reprocessed using the same pipeline as  
662 SPARK. For 91 trios that are not available or incomplete, we collected coding DNVs from  
663 published studies<sup>2,30</sup>. In analysis to associate genetic variants with phenotype severity, we used  
664 standardized Vineland adaptive behavior score to group affected cases because it shows higher  
665 correlation than full-scale IQ with self-reported cognitive impairment in SPARK (**Supplementary**  
666 **Figure S23**). Cases with cognitive impairment in SSC were defined by Vineland score $\leq$ 70, and  
667 cases with no cognitive impairment by score $>$ 70.

#### 668 ASC

669 ASC (Autism Sequencing Consortium) is an international genomics consortium to integrate  
670 heterogenous ASD cohorts and sequencing data from over 30 different studies<sup>66</sup>. Individual  
671 level genetic data are not available. So we included 4,433 published trios (4,082 affected and  
672 351 unaffected) merged from two previous studies<sup>3,8</sup> for DNV analysis. To define low and high  
673 functioning cases, we used binary indicator of intellectual disability which was available for 66%  
674 of cases. Families with multiple affected trios are considered multiplex, others are simplex.

#### 675 MSSNG

676 The MSSNG initiative aims to generate WGS data and detailed phenotypic information of  
677 individuals with ASD and their families<sup>6</sup>. It comprehensively samples families with different  
678 genetic characteristics in order to delineate the full spectrum of risk factors. We included 3,689  
679 trios in DB6 release with whole genome DNV calls are available and passed QC in DNV analysis,  
680 of which 1,754 trios were published in the previous study<sup>6</sup>. A total of 3,404 offspring with a  
681 confirmed clinical diagnosis of ASD were included as cases. Among individuals without a  
682 confirmed ASD diagnosis, 222 who did not show broader or atypical autistic phenotype or other  
683 developmental disorders were used as part of controls. Multiplex families were defined as  
684 families having multiple affected siblings in sequenced trios or in phenotype database.  
685 Information about cognitive impairment was not available at the time of analysis.

686 *Variant calling and quality control*

687 Supplementary Table S11 describes software version and parameter settings for each analysis  
688 below.

689 *Data processing*

690 Sequencing reads were mapped to human genome reference (hg38) using bwa-mem<sup>67</sup> and  
691 stored in CRAM format<sup>68</sup>. Duplicated read pairs in the same sequencing library of each  
692 individual were marked up by MarkupDuplicates of Picard Tools<sup>69</sup>. Additional QC metrics for GC  
693 bias, insert size distribution, hybridization selection were also calculated from mapped reads by  
694 Picard Tools<sup>69</sup>. Mosdepth<sup>70</sup> was used to calculate sequencing depth on exome targets (or 500  
695 bp sliding windows for WGS) and determine callable regions at 10X or 15X coverage. Cross-  
696 sample contamination was tested by VerifyBamID<sup>71</sup> using sequencing only mode. Samples were  
697 excluded if it has insufficient coverage (less than 80% targeted region with  $\geq 20X$ ), shows  
698 evidence of cross-sample contamination (FREEMIX $>5\%$ ), or discordant sex between normalized  
699 X and Y chromosome depth and self/parent reports that cannot be explained by aneuploidy.

700 Variants for each individual were discovered from mapped reads using GATK HaplotypeCaller<sup>72</sup>,  
701 weCall<sup>73</sup>, and DeepVariant<sup>74</sup>. Individual variant calls from GATK and weCall were stored in gVCF  
702 format and jointly genotyped across all samples in each sequencing batch using GLnexus<sup>75</sup>.  
703 Variants were also jointly discovered and genotyped for individuals of the same family using  
704 GATK HaplotypeCaller<sup>72</sup> and freebayes<sup>76</sup>, and then read-backed phased using WhatsHap<sup>77</sup>. To  
705 verify sample relatedness, identify overlapping samples with other cohorts, and verify sample  
706 identity with SNP genotyping data, genotypes of over 110,000 known biallelic SNPs from 1000  
707 Genomes or HapMap projects that have call rate  $>98\%$  and minor allele frequency (MAF)  $>1\%$  in  
708 the cohort were extracted from joint genotyping VCFs. SNP array genotypes were called by  
709 Illumina GenomeStudio. We kept samples with  $>90\%$  non-missing genotype calls and used  
710 genotypes of over 400,000 known SNPs that have call rate  $>98\%$  and MAF $>0.1$  for relatedness  
711 check and ancestry inference.

712 *De novo variants*

713

714 We identified candidate *de novo* SNVs/indels from SPARK and SSC cohorts from per-family VCFs  
715 generated by GATK and freebayes and cohort-wide population VCF by weCall using a set of  
716 heuristic filters that aim to maximize the sensitivity while minimizing false negatives in parents<sup>7</sup>.  
717 We then reevaluated the evidence of all *de novo* candidates from all input sources. Candidate  
718 was removed if there was contradictory evidence against from any input source (“contradiction  
719 filters”, see **Supplementary Table S11**). Further, we only kept candidates if they can be called  
720 by DeepVariant in offspring but have no evidence of variant in parents. For candidates that  
721 were identified in multiple offspring (recurrent), we only kept the ones that passed DeepVariant  
722 filter in all trios. For candidates that were shared by siblings in the same family, we only kept  
723 the ones with *de novo* quality estimated by triodenovo higher than 8 (or 7 for SNVs in CpG  
724 context). Before creating the final cleaned call set, we selected subsets of variants (see  
725 **Supplementary Table S11**) for manual evaluation by IGV to filter out candidates with failed  
726 review. Finally, we merged nearby clustered *de novo* coding variants (within 2bp for SNVs or  
727 50bp for indels) on the same haplotype to form multi-nucleotide variants (MNVs) or complex

728 indels. We removed variants located in regions known to be difficult for variant calling (HLA,  
729 mucin, and olfactory receptors). DNVs in the final call set follow a Poisson distribution with an  
730 average 1.4 coding DNVs per affected and 1.3 per unaffected offspring (**Supplementary Figure**  
731 **S25**). The proportion of different types of DNVs, the mutation spectrum of SNVs, and indel  
732 length distributions were similar between SPARK and SSC (**Supplementary Figure S25**). A small  
733 fraction of variants in the final call set are likely post-zygotic mosaic mutations (**Supplementary**  
734 **Figure S26**).

735

#### 736 *Rare variants*

737 Rare variant genotypes were filtered from cohort-wide population VCFs with QC metrics  
738 collected from individual and family VCFs (**Supplementary Figure S27A**). Briefly, we initially  
739 extracted high quality genotypes for each individual for variants that appear in less than 1% of  
740 families in the cohort. Evidence for the variant genotypes were re-evaluated by DeepVariant  
741 from aligned reads and collapsed over individuals to create site level summary statistics  
742 including fraction of individual genotypes that passed DeepVariant filter and mean genotype  
743 quality over all individuals. For variant genotypes extracted from GLnexus VCFs, we re-  
744 examined variant genotype from per-family VCFs by GATK to collect GATK site level metrics  
745 (including QD, MQ, SOR, etc.) then took read-depth weighted average over families to create  
746 cohort-wide site metrics. For variant genotypes extracted GATK joint genotyping VCFs, these  
747 site metrics were directly available directly from INFO fields.

748 Variant site level QC filters were calibrated using familial transmission information, assuming  
749 that false positive calls are more likely to show Mendelian inheritance error (**Supplementary**  
750 **Figure S27B**). Briefly, we first applied a baseline site level filter that favors high sensitivity, then  
751 optimized thresholds for filters with additional QC metrics. The selected QC metrics were  
752 reviewed first to determine a small number of optional thresholds. Then the final set of QC  
753 parameters were optimized from a grid search over the combinations of available thresholds  
754 such that: 1. presumed neutral variants identified from parents (silent variants or variants in  
755 non-constrained genes) shows equal transmission and non-transmission to offspring; 2. rates of  
756 neutral variants are similar in different sample groups from the same population ancestry; 3.  
757 vast majority variants identified in trio offspring are inherited from parents. In case when  
758 multiple sets of QC thresholds give similar results, priority will be given to the set that also  
759 recovers maximum number of DNV calls in trio offspring. The optimized filtering parameters  
760 were used in final QC filters to generate analysis-ready variants.

761 For a rare coding variant initially annotated as LoF (including stop gained, frameshift, or splice  
762 site), we searched for nearby variants on the same haplotype (within 2bp for SNVs or 50bp for  
763 indels). If nearby variants can be found, they were merged to form MNVs or complex indel and  
764 re-annotated to get the joint functional effect. If the joint effect was not LoF, then the original  
765 variant was removed from LoF analysis.

#### 766 *Variant annotations*

767 The genomic coordinates of QC passed variants were lifted over to hg19 and normalized to the  
768 leftmost positions<sup>78</sup>. Functional effects of coding variants were annotated to protein coding  
769 transcripts in GENCODE V19 Basic set<sup>79</sup> using variant effect predictor<sup>80</sup>. The gene level effect

770 was taken from the most severe consequences among all transcripts (based on the following  
771 priority: LoF>missense>silent>intronic). pExt for each variant can be operationally defined as  
772 the proportion of expression levels of transcripts whose variant effects are the same as gene  
773 effect over all transcripts included in the annotation<sup>27</sup>. We used transcript level expressions in  
774 prenatal brain development from Human Developmental Biology Resource<sup>81</sup> to calculate pExt.  
775 Missense variants were annotated by pathogenicity scores of REVEL<sup>31</sup>, CADD<sup>82</sup>, MPC<sup>83</sup> and  
776 PrimateAI<sup>84</sup>. Population allele frequencies were queried from gnomAD<sup>26</sup> and ExAC<sup>18</sup> using all  
777 population samples. All rare variants were defined by cohort allele frequency <0.001 (or <0.005  
778 for X chromosome variants). To filter for ultra-rare variants, we keep variants with cohort allele  
779 frequency <1.5e-4 (or allele count=1) and population allele frequency <5e-5 in both gnomAD<sup>26</sup>  
780 and ExAC<sup>18</sup>.

781 LoF variants on each coding transcript were further annotated by LOFTEE<sup>26</sup> (v1.0, default  
782 parameters). We also annotated splice site variants by SpliceAI<sup>85</sup>, and removed low confidence  
783 splice site variants with delta score <0.2 from LoF variants. pExt for LoF variants was calculated  
784 by the proportion of expression level of transcripts that harbor HC LoFs evaluated by LOFTEE  
785 over all transcripts included in the analysis. Thus, the pExt filter for LoFs already incorporated  
786 LOFTEE annotations. The baseline filter to analyze rare, inherited LoFs and LoFs of unknown  
787 inheritance is pExt>=0.1. To refine gene-specific pExt threshold in the second stage, we selected  
788 95 known ASD/NDD genes plus a newly significant DNV enriched gene *MARK2* which harbor at  
789 least four *de novo* LoF variants in combined ASD and other NDD trios, and for each gene choose  
790 the pExt threshold from {0.1,0.5,0.9} that can retain all *de novo* LoF variant with pExt>=0.1  
791 (Supplementary Table S1).

#### 792 *Copy number variants*

793 Copy number variants (CNVs) were called from exome read depth using CLAMMS<sup>86</sup>. CNV calling  
794 windows used by CLAMMS were created from exome targets after splitting large exons into  
795 equally sized windows of roughly 500bp. Calling windows were annotated by average  
796 mappability score<sup>87</sup> (100mer) and GC content assuming average insert size of 200. Depths of  
797 coverage for each individual on the windows were calculated using Mosdepth<sup>70</sup> and then  
798 normalized to control for GC-bias and sample's overall average depth. Only windows with GC  
799 content between 0.3 and 0.75 and mappability >=0.75 were included in further analyses. For  
800 each given sample, we used two approaches to reduce the dimension of sample's coverage  
801 profile and automatically selected 100 nearest neighbors of the sample under analysis as  
802 reference samples. The first approach used seven QC metrics calculated by Picard Tools from  
803 aligned reads as recommended by the CLAMMS developer<sup>86</sup>, we further normalized those  
804 metrics in the cohort by its median absolute deviation in the cohort. The second approach used  
805 singular value decomposition of the sample by read-depth matrix to compute the coordinates  
806 of the first 10 principal components for each sample.

807 Model fitting and CNV calling for each individual using custom reference samples were  
808 performed using default parameters. From raw CNV calls, neighboring over-segmented CNVs of  
809 the same type were joined if joined CNVs include over 80% of the calling windows of original  
810 calls. For each sample, we kept CNV calls made from one set of reference samples that have  
811 smaller number of raw CNV calls. Outliers with excessive raw CNV calls (>400) were removed.  
812 For each CNV, we counted the number of CNVs of the same type in parents that overlap >50%

813 of the calling windows. High-quality rare CNVs were defined as <1% carrier frequency among  
814 parents and have Phred-scaled quality of CNV in the interval >90. We queried high-quality rare  
815 copy number deletions to look for additional evidence to support new genes.

816 Genetic analysis

817 *De novo* variants analysis

818 In the discovery stage analysis, the DNV call sets of SPARK and SSC were merged with published  
819 DNVs from ASC<sup>3,8</sup> and MSSNG<sup>6</sup> and additional SSC trios of which we did not have sequencing  
820 data. To infer likely samples overlaps with published trios of which we do not have individual  
821 level data, we tallied the proportion of shared DNVs between all pairs of trios. For a pair of  
822 trios, let  $N_1$  and  $N_2$  be the number of coding DNVs and  $O$  the number of shared DNVs between  
823 pair. To account for mutation hotspots, if a DNV is a SNV within CpG context or a known  
824 recurrent DNVs identified in SPARK and SSC, it contributes 0.5 to the count. Likely overlapping  
825 samples were identified if  $\frac{O}{N_1} \geq 0.5$  or  $\frac{O}{N_2} \geq 0.5$  and they have identical sex.

826 To determine the expected number of DNVs in the cohort, we used a 7-mer mutation rate  
827 model<sup>52</sup> in which the expected haploid mutation rate of each base pair (bp) depends on the 3bp  
828 sequence context on both sides. The per-base mutation rates were adjusted by the fraction of  
829 callable trios at each base pair which was the fraction of trios with  $\geq 10X$  coverage in parents  
830 and  $\geq 15X$  coverage in offspring. For published trios, we used an inhouse WGS data of 300 trios  
831 with average 36X coverage to approximate the callable regions. Gene level haploid mutation  
832 rates for different classes of DNVs were calculated by summing up the depth-adjusted per-base  
833 mutation rate of all possible SNVs of the same class. The rate for frameshift variants was  
834 presumed to be 1.3 times the rate of stop gained SNVs<sup>53</sup>. Mutation rates in haploid X  
835 chromosome regions were adjusted for the observed male-female ratio (4.2) assuming  
836 mutation rates in spermatogenesis is 3.4 times higher than oogenesis<sup>9</sup>. The exome-wide rate of  
837 synonymous DNVs closely matches the observed number of DNVs (**Supplementary Figure S12**).  
838 We also observed similar fold enrichment of damaging DNVs (vs. expected rate) in ASD cases  
839 across four cohorts after accounting for samples with family history (**Supplementary Figure**  
840 **S12**).

841 To perform gene-based test of DNVs, we applied DeNovoWEST<sup>11</sup> a simulation-based approach  
842 to test the enrichment of weighted sum of different classes of DNVs compared to the expected  
843 sum based on per-base mutation rates in each gene. We used empirical burden of DNVs to  
844 derive weights for different variant classes in constrained genes (ExAC pLI $\geq 0.5$ ) and non-  
845 constrained genes separately based on positive predictive values (PPV) (**Supplementary Table**  
846 **S13**). For ASD, we defined *de novo* D-mis variants by REVEL score  $\geq 0.5$ , and the rest of *de novo*  
847 missense variants are taken as benign missense (B-mis). For other NDDs, we defined two  
848 classes of *de novo* D-mis variants by MPC score  $\geq 2$  or MPC  $\leq 2$  and CADD score  $\geq 25$ , and the  
849 remaining *de novo* missense variants are B-mis. We first ran DeNovoWEST to test the  
850 enrichment of all nonsynonymous DNVs (pEnrichAll). To account for risk genes that harbor only  
851 missense variants, we ran DeNovoWEST to test the enrichment of *de novo* missense variants  
852 only and applied a second test for spatial clustering of missense variants using DenovoNear<sup>9</sup>,  
853 then combined evidence of missense enrichment and clustering (pCombMis). The minimal of  
854 pEnrichAll and pCombMis was used as the final p-value for DeNovoWEST. The exome-wide

855 significance threshold was set to  $1.3e-6$  ( $=0.05/(18,000 \text{ genes} * 2 \text{ tests})$ ) to account for the two  
856 tests. The analysis on replication cohort used the same weights as derived from discovery  
857 cohort. Compared with the original publication<sup>11</sup>, our implementation of DeNovoWEST used  
858 different ways to stratify genes, determine variant weights, and calculate per-base mutation  
859 rates. We applied our DeNovoWEST implementation on 31,058 NDD trios and compared with  
860 published results on the same data set. The p-values from re-analysis show high overall  
861 concordance with published results (**Supplementary Figure S28**). We used p-values from our re-  
862 analysis on other NDD trios in comparative analysis with ASD.

863 Gene set enrichment analysis of DNVs was performed by DnEnrich framework<sup>32</sup>. We included  
864 all *de novo* LoF and D-mis variants in 5,754 constrained genes from 16,877 ASD and 5,764  
865 control trios. For each gene set, we calculated the fraction of weighted sums of damaging DNVs  
866 in the set using PPV weights of constrained genes (**Supplementary Table S13**) for cases and  
867 controls respectively. The test statistics for each gene set is the ratio of such fractions in cases  
868 over controls. To determine the distribution of test statistic under the null hypothesis, we  
869 randomly placed mutations onto the exome of all constrained genes, while held the number of  
870 mutations, their tri-nucleotide context and functional impact to be the same as observed in  
871 cases and controls separately. Note that by conditioning on the observed number of damaging  
872 DNVs in cases and controls, we tested enriched gene sets in cases that are not due to an  
873 increased overall burden. At each round of simulation, the permuted test statistic in each gene  
874 set was calculated. Finally, the p-value was calculated as number of times the permuted  
875 statistic is greater than or equal to observed statistic. Fold enrichment (FE) was calculated as  
876 the ratio of between observed and average of test statistics over all permutations. We also  
877 approximated 95% confidence interval for FE by assuming  $\log(\text{FE})$  follows normal distribution  
878 with mean 0 and standard deviation determined by the p-value.

879 In all DNV analyses above, DNVs shared by full or twin siblings represent single mutational  
880 events and were counted only once. When an individual carry multiple DNVs within 100bp  
881 in the same gene, only one variant with most severe effects was included in the analysis.

#### 882 Transmission disequilibrium analysis

883 The effect of inherited LoF variants was analyzed using TDT in each individual genes or in gene  
884 sets. Rare LoF variants were first identified in parents without ASD diagnoses or intellectual  
885 disability who have at least one offspring, then for each parent-offspring pair, the number of  
886 times the LoF variant was transmitted from parents to offspring was tallied. For variants in  
887 (non-PAR part of) X chromosome, we only used rare LoF variants carried by mothers without  
888 ASD diagnoses or intellectual disability and analyzed transmission in different types of mother-  
889 offspring pairs. For TDT analysis of rare, inherited missense variants in selected gene sets,  
890 different D-mis definitions and allele frequency cutoffs were used (**Supplementary Figure S3**).

891 The over-transmission of LoFs to affected offspring was evaluated by a binomial test assuming  
892 transmission equilibrium under the null hypothesis of 50% chance of transmission. In the  
893 discovery stage, ultra-rare LoFs with  $p\text{Ext} \geq 0.1$  were used in exome-wide transmission  
894 disequilibrium and gene set enrichment analysis. For gene-based test, all rare LoFs with  
895  $p\text{Ext} \geq 0.1$  were also used, and TDT statistic<sup>39</sup> for each gene was calculated by  $z = \frac{T-NT}{\sqrt{T+NT}}$ , where  
896  $T(NT)$  is the number of times LoF variants were transmitted (not transmitted) to affected

897 offspring. When offspring include monozygotic twin pairs, only one was kept in the  
898 transmission analysis. We prioritized 244 autosomal genes with  $z > 1$  in top 10% LOEUF or in top  
899 20% LOEUF and  $A\text{-risk} \geq 0.4$ . In the second stage gene-based test, if a gene-specific pExt  
900 threshold is available, we used HC LoF variants passed the gene-specific pExt filter.

901 In gene set enrichment analysis of inherited LoFs, the rate of transmission to affected offspring  
902 in each gene set was compared with the transmission rate in rest of the genes in the  
903 background using chi-squared test.

#### 904 Case control analysis

905 Pseudo-controls are constructed from parents without ASD diagnoses or intellectual disability  
906 in simplex families, using alleles that were not transmitted to affected offspring. Each parent  
907 without ASD diagnoses or intellectual disability contributes sample size of 0.5 to pseudo-  
908 controls. Rare LoFs in ASD cases whose parent data are not available and from other cases that  
909 were not utilized in DNV enrichment or TDT analysis were analyzed in this stage. Specifically, for  
910 each ASD case, we found out all his/her most recent unaffected ancestors without ASD  
911 diagnoses or intellectual disability in the pedigree and calculated the contributing sample size  
912 as 1 minus the summation of kinship coefficients with these ancestors. If the contributing  
913 sample size is greater than 0, then the sample was included in pseudo-cases after removing  
914 alleles that were observed in any unaffected ancestors without ASD diagnoses or intellectual  
915 disability used in TDT and alleles included in DNV analysis if any. Examples of such rare LoFs in  
916 cases and their contributing sample sizes are given in Supplementary Figure S29.

917 Rare LoFs in cases and controls for X chromosome were categorized separately for males and  
918 females. For male controls, because fathers do not transmit X chromosomes to sons, male  
919 controls include all fathers. In contrast, male cases only include those whose mothers do not  
920 have ASD diagnoses or intellectual disability (thus not included in TDT analysis). For females,  
921 because we only include mothers without ASD diagnoses or intellectual disability and affected  
922 sons in TDT, female pseudo-cases include all affected females. Female pseudo-controls were  
923 established from unaffected mothers in simplex families using alleles that do not transmit to  
924 affected sons. Each unaffected mother contributes a sample size of 0.5 to pseudo-controls. In  
925 both sexes, DNVs were removed from pseudo-cases.

926 For gene-based tests in Stage 2, case-control comparisons are not independent of TDT. So we  
927 used population references as controls, including gnomAD exomes<sup>26</sup> (v2.1.1 non-neuro subset),  
928 gnomAD genomes<sup>26</sup> (v3.1 non-neuro subset), and TopMed genomes<sup>88</sup> (Freeze 8). Variants in  
929 the population references were filtered to keep those passed default QC filter in released data.  
930 For variants in gnomAD data set, we further removed variants located in low complexity region,  
931 because such regions are enriched with false positive calls<sup>89</sup> but the default filter does not  
932 effectively remove variants in those regions. QC filters in the inhouse ASD cohort and in  
933 TopMed had already removed most of variants located in such regions. Variants from  
934 population references were re-annotated in the same way as rare variants identified in ASD  
935 cohort. In gene level case-control comparison of LoF burden, we used baseline  $p\text{Ext} \geq 0.1$  filter  
936 or gene-specific pExt threshold if available to define HC LoF variants. For LoF variants in  
937 selected genes, we also extracted curation results by gnomAD to remove curated non-LoF  
938 variants and manually reviewed IGV snapshots from gnomAD browser if available to remove

939 likely variant calling artifacts (Supplementary Data 1). Number of HC LoF variants were obtained  
 940 from the summation of allele count in site level VCF files. Gene level burden of HC LoF variants  
 941 between cases and population controls are tested by comparing the HC LoF variant rates  
 942 between cases and controls using Poisson test. To account for different in depth of coverage,  
 943 sample sizes are multiplied by the fraction of callable coding regions of each gene ( $\geq 15X$  for  
 944 autosomes or female X chromosome,  $\geq 10X$  for male X chromosome) in ASD cases and in  
 945 population controls respectively.

946 To account for sample relatedness in case-control analysis, we created a relationship graph in  
 947 which each node represents an individual and each edge represents a known first or second-  
 948 degree relationship between two individuals. We also add edges to pairs of individuals without  
 949 known familial relationship but have estimated kinship coefficient  $\geq 0.1$ . From the graph, we  
 950 select one individual from each connected component to create unrelated case-control  
 951 samples. For chromosome X, father and sons were treated as unrelated. For population  
 952 controls, only gnomAD data included sex specific allele counts and were used in the sex-specific  
 953 analysis.

954 Meta-analysis was performed for prioritized autosomal genes among top 30% LOEUF. We  
 955 integrated evidence from the enrichment of all DNVs, transmission disequilibrium, and  
 956 increased burden in case compared with population controls by combining p-values using  
 957 Fisher's method<sup>40</sup>. Experiment-wide error rate was set at  $9e-6$  ( $=0.05$  divided by 5340  
 958 autosomal genes at LOEUF 30%). In mega-analysis, we combined all unrelated ASD cases  
 959 together and compared CAFs of HC LoF variants with three population references.

960 Power calculation

961 To calculate statistical power of the current study and to estimate sample size for future gene  
 962 discovery efforts, we adopted the statistical framework by Zuk et al. 2014<sup>41</sup> comparing CAF of  
 963 LoF variants in  $N$  unrelated cases  $f_{\text{case}}$  with CAF  $f$  in natural population. The effect of LoFs in  
 964 the same gene are assumed to be the same and increase ASD risk by  $\gamma$  fold. The population CAF  
 965  $f$  is assumed to be known with high precision from large cohorts. Since we only focus on LoF-  
 966 intolerant genes in the population,  $f$  is assumed to be at selection-mutation equilibrium  $f =$   
 967  $\frac{\mu_{\text{LoF}}}{s}$  where  $\mu_{\text{LoF}}$  is LoF mutation rate and  $s$  is selection coefficient. The test statistic  
 968 asymptotically follows a non-central chi-squared distribution with 1-df and non-centrality  
 969 parameter (NCP):

$$970 \quad \lambda = 4N \left[ \gamma f \ln \gamma + (1 - \gamma) \ln \frac{1 - \gamma f}{1 - f} \right]$$

971 Given the significance threshold  $\alpha$ , power can be calculated analytically by

$$972 \quad 1 - \beta = 1 - F(F^{-1}(1 - \alpha, 0), \lambda)$$

973 where  $F(x, \lambda)$  is the cumulative distribution of  $\chi_1^2$  with NCP  $\lambda$ .

974 To calculate sample size to achieve desired power  $1 - \beta$  at significance level  $\alpha$ , we first solve  
 975 NCP  $\lambda_{\alpha, \beta}$  from the above equation. Then sample size can be approximated by:

$$976 \quad n_{\alpha, \beta} \approx \frac{\lambda_{\alpha, \beta}}{4f[\gamma \ln \gamma - (\gamma - 1)]}$$

977 For current study in ASD, sample size is  $N=31,976$  unrelated cases, experimental wide error  
978 rate is  $\alpha=9e-6$ . Given continuing expansion of population reference, treating  $f$  as known  
979 without error is a reasonable assumption for future studies. To calculate power for new genes  
980 identified in this study, we used point estimates of  $\gamma$  and  $f$  from mega-analysis using gnomAD  
981 exomes as population controls, and used  $\mu_{LoF}$  computed from the 7mer context dependent  
982 mutation rate model<sup>52</sup> to convert  $f$  to  $s = \frac{\mu_{LoF}}{f}$ . The required sample sizes were calculated to  
983 achieve 90% of power.

984 Power and sample size are both calculated as a function of relative risk for ASD ( $\gamma$ ) and  
985 selection coefficient ( $s$ ) across different haploid LoF mutation rates ( $\mu_{LoF}$ ). We only considered  
986  $s$  between 0.01 and 0.5, because most prioritized genes have point estimates of  $s>0.01$ (Error!  
987 Reference source not found.) and genes with  $s>0.5$  are expected to harbor to *de novo* than  
988 inherited LoF variants and can to be identified from the enrichment of DNVs. Relative risk to  
989 ASD ( $\gamma$ ) was constrained between 1 and 20 since we are mainly interested in discovering genes  
990 with moderate to small effects. The reduction in fitness  $s$  is correlated with the increases in ASD  
991 risk  $\gamma$  by  $s = \gamma\pi s_D$  under the assumption of no pleiotropic effect, where  $\pi$  is ASD prevalence  
992 and  $s_D$  is decreased reproductive fitness of ASD cases. Based on epidemiological studies,  
993 current estimated prevalence of ASD is  $\hat{\pi}=1/54$ <sup>90</sup>, estimated  $s_D$  is for 0.75 male and for 0.52  
994 female<sup>91</sup> so sex averaged  $\hat{s}_D=0.71$  (assuming male-to-female ratio of 4.2). In reality, most  
995 known ASD genes also show pleiotropic effects with other NDDs or associated with prenatal  
996 death and therefore  $s \geq \gamma\pi s_D \approx \gamma\hat{\pi}\hat{s}_D = 0.013\gamma$ . So we only considered combinations of  $(s, \gamma)$   
997 that satisfy the condition:  $s \geq 0.013\gamma$ .

#### 998 Gene sets

999 To evaluate the contribution of known ASD risk genes to the burdens of DNVs and inherited LoF  
1000 variants identified in this study, we collected 618 known dominant ASD/NDD genes from the  
1001 following sources:

- 1002 1. Known developmental disorder genes from DDG2P<sup>92</sup> (2020-02) that are dominant or X-  
1003 linked and have organ specificity list includes brain or cause multi-system syndrome.
- 1004 2. High confidence ASD genes collected by SFARI<sup>93</sup> (2019-08) with score of 1 or 2 excluding  
1005 known recessive genes.
- 1006 3. Newly emerging dominant ASD genes reported in recent literatures and included in  
1007 SPARK genes list<sup>94</sup> (2020-07).

1008 To evaluate the gene sets enriched by damaging DNVs or inherited HC LoFs, we used all  
1009 constrained genes by ExAC pLI $\geq$ 0.5 or in top 20% of LOEUF as the background. Gene sets of the  
1010 following five categories were collected for gene sets enrichment analysis.

#### 1011 Transcriptome and proteome

- 1012 • For genes with brain-specific expression, we used processed RNA-seq data from  
1013 Fagerberg *et al.* 2014<sup>95</sup> and selected genes with average reads per kilobase of transcript  
1014 per million mapped reads (RPKM) $>$ 1 in brain and over four times of median RPKM of 27  
1015 tissues.
- 1016 • Genes in co-expression modules M2 and M3 derived from weighted gene correlation  
1017 network analysis (WGCNA) analysis of BrainSpan developmental RNAseq data were

1018 previously reported to enrich for known ASD genes<sup>33</sup> and collected from Table S1 from  
1019 that reference.

- 1020 • To find genes expressed in excitatory or inhibitory neurons, we selected genes from Mo  
1021 *et al.* 2015<sup>96</sup> that have average transcripts per million (TPM) greater than 100 in  
1022 excitatory and inhibitory neurons respectively.
- 1023 • Synaptic genes including those encode presynaptic proteins, presynaptic active zone,  
1024 synaptic vesicles, and postsynaptic density were collected from SynaptomeDB<sup>97</sup>.

1025 **Neuronal regulome**

- 1026 • Putative CELF4 target genes are defined as genes whose iCLIP occupancy>0.2 in Wagnon  
1027 *et al.* 2012<sup>98</sup>.
- 1028 • CHD8 target genes are defined as genes whose promoter or enhancer region overlap  
1029 with CHD8 binding peaks in human neural stem cells or mid-fetal brain in Cotney *et al.*  
1030 2015<sup>36</sup>.
- 1031 • FMRP target genes in mouse were first collected from Table S2C of Darnell *et al.* 2011<sup>35</sup>  
1032 with FDR<0.1. They were then mapped to orthologous human genes using homology  
1033 mapping provided by MGI<sup>99</sup> (2018-07).
- 1034 • Genes targeted by RBFOX2 were selected from Weyn-Vanhentenryck *et al.* 2014<sup>34</sup> to  
1035 have Rbfox2 tag counts greater 8. Due to high correlations between RBFOX1 and  
1036 RBFOX3, targeted genes by the two RNA binding proteins were merged in one gene set  
1037 and selected to have total tag counts of Rbfox1 and Rbfox3 greater than 24. Selected  
1038 mouse genes symbols were then mapped to orthologous human genes using homology  
1039 mapping provided by MGI.

1040 **Autism gene predictions**

- 1041 • ForecASD is an ensemble classifier that integrates brain gene expression, heterogeneous  
1042 network data, and previous gene-level predictors of autism association to yield a single  
1043 prediction score<sup>37</sup>. We created two sets of genes with forecASD prediction score greater  
1044 than 0.4 or 0.5.
- 1045 • A-risk is a classifier that uses a used gradient boosting tree to predict autism candidate  
1046 genes using cell-type specific expression signatures in fetal brain<sup>38</sup>. We created three  
1047 sets of genes with prediction score greater 0.4, 0.5 or 0.6.

1048 **Genetic evidence**

- 1049 • For genes enriched by DNVs in ASD, we selected genes showing nominal statistical  
1050 evidence (P<0.01 or P<0.05 by DeNovoWEST) in discovery cohort of 16,877 trios.
- 1051 • For genes implicated by in other NDD, we selected genes nominally enriched by DNVs in  
1052 31,058 NDDs<sup>11</sup> (P<0.01 or P<0.01 by DeNovoWEST using our implementation).
- 1053 • For genes in implicated in schizophrenia, we selected genes nominally significant  
1054 (P<0.05) by gene-based test in latest schizophrenia case-control study of 24,248 cases  
1055 and 97,322 controls<sup>47</sup>.

1056 **Archetypal analysis:** STRING v11<sup>100</sup> clusters and Human Phenotype Ontology (HPO)<sup>101</sup>  
1057 terms were formatted as gene-by-term binary matrices. The working gene list was taken as  
1058 the union of forecASD top decile genes and the 62 autism-associated gene from this study  
1059 (total 1,776 genes). A total of 583 genes from this set had annotations in both STRING and

1060 HPO, and using these genes, a canonical correlation analysis (CCA) was carried out using the  
1061 RGCCA package for R (<https://cran.r-project.org/web/packages/RGCCA/index.html>) using  
1062 five components and sparsity parameter  $c1$  set to 0.8 for both the HPO and STRING  
1063 matrices. Component scores for all 1,776 genes were calculated using the STRING cluster  
1064 annotations and the corresponding coefficients from the CCA. This 1,776 gene by 5 CC  
1065 component matrix was used as input for archetypal analysis<sup>102</sup>, and the optimal  $k$  (number  
1066 of archetypes) was selected using the elbow plot heuristic<sup>103</sup>, with the residual sums of  
1067 squares (RSS) plotted as a function of  $k$ . We displayed the archetypal embedding using the  
1068 `simplexplot()` function of the archetypes R package. Genes were identified as “archetypal” if  
1069 their top archetype coefficient was  $> 2x$  the next highest archetypal coefficient. Those genes  
1070 that did not fulfill this criterion were classified as “mixed”, while those that did were  
1071 assigned to their maximally-scoring archetype. Each of the six identified archetypes were  
1072 given a human-readable summary description based on review of the top associated  
1073 STRING clusters (Figure 7). Further cluster/term association results are available in  
1074 Supplementary Table S10. Representative genes for each archetype were chosen from  
1075 among the list of 62 risk genes identified in this study, using the top 6 genes for each  
1076 archetype (note that these genes do not necessarily fulfill the “archetypal” criterion  
1077 described above, but are simply the top six of the 62 for each archetype).

1078 **Author Contributions** I.Astrovskaya, J.B.H., J.J.M., N.V., P.F., C.Shu, T.W., W.K.C., X.Z., and Y.S.  
1079 designed and conceived this study. A.Adams, A.Andrus, A.Berman, A.Brown, A.C., A.C.G., A.D.S.,  
1080 A.E., A.Fanta, A.Fatemi, A.Fish, A.Goler, A.Gonzalez, A.Gutierrez, Jr., A.Hardan, A.Hess,  
1081 A.Hirshman, A.Holbrook, A.J.A., A.J.Griswold, A.Jarratt, A.Jelinek, A.Jorgenson, A.Juarez, A.Kim,  
1082 A.Kitaygorodsky, A.L., A.L.R., A.L.W., A.M.D., A.Mankar, A.Mason, A.Miceli, A.Milliken, A.M.-L.,  
1083 A.N.S., A.Nguyen, A.Nicholson, A.Nishida, A.P., A.P.M., A.R.G., A.Raven, A.Rhea, A.Simon,  
1084 A.Swanson, A.Sziklay, A.Tallbull, A.Tesng, A.W., A.Z., B.A.H., B.B., B.E., B.E.R., B.Hauf, B.J.O., B.L.,  
1085 B.M.V., B.S., B.V., C.A.E., C.A.W.S., C.Albright, C.Anglo, C.B., C.C.B., C.C.-S., C.Cohen, C.Colombi,  
1086 C.D., C.E., C.E.R., C.Fassler, C.Gray, C.Gunter, C.H.W., C.K., C.Leonczyk, C.L.M., C.Lord, C.M.T.,  
1087 C.M., C.O.-L., C.Ortiz, C.P., C.R.R., C.Roche, C.Shrier, C.Smith, C.V., C.W.-L., C.Zaro, C.Zha, D.B.,  
1088 Dan.Cho, D.Correa, D.E.S., D.G., D.G.A., D.H., D.I., D.L.C., D.Li, D.Limon, D.Limpoco, D.P.,  
1089 D.Rambeck, D.Rojas, D.Srishyla, D.Stamps, E.A.F., E.Bahl, E.B.-K., E.Blank, E.Bower, E.Brooks,  
1090 E.C., E.Dillon, E.Doyle, E.Given, E.Grimes, E.J., E.J.F., E.K., E.L.W., E.Lamarche, E.Lampert, E.M.B.,  
1091 E.O'Connor, E.Ocampo, E.Orrick, E.P., E.R., E.S., E.T.M., E.V.P., F.F., F.K.M., G.A., G.B., G.D., G.H.,  
1092 G.M., G.S., G.S.D., G.T., H.C., H.E.K., H.G., H.H., H.K., H.L.S., H.Lechniak, H.Li, H.M., H.R.,  
1093 H.Zaydens, I.Arriaga, I.F.T., J.A., J.A.G., J.Beeson, J.Brown, J.Comitre, J.Cordova, J.D., J.F.C.,  
1094 J.F.H., J.Gong, J.Gunderson, J.H., J.J.M., J.Judge, J.Jurayj, J.Manoharan, J.Montezuma, J.N., J.O.,  
1095 J.Pandey, J.Piven, J.Polanco, J.Polite, J.R., J.S., J.S.S., J.T.M., J.Tjernagel, J.Toroney, J.V.-V.,  
1096 J.Wang, J.Wright, K.A., K.A.S., K.Baalman, K.Beard, K.Callahan, K.Coleman, K.D.F., K.Dent,  
1097 K.Diehl, K.G., K.G.P., K.H., K.L., K.L.P., K.Murillo, K.Murray, K.N., K.O., K.Pama, K.R., K.Singer,  
1098 K.Smith, K.Stephenson, K.T., L.A., L.A.C., L.Beeson, L.Carpenter, L.Casten, L.Coppola, L.Cordiero,  
1099 L.D., L.D.P., L.F.C., L.G.S., L.H.S., L.K.W., L.L., L.M.H., L.M.P., L.Malloch, L.Mann, L.P.G., L.S.,  
1100 L.V.S., L.W., L.Y., L.Y.-H., M.A., M.Baer, M.Beckwith, M.Casseus, M.Coughlin, M.Currin, M.Cutri,  
1101 M.DuBois, M.Dunlevy, M.F., M.F.G., M.G., M.Haley, M.Heyman, M.Hojlo, M.J., M.J.M.,  
1102 M.Kowanda, M.Koza, M.L., M.M., M.N., M.N.H., M.O., M.P., M.R., M.Sabiha, M.Sahin, M.Sarris,  
1103 M.Shir, M.Siegel, M.Steele, M.Sweeney, M.T., M.V.-M., M.Verdi, M.Y.D., N.A., N.Bardett,  
1104 N.Berger, N.C., N.D., N.G., N.H., N.Lillie, N.Long, N.M.R.-P., N.Madi, N.Mccoy, N.N., N.Rodriguez,  
1105 N.Russell, N.S., N.Takahashi, N.Targalia, N.V., O.N., O.Y.O., P.F., P.H., P.M., P.S.C., R.A.B., R.A.G.,  
1106 R.C.S., R.D.A., R.D.C., R.J., R.J.L., R.K.E., R.L., R.P.G.-K., R.Remington, R.S., R.T.S., S.A., S.Birdwell,  
1107 S.Boland, S.Booker, S.Carpenter, S.Chintalapalli, S.Conyers, S.D., S.D.B., S.E., S.F., S.G.,  
1108 S.Hepburn, S.Horner, S.Hunter, S.J.B., S.J.L., S.Jacob, S.Jean, S.Kim, S.Kramer, S.L.F., S.Licon,  
1109 S.Littlefield, S.M.K., S.Mastel, S.Mathai, S.Melnyk, S.Michaels, S.Mohiuddin, S.Palmer, S.Plate,  
1110 S.Q., S.R., S.Sandhu, S.Santangelo, S.Skinner, S.T., S.Xu, S.Xiao, Sa.White, St.White, T.C., T.G.,  
1111 T.H., T.I., T.K., T.P., T.R., T.S., T.Thomas, T.Tran, V.Galbraith, V.Gazestani, V.J.M., V.R., V.S.,  
1112 W.C.W., W.Cal, W.K.C., W.S.Y., Y.C. and Z.E.W. recruited participants and collected clinical data  
1113 and biospecimens. A.Amatya, A.Bashar, A.E.L., A.Mankar, A.Nguyen, B.J., C.Rigby, Dav.Cho,  
1114 D.V.M., E.O'Connor, J.A., M.D.M., M.E.B., N.Lawson, N.Lo, N.V., R.M., R.Rana, S.G., S.Jean,  
1115 S.Shah and W.Chin built and supported the SPARKforAutism.org website, software, databases  
1116 and systems, and managed SPARK data. A.D.K., A.J.Gruber, A.Nishida, B.Han, B.J.O., C.Fleisch,  
1117 C.Shu, D.V.M., E.Brooks, G.J.F., I.Astrovskaya, J.B.H., J.J.M., J.U.O., J.Wright, L.Brueggeman,  
1118 L.G.S., M.A.P., N.Lo, N.V., O.M., P.F., S.D.B., S.Murali, S.X.X., T.N.T., T.S.C., T.W., W.H., X.Z. and  
1119 Y.S. performed analyses, processed biospecimens and sequenced DNA samples. A.Fatemi,  
1120 A.Kitaygorodsky, A.Soucy, C.Shu, D.H.G., E.B.-K., E.E.E., E.R., H.Q., H.Zhang, H.Zhao,  
1121 I.Astrovskaya, J.B.H., J.J.M., J.U.O., J.Wright, L.Brueggeman, L.G.S., M.Y.D., N.V., O.M., P.F.,

1122 R.N.D., S.D.B., S.Murali, S.X.X., T.K., T.N.T., T.P., T.S., T.Thomas, T.W., T.W.Y., W.K.C., X.Z. and  
1123 Y.S. helped with data interpretation. A.E.L., E.E.E., J.J.M., N.V., P.F., W.K.C. and X.Z. supervised  
1124 the work. B.J.O., I.Astrovskaya, J.B.H., J.J.M., J.U.O., C.Shu, J.Wright, N.V., P.F., T.N.T., T.W.,  
1125 W.K.C., X.Z. and Y.S. wrote this paper.

1126

1127 Competing interests

1128 D.H.G. has received research funding from Takeda Pharmaceuticals, and consulting fees or  
1129 equity participation for scientific advisory board work from Ovid Therapeutics, Axial Bio-  
1130 therapeutics, Acurastem, and Falcon Computing. E.E.E. is on the Scientific Advisory Board (SAB)  
1131 of DNAnexus, Inc. M.Sahin has received research funding from Novartis, Roche, Biogen,  
1132 Astellas, Aeovian, Bridgebio, Aucta and Quadrant Biosciences and has served on Scientific  
1133 Advisory Boards for Roche, Celgene, Regenxbio and Takeda. W.K.C. serves on the Regeneron  
1134 Genetics Center Scientific Advisory Board and is the Director of Clinical Research for SFARI.  
1135 A.D.K. is an employee of PreventionGenetics and a member of PrevGen Employees LLC, which  
1136 owns units in PreventionGenetics. Z.E.W. serves as a consultant for Roche and receive research  
1137 support from Adaptive Technology Consulting. All other authors declare no competing  
1138 interests.

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148 **The SPARK Consortium:** Adrienne Adams<sup>13</sup>, Alpha Amatya<sup>2</sup>, Alicia Andrus<sup>14</sup>, Asif Bashar<sup>2</sup>, Anna  
1149 Berman<sup>15</sup>, Alison Brown<sup>16</sup>, Alexies Camba<sup>17</sup>, Amanda C. Gulsrud<sup>17</sup>, Anthony D. Krentz<sup>18</sup>, Amanda  
1150 D. Shocklee<sup>19</sup>, Amy Esler<sup>20</sup>, Alex E. Lash<sup>2</sup>, Anne Fanta<sup>21</sup>, Ali Fatemi<sup>22</sup>, Angela Fish<sup>23</sup>, Alexandra  
1151 Goler<sup>2</sup>, Antonio Gonzalez<sup>24</sup>, Anibal Gutierrez, Jr.<sup>24</sup>, Antonio Hardan<sup>25</sup>, Amy Hess<sup>26</sup>, Anna  
1152 Hirshman<sup>13</sup>, Alison Holbrook<sup>2</sup>, Andrea J. Ace<sup>2</sup>, Anthony J. Griswold<sup>27</sup>, Angela J. Gruber<sup>18</sup>, Andrea  
1153 Jarratt<sup>28</sup>, Anna Jelinek<sup>29</sup>, Alissa Jorgenson<sup>28</sup>, A. Pablo Juarez<sup>30</sup>, Annes Kim<sup>23</sup>, Alex  
1154 Kitaygorodsky<sup>31</sup>, Addie Luo<sup>32</sup>, Angela L. Rachubinski<sup>33</sup>, Allison L. Wainer<sup>13</sup>, Amy M. Daniels<sup>2</sup>,  
1155 Anup Mankar<sup>2</sup>, Andrew Mason<sup>34</sup>, Alexandra Miceli<sup>15</sup>, Anna Milliken<sup>35</sup>, Amy Morales-Lara<sup>36</sup>,  
1156 Alexandra N. Stephens<sup>2</sup>, Ai Nhu Nguyen<sup>2</sup>, Amy Nicholson<sup>30</sup>, Anna Marie Paolicelli<sup>37</sup>, Alexander P.  
1157 McKenzie<sup>16</sup>, Abha R. Gupta<sup>21</sup>, Ashley Raven<sup>23</sup>, Anna Rhea<sup>38</sup>, Andrea Simon<sup>39</sup>, Aubrie Soucy<sup>40</sup>,  
1158 Amy Swanson<sup>15</sup>, Anthony Sziklay<sup>34</sup>, Amber Tallbull<sup>33</sup>, Angela Tesng<sup>28</sup>, Audrey Ward<sup>38</sup>, Allyson  
1159 Zick<sup>23</sup>, Brittani A. Hilscher<sup>41</sup>, Brandi Bell<sup>38</sup>, Barbara Enright<sup>42</sup>, Beverly E. Robertson<sup>2</sup>, Brenda  
1160 Hauf<sup>43</sup>, Bill Jensen<sup>2</sup>, Brandon Lobisi<sup>24</sup>, Brianna M. Vernoia<sup>2</sup>, Brady Schwind<sup>2</sup>, Bonnie VanMetre<sup>16</sup>,  
1161 Craig A. Erickson<sup>29</sup>, Catherine A.W. Sullivan<sup>21</sup>, Charles Albright<sup>26</sup>, Claudine Anglo<sup>41</sup>, Cate  
1162 Buescher<sup>4</sup>, Catherine C. Bradley<sup>38</sup>, Claudia Campo-Soria<sup>28</sup>, Cheryl Cohen<sup>2</sup>, Costanza Colombi<sup>23</sup>,  
1163 Chris Diggins<sup>2</sup>, Catherine Edmonson<sup>16</sup>, Catherine E. Rice<sup>45</sup>, Carrie Fassler<sup>29</sup>, Catherine Gray<sup>43</sup>,  
1164 Chris Gunter<sup>46</sup>, Corrie H. Walston<sup>43</sup>, Cheryl Klaiman<sup>46</sup>, Caroline Leonczyk<sup>13</sup>, Christa Lese  
1165 Martin<sup>47</sup>, Catherine Lord<sup>17</sup>, Cora M. Taylor<sup>47</sup>, Caitlin McCarthy<sup>38</sup>, Cesar Ochoa-Lubinoff<sup>48</sup>, Crissy  
1166 Ortiz<sup>38</sup>, Cynthia Pierre<sup>13</sup>, Cordelia R. Rosenberg<sup>33</sup>, Chris Rigby<sup>2</sup>, Casey Roche<sup>38</sup>, Clara Shrier<sup>38</sup>,  
1167 Chris Smith<sup>34</sup>, Candace Van Wade<sup>38</sup>, Casey White-Lehman<sup>2</sup>, Christopher Zaro<sup>35</sup>, Cindy Zha<sup>24</sup>,  
1168 Dawn Bentley<sup>14</sup>, Dahriana Correa<sup>24</sup>, Dustin E. Sarver<sup>49</sup>, David Giancarla<sup>24</sup>, David G. Amaral<sup>41</sup>,  
1169 Dain Howes<sup>28</sup>, Dalia Istephanous<sup>28</sup>, Daniel Lee Coury<sup>26</sup>, Deana Li<sup>41</sup>, Danica Limon<sup>39</sup>, Desi  
1170 Limpoco<sup>32</sup>, Diamond Phillips<sup>13</sup>, Desiree Rambeck<sup>28</sup>, Daniela Rojas<sup>25</sup>, Diksha Srishyla<sup>28</sup>, Danielle  
1171 Stamps<sup>49</sup>, Dennis Vasquez Montes<sup>2</sup>, Daniel Cho<sup>50</sup>, Dave Cho<sup>2</sup>, Emily A. Fox<sup>50</sup>, Ethan Bahl<sup>4</sup>,  
1172 Elizabeth Berry-Kravis<sup>48</sup>, Elizabeth Blank<sup>29</sup>, Erin Bower<sup>34</sup>, Elizabeth Brooks<sup>2</sup>, Eric Courchesne<sup>34</sup>,  
1173 Emily Dillon<sup>16</sup>, Erin Doyle<sup>38</sup>, Erin Given<sup>26</sup>, Ellen Grimes<sup>15</sup>, Erica Jones<sup>2</sup>, Eric J. Fombonne<sup>32</sup>,  
1174 Elizabeth Kryszak<sup>26</sup>, Ericka L. Wodka<sup>16</sup>, Elena Lamarche<sup>43</sup>, Erica Lampert<sup>29</sup>, Eric M. Butter<sup>26</sup>,  
1175 Eirene O'Connor<sup>2</sup>, Edith Ocampo<sup>13</sup>, Elizabeth Orrick<sup>25</sup>, Esmeralda Perez<sup>2</sup>, Elizabeth Ruzzo<sup>5</sup>, Emily  
1176 Singer<sup>2</sup>, Emily T. Matthews<sup>35</sup>, Ernest V. Pedapati<sup>29</sup>, Faris Fazal<sup>32</sup>, Fiona K. Miller<sup>23</sup>, Gabriella  
1177 Aberbach<sup>35</sup>, Gabriele Baraghoshi<sup>14</sup>, Gabrielle Duhon<sup>39</sup>, Gregory Hooks<sup>28</sup>, Gregory J. Fischer<sup>18</sup>,  
1178 Gabriela Marzano<sup>39</sup>, Gregory Schoonover<sup>14</sup>, Gabriel S. Dichter<sup>43</sup>, Gabrielle Tiede<sup>26</sup>, Hannah  
1179 Cottrell<sup>19</sup>, Hannah E. Kaplan<sup>34</sup>, Haidar Ghina<sup>50</sup>, Hanna Hutter<sup>16</sup>, Hope Koene<sup>21</sup>, Hoa Lam  
1180 Schneider<sup>24</sup>, Holly Lechniak<sup>13</sup>, Hai Li<sup>48</sup>, Hadley Morotti<sup>32</sup>, Hongjian Qi<sup>31</sup>, Harper Richardson<sup>38</sup>,  
1181 Hana Zaydens<sup>2</sup>, Haicang Zhang<sup>31</sup>, Haoquan Zhao<sup>31</sup>, Ivette Arriaga<sup>17</sup>, Ivy F. Tso<sup>51</sup>, John  
1182 Acampado<sup>2</sup>, Jennifer A. Gerds<sup>50</sup>, Josh Beeson<sup>48</sup>, Jennylyn Brown<sup>2</sup>, Joaquin Comitre<sup>24</sup>, Jeanette  
1183 Cordova<sup>33</sup>, Jennifer Delaporte<sup>19</sup>, Joseph F. Cubells<sup>45</sup>, Jill F. Harris<sup>42</sup>, Jared Gong<sup>25</sup>, Jaclyn  
1184 Gunderson<sup>28</sup>, Jessica Hernandez<sup>2</sup>, Jessyca Judge<sup>23</sup>, Jane Jurayj<sup>21</sup>, Julie Manoharan<sup>2</sup>, Jessie  
1185 Montezuma<sup>38</sup>, Jason Neely<sup>16</sup>, Jessica Orobio<sup>39</sup>, Juhi Pandey<sup>52</sup>, Joseph Piven<sup>43</sup>, Jose Polanco<sup>29</sup>,  
1186 Jibrielle Polite<sup>2</sup>, Jacob Rosewater<sup>24</sup>, Jessica Scherr<sup>26</sup>, James S. Sutcliffe<sup>53</sup>, James T. McCracken<sup>17</sup>,  
1187 Jennifer Tjernagel<sup>2</sup>, Jaimie Toroney<sup>2</sup>, Jeremy Veenstra-Vanderweele<sup>54</sup>, Jiayao Wang<sup>31</sup>, Katie  
1188 Ahlers<sup>50</sup>, Kathryn A. Schweers<sup>13</sup>, Kelli Baalman<sup>39</sup>, Katie Beard<sup>28</sup>, Kristen Callahan<sup>49</sup>, Kendra  
1189 Coleman<sup>34</sup>, Kate D. Fitzgerald<sup>23</sup>, Kate Dent<sup>47</sup>, Katharine Diehl<sup>2</sup>, Kelsey Gonring<sup>48</sup>, Katherine G.  
1190 Pawlowski<sup>35</sup>, Kathy Hirst<sup>19</sup>, Kiely Law<sup>2</sup>, Karen L. Pierce<sup>34</sup>, Karla Murillo<sup>17</sup>, Kailey Murray<sup>43</sup>, Kerri  
1191 Nowell<sup>19</sup>, Kaela O'Brien<sup>29</sup>, Katrina Pama<sup>16</sup>, Kelli Real<sup>43</sup>, Kaitlyn Singer<sup>47</sup>, Kaitlin Smith<sup>41</sup>, Kevin

1192 Stephenson<sup>26</sup>, Katherine Tsai<sup>17</sup>, Leonard Abbeduto<sup>41</sup>, Lindsey A. Cartner<sup>2</sup>, Landon Beeson<sup>32</sup>,  
1193 Laura Carpenter<sup>38</sup>, Lucas Casten<sup>4</sup>, Leigh Coppola<sup>32</sup>, Lisa Cordiero<sup>33</sup>, Lindsey DeMarco<sup>52</sup>, Lillian D.  
1194 Pacheco<sup>32</sup>, Lorena Ferreira Corzo<sup>48</sup>, Lisa H. Shulman<sup>36</sup>, Lauren Kasperson Walsh<sup>47</sup>, Laurie  
1195 Leshner<sup>14</sup>, Lynette M. Herbert<sup>24</sup>, Lisa M. Prock<sup>35</sup>, Lacy Malloch<sup>49</sup>, Lori Mann<sup>2</sup>, Luke P. Grosvenor<sup>2</sup>,  
1196 Laura Simon<sup>28</sup>, Latha V. Soorya<sup>13</sup>, Lucy Wasserburg<sup>28</sup>, Lisa Yeh<sup>13</sup>, Lark Y. Huang-Storms<sup>32</sup>,  
1197 Michael Alessandri<sup>24</sup>, Marc A. Popp<sup>18</sup>, Melissa Baer<sup>48</sup>, Malia Beckwith<sup>42</sup>, Myriam Casseus<sup>42</sup>,  
1198 Michelle Coughlin<sup>35</sup>, Mary Currin<sup>43</sup>, Michele Cutri<sup>24</sup>, Malcolm D. Mallardi<sup>2</sup>, Megan DuBois<sup>28</sup>,  
1199 Megan Dunlevy<sup>46</sup>, Martin E. Butler<sup>2</sup>, Margot Frayne<sup>25</sup>, McLeod F. Gwynette<sup>55</sup>,  
1200 Mohammad Ghaziuddin<sup>23</sup>, Monica Haley<sup>17</sup>, Michelle Heyman<sup>37</sup>, Margaret Hojlo<sup>35</sup>, Michelle  
1201 Jordy<sup>43</sup>, Michael J. Morrier<sup>45</sup>, Misia Kowanda<sup>2</sup>, Melinda Koza<sup>16</sup>, Marilyn Lopez<sup>42</sup>, Megan  
1202 McTaggart<sup>16</sup>, Megan Norris<sup>26</sup>, Melissa N. Hale<sup>24</sup>, Molly O'Neil<sup>36</sup>, Madison Printen<sup>13</sup>, Madelyn  
1203 Rayos<sup>24</sup>, Mahfuza Sabiha<sup>2</sup>, Mustafa Sahin<sup>56</sup>, Marina Sarris<sup>2</sup>, Mojeeb Shir<sup>34</sup>, Matthew Siegel<sup>57</sup>,  
1204 Morgan Steele<sup>25</sup>, Megan Sweeney<sup>19</sup>, Maira Tafolla<sup>17</sup>, Maria Valicenti-McDermott<sup>36</sup>, Mary  
1205 Verdi<sup>57</sup>, Megan Y. Dennis<sup>58</sup>, Nicolas Alvarez<sup>16</sup>, Nicole Bardett<sup>15</sup>, Natalie Berger<sup>13</sup>, Norma  
1206 Calderon<sup>13</sup>, Nickelle Decius<sup>24</sup>, Natalia Gonzalez<sup>42</sup>, Nina Harris<sup>15</sup>, Noah Lawson<sup>2</sup>, Natasha Lillie<sup>28</sup>,  
1207 Nathan Lo<sup>2</sup>, Nancy Long<sup>26</sup>, Nicole M. Russo-Ponsaran<sup>13</sup>, Natalie Madi<sup>29</sup>, Nicole McCoy<sup>29</sup>, Natalie  
1208 Nagpal<sup>2</sup>, Nicki Rodriguez<sup>41</sup>, Nicholas Russell<sup>26</sup>, Neelay Shah<sup>2</sup>, Nicole Takahashi<sup>19</sup>, Nicole  
1209 Targalia<sup>33</sup>, Olivia Newman<sup>28</sup>, Opal Y. Ousley<sup>45</sup>, Peter Heydemann<sup>48</sup>, Patricia Manning<sup>29</sup>, Paul S.  
1210 Carbone<sup>14</sup>, Raphael A. Bernier<sup>50</sup>, Rachel A. Gordon<sup>13</sup>, Rebecca C. Shaffer<sup>29</sup>, Robert D. Annett<sup>49</sup>,  
1211 Renee D. Clark<sup>43</sup>, Roger Jou<sup>21</sup>, Rebecca J. Landa<sup>16</sup>, Rachel K. Earl<sup>50</sup>, Robin Libove<sup>25</sup>, Richard  
1212 Marini<sup>2</sup>, Ryan N. Doan<sup>40</sup>, Robin P. Goin-Kochel<sup>39</sup>, Rishiraj Rana<sup>2</sup>, Richard Remington<sup>2</sup>, Roman  
1213 Shikov<sup>16</sup>, Robert T. Schultz<sup>52</sup>, Shelley Aberle<sup>32</sup>, Shelby Birdwell<sup>19</sup>, Sarah Boland<sup>21</sup>, Stephanie  
1214 Booker<sup>29</sup>, S. Carpenter<sup>15</sup>, Sharmista Chintalapalli<sup>23</sup>, Sarah Conyers<sup>38</sup>, Sophia D'Ambrosi<sup>38</sup>, Sara  
1215 Eldred<sup>26</sup>, Sunday Francis<sup>28</sup>, Swami Ganesan<sup>2</sup>, Susan Hepburn<sup>33</sup>, Susannah Horner<sup>52</sup>, Samantha  
1216 Hunter<sup>19</sup>, Stephanie J. Brewster<sup>59</sup>, Soo J. Lee<sup>13</sup>, Suma Jacob<sup>28</sup>, Stanley Jean<sup>2</sup>, So Hyun Kim<sup>60</sup>,  
1217 Sydney Kramer<sup>4</sup>, Sandra L. Friedman<sup>33</sup>, Sarely Licon<sup>13</sup>, Sandy Littlefield<sup>25</sup>, Stephen M.  
1218 Kanne<sup>19, 61</sup>, Sarah Mastel<sup>32</sup>, Sheena Mathai<sup>46</sup>, Sophia Melnyk<sup>29</sup>, Sarah Michaels<sup>16</sup>, Sarah  
1219 Mohiuddin<sup>23</sup>, Samiza Palmer<sup>52</sup>, Samantha Plate<sup>52</sup>, Shanping Qiu<sup>37</sup>, Shelley Randall<sup>29</sup>, Sophia  
1220 Sandhu<sup>17</sup>, Susan Santangelo<sup>57</sup>, Swapnil Shah<sup>2</sup>, Steve Skinner<sup>62</sup>, Samantha Thompson<sup>41</sup>, Sabrina  
1221 Xiao<sup>2</sup>, Sidi Xu<sup>34</sup>, Sabrina White<sup>49</sup>, Stormi White<sup>46</sup>, Tia Chen<sup>34</sup>, Tunisia Greene<sup>2</sup>, Theodore Ho<sup>50</sup>,  
1222 Teresa Ibanez<sup>26</sup>, Tanner Koomar<sup>4</sup>, Tiziano Pramparo<sup>34</sup>, Tara Rutter<sup>50</sup>, Tamim Shaikh<sup>33</sup>, Taylor  
1223 Thomas<sup>4</sup>, Thao Tran<sup>48</sup>, Timothy W. Yu<sup>40</sup>, Virginia Galbraith<sup>38</sup>, Vahid Gazestani<sup>63</sup>, Vincent J.  
1224 Myers<sup>2</sup>, Vaikunt Ranganathan<sup>52</sup>, Vini Singh<sup>16</sup>, William Curtis Weaver<sup>47</sup>, Wenteng Cai<sup>28</sup>, Wubin  
1225 Chin<sup>2</sup>, Wha S. Yang<sup>17</sup>, YB Choi<sup>60</sup>, Zachary E. Warren<sup>30</sup>

1226 <sup>13</sup>Department of Psychiatry and Behavioral Sciences, Rush University Medical Center, Chicago,  
1227 Illinois 60612; <sup>14</sup>Department of Pediatrics, University of Utah, Salt Lake City, Utah 84108;  
1228 <sup>15</sup>Vanderbilt Kennedy Center, Vanderbilt University Medical Center, Nashville, Tennessee  
1229 37232; <sup>16</sup>Center for Autism and Related Disorders, Kennedy Krieger Institute, Baltimore,  
1230 Maryland 21211; <sup>17</sup>Department of Psychiatry and Biobehavioral Sciences, University of  
1231 California, Los Angeles, Los Angeles, California 90095; <sup>18</sup>PreventionGenetics, Marshfield,  
1232 Wisconsin 54449; <sup>19</sup>Thompson Center for Autism and Neurodevelopmental Disorders,  
1233 University of Missouri, Columbia, Missouri 65211; <sup>20</sup>Department of Pediatrics, University of  
1234 Minnesota, Minneapolis, Minnesota 55414; <sup>21</sup>Child Study Center, Yale School of Medicine, New  
1235 Haven, Connecticut 06519; <sup>22</sup>Department of Neurogenetics, Kennedy Krieger Institute,  
1236 Baltimore, Maryland 21205; <sup>23</sup>Department of Psychiatry, University of Michigan, Ann Arbor,  
1237 Michigan 48109; <sup>24</sup>Department of Psychology, University of Miami's Center for Autism and  
1238 Related Disabilities (UM-CARD), Coral Gables, Florida 33146; <sup>25</sup>Department of Psychiatry and  
1239 Behavioral Sciences, Stanford University, Stanford, California 94305; <sup>26</sup>Division of Pediatric  
1240 Psychology and Neuropsychology, Nationwide Children's Hospital (Child Development Center),  
1241 Columbus, Ohio 43205; <sup>27</sup>John P. Hussman Institute for Human Genomics, University of Miami  
1242 Miller School of Medicine, Miami, Florida 33136; <sup>28</sup>Department of Psychiatry, University of  
1243 Minnesota, Minneapolis, Minnesota 55455; <sup>29</sup>Department of Psychiatry and Behavioral  
1244 Neuroscience, Cincinnati Children's Hospital Medical Center - Research Foundation, Cincinnati,  
1245 Ohio 45229; <sup>30</sup>Department of Pediatrics, Vanderbilt University Medical Center, Nashville,  
1246 Tennessee 37232; <sup>31</sup>Department of Systems Biology, Columbia University Medical Center, New  
1247 York, NY 10032; <sup>32</sup>Department of Psychiatry, Oregon Health & Science University, Portland,  
1248 Oregon 97239; <sup>33</sup>Department of Pediatrics, JFK Partners/University of Colorado School of  
1249 Medicine, Aurora, Colorado 80045; <sup>34</sup>Department of Neurosciences, University of California,  
1250 San Diego and SARRC Phoenix, La Jolla, California 92037; <sup>35</sup>Department of Pediatrics, Boston  
1251 Children's Hospital, Boston, Massachusetts 02115; <sup>36</sup>Department of Pediatrics, Montefiore  
1252 Medical Center and The Albert Einstein College of Medicine, Bronx, New York 10461;  
1253 <sup>37</sup>Department of Psychiatry, Weill Cornell Medicine, White Plains, New York 10605;  
1254 <sup>38</sup>Department of Pediatrics, Medical University of South Carolina, Charleston, South Carolina  
1255 29425; <sup>39</sup>Department of Pediatrics, Texas Children's Hospital (Baylor College of Medicine),  
1256 Houston, Texas 77030; <sup>40</sup>Department of Medicine, Boston Children's Hospital, Boston,  
1257 Massachusetts 02115; <sup>41</sup>MIND Institute and Department of Psychiatry and Behavioral Sciences,  
1258 University of California, Davis, Sacramento, California 95817; <sup>42</sup>Children's Specialized Hospital,  
1259 Toms River, New Jersey 08755; <sup>43</sup>Department of Psychiatry, University of North Carolina (UNC,  
1260 TEACCH, CIDD), Chapel Hill, North Carolina 27599; <sup>44</sup>Department of Molecular and Medical  
1261 Genetics, Oregon Health & Science University, Portland, Oregon 97239; <sup>45</sup>Department of  
1262 Psychiatry and Behavioral Sciences, Emory University and Marcus Autism Center, Atlanta,  
1263 Georgia 30033; <sup>46</sup>Department of Pediatrics, Emory University and Marcus Autism Center,  
1264 Atlanta, Georgia 30329; <sup>47</sup>Geisinger Autism & Developmental Medicine Institute, Lewisburg,  
1265 Pennsylvania 17837; <sup>48</sup>Department of Pediatrics, Rush University Medical Center, Chicago,  
1266 Illinois 60612; <sup>49</sup>Department of Pediatrics, University of Mississippi Medical Center, Jackson,  
1267 Mississippi 39110; <sup>50</sup>Department of Psychiatry and Behavioral Sciences, University of  
1268 Washington/Seattle Children's Autism Center, Seattle, Washington 98195; <sup>51</sup>Department of  
1269 Psychology, University of Michigan, Ann Arbor, Michigan 48109; <sup>52</sup>Center for Autism Research,

1270 Children's Hospital of Philadelphia, Philadelphia, Pennsylvania 19146; <sup>53</sup>Department  
1271 of Molecular Physiology and Biophysics, Vanderbilt University, Nashville, Tennessee 37232;  
1272 <sup>54</sup>Department of Psychiatry, Columbia University Medical Center, New York, NY 10032;  
1273 <sup>55</sup>Department of Psychiatry and Behavioral Sciences, Medical University of South Carolina,  
1274 Charleston, South Carolina 29425; <sup>56</sup>Department of Neurology, Boston Children's Hospital,  
1275 Boston, Massachusetts 02115; <sup>57</sup>Maine Medical Center Research Institute, Scarborough, Maine  
1276 04074; <sup>58</sup>Genome Center, MIND Institute, Department of Biochemistry and Molecular  
1277 Medicine, University of California, Davis, Sacramento, California 95616; <sup>59</sup>Translational  
1278 Neuroscience Center, Boston Children's Hospital, Boston, Massachusetts 02115; <sup>60</sup>Center for  
1279 Autism and the Developing Brain (CADB), Weill Cornell Medicine, White Plains, New York  
1280 10605; <sup>61</sup>Department of Health Psychology, University of Missouri, Columbia, Missouri 65211;  
1281 <sup>62</sup>Greenwood Genetic Center, Greenwood, South Carolina 29646; <sup>63</sup>Department of Pediatrics,  
1282 University of California, San Diego and SARRC Phoenix, La Jolla, California 92037

1283 References

1284

- 1285 1. Lord, C. *et al.* Autism spectrum disorder. *Nat Rev Dis Primers* **6**, 5 (2020).
- 1286 2. Iossifov, I. *et al.* The contribution of de novo coding mutations to autism spectrum  
1287 disorder. *Nature* **515**, 216-21 (2014).
- 1288 3. De Rubeis, S. *et al.* Synaptic, transcriptional and chromatin genes disrupted in autism.  
1289 *Nature* **515**, 209-15 (2014).
- 1290 4. O'Roak, B.J. *et al.* Recurrent de novo mutations implicate novel genes underlying  
1291 simplex autism risk. *Nat Commun* **5**, 5595 (2014).
- 1292 5. Yuen, R.K. *et al.* Whole-genome sequencing of quartet families with autism spectrum  
1293 disorder. *Nat Med* **21**, 185-91 (2015).
- 1294 6. Yuen, R.K. *et al.* Whole genome sequencing resource identifies 18 new candidate genes  
1295 for autism spectrum disorder. *Nat Neurosci* **20**, 602-611 (2017).
- 1296 7. Feliciano, P. *et al.* Exome sequencing of 457 autism families recruited online provides  
1297 evidence for autism risk genes. *NPJ Genom Med* **4**, 19 (2019).
- 1298 8. Satterstrom, F.K. *et al.* Large-Scale Exome Sequencing Study Implicates Both  
1299 Developmental and Functional Changes in the Neurobiology of Autism. *Cell* **180**, 568-  
1300 584 e23 (2020).
- 1301 9. Study, D.D.D. Large-scale discovery of novel genetic causes of developmental disorders.  
1302 *Nature* **519**, 223-8 (2015).
- 1303 10. Study, D.D.D. Prevalence and architecture of de novo mutations in developmental  
1304 disorders. *Nature* **542**, 433-438 (2017).
- 1305 11. Kaplanis, J. *et al.* Evidence for 28 genetic disorders discovered by combining healthcare  
1306 and research data. *Nature* **586**, 757-762 (2020).
- 1307 12. He, X. *et al.* Integrated model of de novo and inherited genetic variants yields greater  
1308 power to identify risk genes. *PLoS Genet* **9**, e1003671 (2013).
- 1309 13. Nguyen, H.T. *et al.* Integrated Bayesian analysis of rare exonic variants to identify risk  
1310 genes for schizophrenia and neurodevelopmental disorders. *Genome Med* **9**, 114 (2017).
- 1311 14. Gaugler, T. *et al.* Most genetic risk for autism resides with common variation. *Nat Genet*  
1312 **46**, 881-5 (2014).
- 1313 15. Sandin, S. *et al.* The familial risk of autism. *JAMA* **311**, 1770-7 (2014).
- 1314 16. Sandin, S. *et al.* The Heritability of Autism Spectrum Disorder. *JAMA* **318**, 1182-1184  
1315 (2017).
- 1316 17. Grove, J. *et al.* Identification of common genetic risk variants for autism spectrum  
1317 disorder. *Nat Genet* **51**, 431-444 (2019).
- 1318 18. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**,  
1319 285-91 (2016).
- 1320 19. Krumm, N. *et al.* Excess of rare, inherited truncating mutations in autism. *Nat Genet* **47**,  
1321 582-8 (2015).
- 1322 20. Kosmicki, J.A. *et al.* Refining the role of de novo protein-truncating variants in  
1323 neurodevelopmental disorders by using population reference samples. *Nat Genet* **49**,  
1324 504-510 (2017).
- 1325 21. Ruzzo, E.K. *et al.* Inherited and De Novo Genetic Risk for Autism Impacts Shared  
1326 Networks. *Cell* **178**, 850-866 e26 (2019).

- 1327 22. Wilfert, A.B. *et al.* Recent ultra-rare inherited variants implicate new autism candidate  
1328 risk genes. *Nat Genet* **53**, 1125-1134 (2021).
- 1329 23. pfeliciano@simonsfoundation.org, S.C.E.a. & Consortium, S. SPARK: A US Cohort of  
1330 50,000 Families to Accelerate Autism Research. *Neuron* **97**, 488-493 (2018).
- 1331 24. MacArthur, D.G. & Tyler-Smith, C. Loss-of-function variants in the genomes of healthy  
1332 humans. *Hum Mol Genet* **19**, R125-30 (2010).
- 1333 25. MacArthur, D.G. *et al.* A systematic survey of loss-of-function variants in human protein-  
1334 coding genes. *Science* **335**, 823-8 (2012).
- 1335 26. Karczewski, K.J. *et al.* The mutational constraint spectrum quantified from variation in  
1336 141,456 humans. *Nature* **581**, 434-443 (2020).
- 1337 27. Cummings, B.B. *et al.* Transcript expression-aware annotation improves rare variant  
1338 interpretation. *Nature* **581**, 452-458 (2020).
- 1339 28. Cassa, C.A. *et al.* Estimating the selective effects of heterozygous protein-truncating  
1340 variants from human exome data. *Nat Genet* **49**, 806-810 (2017).
- 1341 29. Fuller, Z.L., Berg, J.J., Mostafavi, H., Sella, G. & Przeworski, M. Measuring intolerance to  
1342 mutation in human genetics. *Nat Genet* **51**, 772-776 (2019).
- 1343 30. An, J.Y. *et al.* Genome-wide de novo risk score implicates promoter variation in autism  
1344 spectrum disorder. *Science* **362**(2018).
- 1345 31. Ioannidis, N.M. *et al.* REVEL: An Ensemble Method for Predicting the Pathogenicity of  
1346 Rare Missense Variants. *Am J Hum Genet* **99**, 877-885 (2016).
- 1347 32. Fromer, M. *et al.* De novo mutations in schizophrenia implicate synaptic networks.  
1348 *Nature* **506**, 179-84 (2014).
- 1349 33. Parikshak, N.N. *et al.* Integrative functional genomic analyses implicate specific  
1350 molecular pathways and circuits in autism. *Cell* **155**, 1008-21 (2013).
- 1351 34. Weyn-Vanhentenryck, S.M. *et al.* HITS-CLIP and integrative modeling define the Rbfox  
1352 splicing-regulatory network linked to brain development and autism. *Cell Rep* **6**, 1139-  
1353 1152 (2014).
- 1354 35. Darnell, J.C. *et al.* FMRP stalls ribosomal translocation on mRNAs linked to synaptic  
1355 function and autism. *Cell* **146**, 247-61 (2011).
- 1356 36. Cotney, J. *et al.* The autism-associated chromatin modifier CHD8 regulates other autism  
1357 risk genes during human neurodevelopment. *Nat Commun* **6**, 6404 (2015).
- 1358 37. Brueggeman, L., Koomar, T. & Michaelson, J.J. Forecasting risk gene discovery in autism  
1359 with machine learning and genome-scale data. *Sci Rep* **10**, 4569 (2020).
- 1360 38. Chen, S. *et al.* Dissecting Autism Genetic Risk Using Single-cell RNA-seq Data. *bioRxiv*,  
1361 2020.06.15.153031 (2020).
- 1362 39. Ewens, W.J. & Spielman, R.S. The transmission/disequilibrium test: history, subdivision,  
1363 and admixture. *Am J Hum Genet* **57**, 455-64 (1995).
- 1364 40. Fisher, R.A. *Statistical methods for research workers, 11th ed. rev.* (Edinburgh, Oliver and  
1365 Boyd, 1925).
- 1366 41. Zuk, O. *et al.* Searching for missing heritability: designing rare variant association  
1367 studies. *Proc Natl Acad Sci U S A* **111**, E455-64 (2014).
- 1368 42. Wright, S. Evolution in Mendelian Populations. *Genetics* **16**, 97-159 (1931).

- 1369 43. Skene, N.G. & Grant, S.G. Identification of Vulnerable Cell Types in Major Brain Disorders  
1370 Using Single Cell Transcriptomes and Expression Weighted Cell Type Enrichment. *Front*  
1371 *Neurosci* **10**, 16 (2016).
- 1372 44. Willsey, A.J. *et al.* Coexpression networks implicate human midfetal deep cortical  
1373 projection neurons in the pathogenesis of autism. *Cell* **155**, 997-1007 (2013).
- 1374 45. Wang, S. *et al.* De Novo Sequence and Copy Number Variants Are Strongly Associated  
1375 with Tourette Disorder and Implicate Cell Polarity in Pathogenesis. *Cell Rep* **24**, 3441-  
1376 3454 e12 (2018).
- 1377 46. Heyne, H.O. *et al.* De novo variants in neurodevelopmental disorders with epilepsy. *Nat*  
1378 *Genet* **50**, 1048-1053 (2018).
- 1379 47. Singh, T., Neale, B.M., Daly, M.J. & Consortium, o.b.o.t.S.E.M.-A. Exome sequencing  
1380 identifies rare coding variants in 10 genes which confer substantial risk for  
1381 schizophrenia. *medRxiv*, 2020.09.18.20192815 (2020).
- 1382 48. Shu, C., Snyder, L.G., Shen, Y., Chung, W.K. & Consortium, o.b.o.t.S. Imputing cognitive  
1383 impairment in SPARK, a large autism cohort. *medRxiv*, 2021.08.25.21262613 (2021).
- 1384 49. Arnheim, N. & Calabrese, P. Understanding what determines the frequency and pattern  
1385 of human germline mutations. *Nat Rev Genet* **10**, 478-88 (2009).
- 1386 50. Rees, E., Moskvina, V., Owen, M.J., O'Donovan, M.C. & Kirov, G. De novo rates and  
1387 selection of schizophrenia-associated copy number variants. *Biol Psychiatry* **70**, 1109-14  
1388 (2011).
- 1389 51. Raychaudhuri, S. *et al.* Accurately assessing the risk of schizophrenia conferred by rare  
1390 copy-number variation affecting genes with brain function. *PLoS Genet* **6**, e1001097  
1391 (2010).
- 1392 52. Carlson, J. *et al.* Extremely rare variants reveal patterns of germline mutation rate  
1393 heterogeneity in humans. *Nat Commun* **9**, 3753 (2018).
- 1394 53. Samocha, K.E. *et al.* A framework for the interpretation of de novo mutation in human  
1395 disease. *Nat Genet* **46**, 944-50 (2014).
- 1396 54. Spielman, R.S. & Ewens, W.J. The TDT and other family-based tests for linkage  
1397 disequilibrium and association. *Am J Hum Genet* **59**, 983-9 (1996).
- 1398 55. Chen, W.M., Manichaikul, A. & Rich, S.S. A generalized family-based association test for  
1399 dichotomous traits. *Am J Hum Genet* **85**, 364-76 (2009).
- 1400 56. Staples, J. *et al.* PRIMUS: rapid reconstruction of pedigrees from genome-wide estimates  
1401 of identity by descent. *Am J Hum Genet* **95**, 553-64 (2014).
- 1402 57. Genomes Project, C. *et al.* A global reference for human genetic variation. *Nature* **526**,  
1403 68-74 (2015).
- 1404 58. Li, J.Z. *et al.* Worldwide human relationships inferred from genome-wide patterns of  
1405 variation. *Science* **319**, 1100-4 (2008).
- 1406 59. Bergstrom, A. *et al.* Insights into human genetic variation and population history from  
1407 929 diverse genomes. *Science* **367**(2020).
- 1408 60. Patterson, N., Price, A.L. & Reich, D. Population structure and eigenanalysis. *PLoS Genet*  
1409 **2**, e190 (2006).
- 1410 61. Chen, C.Y. *et al.* Improved ancestry inference using weights from external reference  
1411 panels. *Bioinformatics* **29**, 1399-406 (2013).

- 1412 62. Bishop, S.L., Farmer, C. & Thurm, A. Measurement of nonverbal IQ in autism spectrum  
1413 disorder: scores in young adulthood compared to early childhood. *Journal of autism and*  
1414 *developmental disorders* **45**, 966-974 (2015).
- 1415 63. Munson, J. *et al.* Evidence for latent classes of IQ in young children with autism  
1416 spectrum disorder. *American journal of mental retardation : AJMR* **113**, 439-452 (2008).
- 1417 64. Fischbach, G.D. & Lord, C. The Simons Simplex Collection: a resource for identification of  
1418 autism genetic risk factors. *Neuron* **68**, 192-5 (2010).
- 1419 65. Werling, D.M. *et al.* An analytical framework for whole-genome sequence association  
1420 studies and its implications for autism spectrum disorder. *Nat Genet* **50**, 727-736 (2018).
- 1421 66. Buxbaum, J.D. *et al.* The autism sequencing consortium: large-scale, high-throughput  
1422 sequencing in autism spectrum disorders. *Neuron* **76**, 1052-6 (2012).
- 1423 67. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.  
1424 *arXiv* (2013).
- 1425 68. Fritz, M.H.Y., Leinonen, R., Cochrane, G. & Birney, E. Efficient storage of high throughput  
1426 DNA sequencing data using reference-based compression. *Genome Res* **21**, 734-740  
1427 (2011).
- 1428 69. Institute, P.T.-B.B. <https://broadinstitute.github.io/picard/>.
- 1429 70. Pedersen, B.S. & Quinlan, A.R. Mosdepth: quick coverage calculation for genomes and  
1430 exomes. *Bioinformatics* **34**, 867-868 (2018).
- 1431 71. Jun, G. *et al.* Detecting and estimating contamination of human DNA samples in  
1432 sequencing and array-based genotype data. *Am J Hum Genet* **91**, 839-48 (2012).
- 1433 72. Poplin, R. *et al.* Scaling accurate genetic variant discovery to tens of thousands of  
1434 samples. *bioRxiv*, 201178 (2018).
- 1435 73. GitHub - Genomicsplc/wecall: Fast, accurate and simple to use command line tool for  
1436 variant detection in NGS data. <https://github.com/Genomicsplc/wecall>.
- 1437 74. Poplin, R. *et al.* A universal SNP and small-indel variant caller using deep neural  
1438 networks. *Nat Biotechnol* **36**, 983-987 (2018).
- 1439 75. Lin, M.F. *et al.* GLnexus: joint variant calling for large cohort sequencing. *bioRxiv*, 343970  
1440 (2018).
- 1441 76. Garrison, E. & Marth, G. Haplotype-based variant detection from short-read sequencing.  
1442 *arXiv* (2012).
- 1443 77. Martin, M. *et al.* WhatsHap: fast and accurate read-based phasing. *bioRxiv*, 085050  
1444 (2016).
- 1445 78. Tan, A., Abecasis, G.R. & Kang, H.M. Unified representation of genetic variants.  
1446 *Bioinformatics* **31**, 2202-4 (2015).
- 1447 79. Frankish, A. *et al.* GENCODE reference annotation for the human and mouse genomes.  
1448 *Nucleic Acids Res* **47**, D766-d773 (2019).
- 1449 80. Autism Spectrum Disorders Working Group of The Psychiatric Genomics, C. Meta-  
1450 analysis of GWAS of over 16,000 individuals with autism spectrum disorder highlights a  
1451 novel locus at 10q24.32 and a significant overlap with schizophrenia. *Molecular autism*  
1452 **8**, 21-21 (2017).
- 1453 81. Lindsay, S.J. *et al.* HDBR Expression: A Unique Resource for Global and Individual Gene  
1454 Expression Studies during Early Human Brain Development. *Front Neuroanat* **10**, 86  
1455 (2016).

- 1456 82. Kircher, M. *et al.* A general framework for estimating the relative pathogenicity of  
1457 human genetic variants. *Nat Genet* **46**, 310-5 (2014).
- 1458 83. Samocha, K.E. *et al.* Regional missense constraint improves variant deleteriousness  
1459 prediction. *bioRxiv*, 148353 (2017).
- 1460 84. Sundaram, L. *et al.* Predicting the clinical impact of human mutation with deep neural  
1461 networks. *Nat Genet* **50**, 1161-1170 (2018).
- 1462 85. Jaganathan, K. *et al.* Predicting Splicing from Primary Sequence with Deep Learning. *Cell*  
1463 **176**, 535-548 e24 (2019).
- 1464 86. Packer, J.S. *et al.* CLAMMS: a scalable algorithm for calling common and rare copy  
1465 number variants from exome sequencing data. *Bioinformatics* **32**, 133-5 (2016).
- 1466 87. Koehler, R., Issac, H., Cloonan, N. & Grimmond, S.M. The uniqueome: a mappability  
1467 resource for short-tag sequencing. *Bioinformatics* **27**, 272-4 (2011).
- 1468 88. Taliun, D. *et al.* Sequencing of 53,831 diverse genomes from the NHLBI TOPMed  
1469 Program. *Nature* **590**, 290-299 (2021).
- 1470 89. Li, H. Toward better understanding of artifacts in variant calling from high-coverage  
1471 samples. *Bioinformatics* **30**, 2843-51 (2014).
- 1472 90. Maenner, M.J. *et al.* Prevalence of Autism Spectrum Disorder Among Children Aged 8  
1473 Years — Autism and Developmental Disabilities Monitoring Network, 11 Sites, United  
1474 States, 2016. *MMWR Surveill Summ* **69**, 1-12 (2020).
- 1475 91. Power, R.A. *et al.* Fecundity of patients with schizophrenia, autism, bipolar disorder,  
1476 depression, anorexia nervosa, or substance abuse vs their unaffected siblings. *JAMA*  
1477 *Psychiatry* **70**, 22-30 (2013).
- 1478 92. Wright, C.F. *et al.* Genetic diagnosis of developmental disorders in the DDD study: a  
1479 scalable analysis of genome-wide research data. *Lancet* **385**, 1305-14 (2015).
- 1480 93. Abrahams, B.S. *et al.* SFARI Gene 2.0: a community-driven knowledgebase for the autism  
1481 spectrum disorders (ASDs). *Mol Autism* **4**, 36 (2013).
- 1482 94. Singer, E. What Makes an Autism Gene? . (SPARK for Autism).
- 1483 95. Fagerberg, L. *et al.* Analysis of the human tissue-specific expression by genome-wide  
1484 integration of transcriptomics and antibody-based proteomics. *Mol Cell Proteomics* **13**,  
1485 397-406 (2014).
- 1486 96. Mo, A. *et al.* Epigenomic Signatures of Neuronal Diversity in the Mammalian Brain.  
1487 *Neuron* **86**, 1369-84 (2015).
- 1488 97. Pirooznia, M. *et al.* SynaptomeDB: an ontology-based knowledgebase for synaptic  
1489 genes. *Bioinformatics* **28**, 897-9 (2012).
- 1490 98. Wagnon, J.L. *et al.* CELF4 regulates translation and local abundance of a vast set of  
1491 mRNAs, including genes associated with regulation of synaptic function. *PLoS Genet* **8**,  
1492 e1003067 (2012).
- 1493 99. Eppig, J.T. Mouse Genome Informatics (MGI) Resource: Genetic, Genomic, and  
1494 Biological Knowledgebase for the Laboratory Mouse. *Ilar j* **58**, 17-41 (2017).
- 1495 100. Szklarczyk, D. *et al.* STRING v11: protein-protein association networks with increased  
1496 coverage, supporting functional discovery in genome-wide experimental datasets.  
1497 *Nucleic Acids Res* **47**, D607-d613 (2019).
- 1498 101. Köhler, S. *et al.* The Human Phenotype Ontology in 2021. *Nucleic Acids Res* **49**, D1207-  
1499 d1217 (2021).

- 1500 102. Cutler, A. & Breiman, L. Archetypal Analysis. *Technometrics* **36**, 338-347 (1994).  
1501 103. Thorndike, R.L. Who belongs in the family? *Psychometrika* **18**, 267-276 (1953).  
1502