

1 **Oligogenic combinations of rare variants influence specific phenotypes in**  
2 **complex disorders**

3

4 Vijay Kumar Pounraja<sup>1,2</sup> and Santhosh Girirajan<sup>1,2,3</sup>

5

6 1. Department of Biochemistry and Molecular Biology, Pennsylvania State University,  
7 University Park, PA 16802

8 2. Bioinformatics and Genomics Graduate Program, The Huck Institute of the Life  
9 Sciences, University Park, PA 16802

10 3. Department of Anthropology, Pennsylvania State University, University Park, PA 16802

11

12

13 **Correspondence:**

14 Santhosh Girirajan

15 205A Life Sciences Building

16 Pennsylvania State University

17 University Park, PA 16802

18 E-mail: [sxg47@psu.edu](mailto:sxg47@psu.edu)

19

20

21 **Classification:** Biological Sciences

22 **Keywords:** oligogenic, complex disorders, autism, rare variants

## 23 **ABSTRACT**

24 Genetic studies of complex disorders such as autism and intellectual disability (ID) are often  
25 based on enrichment of individual rare variants or their aggregate burden in affected individuals  
26 compared to controls. However, these studies overlook the influence of combinations of rare  
27 variants that may not be deleterious on their own due to statistical challenges resulting from  
28 rarity and combinatorial explosion when enumerating variant combinations, limiting our ability  
29 to study oligogenic basis for these disorders. We present a framework that combines the apriori  
30 algorithm and statistical inference to identify specific combinations of mutated genes associated  
31 with complex phenotypes. Our approach overcomes computational barriers and exhaustively  
32 evaluates variant combinations to identify non-additive relationships between simultaneously  
33 mutated genes. Using this approach, we analyzed 6,189 individuals with autism and identified  
34 718 combinations significantly associated with ID, and carriers of these combinations showed  
35 lower IQ than expected in an independent cohort of 1,878 individuals. These combinations were  
36 enriched for nervous system genes such as *NIN* and *NGF*, showed complex inheritance patterns,  
37 and were depleted in unaffected siblings. We found that an affected individual can carry many  
38 oligogenic combinations, each contributing to the same phenotype or distinct phenotypes at  
39 varying effect sizes. We also used this framework to identify combinations associated with  
40 multiple comorbid phenotypes, including mutations of *COL28A1* and *MFSD2B* for ID and  
41 schizophrenia and *ABCA4*, *DNAH10* and *MC1R* for ID and anxiety/depression. Our framework  
42 identifies a key component of missing heritability and provides a novel paradigm to untangle the  
43 genetic architecture of complex disorders.

44

## 45 **SIGNIFICANCE**

46 While rare mutations in single genes or their collective burden partially explain the genetic basis  
47 for complex disorders, the role of specific combinations of rare variants is not completely  
48 understood. This is because combinations of rare variants are rarer and evaluating all possible  
49 combinations would result in a combinatorial explosion, creating difficulties for statistical and  
50 computational analysis. We developed a data mining approach that overcomes these limitations  
51 to precisely quantify the influence of combinations of two or more mutated genes on a specific  
52 clinical feature or multiple co-occurring features. Our framework provides a new paradigm for

- 53 dissecting the genetic causes of complex disorders and provides an impetus for its utility in  
54 clinical diagnosis.

## 55 INTRODUCTION

56 Recent human population growth has led to a rapid increase in the load of rare variants affecting  
57 functionally important regions of the genome<sup>1-3</sup>. Thus, rare variants are collectively more  
58 abundant in the population compared to common variants, many of which confer significant risk  
59 for neurodevelopmental disorders such as autism and intellectual disability<sup>4</sup>. In fact, recent  
60 studies have directly implicated rare damaging mutations that are very recent or *de novo* in more  
61 than one hundred genes towards neurodevelopmental disorders<sup>5-7</sup>. The ability to establish robust  
62 associations between rare variants of high effect size and complex disease has made this class of  
63 variants the primary focus of recent studies. However, a much larger class of rare and variably  
64 expressive variants that are individually less deleterious but, in combination, exert large effects  
65 towards disease is often overlooked. Variants in this category are often transmitted across  
66 generations without adverse effects on their carriers until they encounter other similar variants  
67 that, when combined, lead to genetic interactions conferring a higher risk for disease than their  
68 individual risks<sup>8,9</sup>. While this phenomenon underpins oligogenic models proposed over the years,  
69 studies so far have not focused on detecting combinatorial effects of specific sets of rare variants  
70 towards disease phenotypes<sup>10-13</sup>.

71 Identifying the effects of specific combinations of rare variants towards disease etiology  
72 has been challenging for many reasons. *First*, combinations of rare variants are rarer, and  
73 extremely large cohorts are required to observe even a few recurrent instances of specific variant  
74 combinations<sup>14</sup>. Prior studies of oligogenic models for rare variants evaded this problem by  
75 aggregating variant information at the sample level and comparing the overall burden of rare  
76 variants between groups of individuals (such as cases and controls)<sup>6,7,15,16</sup>. *Second*, the  
77 combinatorial explosion resulting from even a small set of rare variants makes it difficult to  
78 exhaustively evaluate all combinations. While sophisticated frameworks such as network  
79 analysis and machine learning provide powerful tools to model the composite effects of  
80 thousands of variables on a complex system and predict emergent behaviors and quantitative  
81 outcomes, adapting them to exhaustively search and delineate the effects of specific  
82 combinations of variables is daunting<sup>17,18</sup>. Furthermore, incorporating an efficient search tool  
83 into these frameworks and extending them to detect higher-order combinatorial effects would be  
84 nearly impossible. *Third*, even when all combinations of rare variants could be exhaustively  
85 evaluated within a large cohort, there is a lack of methods that are sensitive enough to detect

86 small differences between comparison groups to establish statistical significance. Therefore, an  
87 alternate approach that is highly flexible, scalable, and sensitive is necessary to address  
88 computational and statistical challenges associated with assessing rare variant combinations.

89 Here, we present a combinatorial framework called *RareComb* that couples the apriori  
90 algorithm<sup>19</sup> with binomial tests to overcome the limitations of data sparsity and high  
91 dimensionality, and systematically analyzes patterns of rare variants between groups of interest  
92 to identify specific combinations that are significantly associated with phenotypes<sup>20</sup>. We apply  
93 our analysis framework to a discovery cohort of 6,189 children with autism to identify genetic  
94 interactions involving pairs and triplets of mutated genes that are enriched in individuals with  
95 intellectual disability compared to individuals without intellectual disability. We demonstrate  
96 that the carriers of mutations in these specific gene pairs and triplets within an independent  
97 cohort of 1,878 children have significantly lower-than-expected intelligence quotient (IQ) scores.  
98 We also demonstrate the adaptability of our framework by leveraging it to identify mutated gene  
99 pairs and triplets significantly associated with two or more comorbid phenotypes among children  
100 with autism. Finally, we show how this generalizable and modular framework can be easily  
101 extended to identify higher order interactions beyond pairs and triplets of variants. Our stand-  
102 alone framework does not depend on *a priori* knowledge and can detect rare patterns from high-  
103 dimensional genetic data to generate interpretable results, making it readily applicable for  
104 analyzing cohorts of all size ranges to dissect the genetic basis of complex disorders.

105

106

107

108

109

110

111

112

## 113 RESULTS

114 We hypothesized that two or more genes disrupted simultaneously by rare deleterious mutations  
115 contribute to a highly penetrant phenotype, as in an oligogenic model, or lead to a more severe  
116 phenotype than when each of the same genes are disrupted individually. We developed  
117 *RareComb* as a framework that combines data mining and statistical analysis to identify specific  
118 combinations (such as pairs, triplets, etc.) of rare variants that show significant associations with  
119 one or more phenotypes. *RareComb* analyzes an ' $n \times p$ ' sparse Boolean matrix with ' $p$ ' genes in  
120 ' $n$ ' individuals in three discrete steps (**Figure 1**). *First*, it applies the apriori algorithm  
121 independently in cases and controls to enumerate the frequency of all simultaneously mutated  
122 combinations that meet a pre-set minimum frequency threshold (**Supp. Figure 1**). *Second*, for  
123 each qualifying combination of variants, the method derives the expected frequency of  
124 simultaneously observing mutations in the constituent genes under the assumption of  
125 independence. It then independently quantifies the magnitude of deviation of the observed from  
126 the expected frequencies using binomial tests in cases and controls, and uses multiple-testing  
127 adjusted p-values to identify combinations that are statistically enriched in cases but not in  
128 controls. *Finally*, the method calculates effect sizes using Cohen's d and statistical power at 1%  
129 and 5% significance thresholds, to enable prioritization of a high confidence set of combinations  
130 that contribute to the phenotype in an oligogenic manner.

131

### 132 ***RareComb* identifies oligogenic combinations associated with ID and autism**

133 We sought to identify pairs and triplets of mutated genes that are significantly associated with  
134 intellectual disability (ID) phenotypes by analyzing 6,189 affected individuals from the  
135 SPARK<sup>21</sup> cohort for discovery and 1,878 affected individuals from the SSC<sup>22</sup> cohort for  
136 validation. To facilitate cross-cohort comparison, we identified 10,217 rare variants ( $MAF \leq 1\%$ )  
137 that were predicted to be deleterious by multiple methods and observed in both cohorts, and  
138 aggregated these variants to genes for the analysis (see **Methods**). We first categorized 1,215  
139 probands from the SPARK cohort diagnosed with ID as cases and 4,974 probands without ID as  
140 controls (**Figure 2A**). We then applied *RareComb* to cases after constraining it to only evaluate  
141 those gene combinations in which simultaneous mutations are observed in at least five probands.  
142 We identified 25,602 pairs involving 1,956 mutated genes in cases that were observed at a higher  
143 frequency than expected under the assumption of independence. Similarly, analyzing the controls

144 using only the 1,956 genes mutated in cases, *RareComb* identified 148 pairs of mutated genes  
145 that were significantly enriched in cases but not in controls (**Supp. Table 1**), with moderate to  
146 high effect sizes (Cohen's  $d$ , 0.08-0.15) and adequate statistical power (70%-100% at 5%  
147 significance threshold) (**Supp. Figure 2**). These 148 gene pairs belonged to 142 probands, with  
148 74% (105/142) of them carrying more than one significant pair. These observations suggest that  
149 an individual can carry multiple combinations, each contributing to the same phenotype at  
150 varying effect sizes (**Supp. Figure 3**).

151 We next sought to validate the association of these 148 mutated gene pairs towards  
152 intellectual disability. We hypothesized that if the association of the gene pairs with ID in the  
153 SPARK cohort were truly significant, carriers of mutations in those gene pairs would tend to  
154 have lower than average IQ scores in the independent SSC cohort. We found that 90 of the 148  
155 significant pairs identified in the SPARK cohort were observed in at least one proband in the  
156 SSC cohort. These 90 mutated gene pairs were carried by 91 unique probands, whose average  
157 full-scale IQ scores (average IQ=68.52) were lower than those of all ascertained probands in the  
158 SSC cohort (average IQ=86). To assess the significance of this result, we performed 10,000  
159 random draws of 91 probands from the SSC cohort to generate a simulated distribution of their  
160 average IQ scores. The average IQ of carriers of mutated gene pairs (average IQ=68.52) was  
161 significantly lower than the overall distribution of average IQ derived from simulations (average  
162 IQ ranged from 73 to 92; empirical  $p=0$ ) (**Figure 2B**). Furthermore, the average IQ of the 91  
163 SSC probands with both mutated genes was significantly lower than the average IQ of 1,252  
164 carriers of mutations in only one of the two genes (68.5 versus 82.8; Kolmogorov-Smirnov  $p =$   
165  $1.302 \times 10^{-16}$ ) (**Figure 2C**). When each of the 90 combinations was evaluated individually,  
166 carriers of mutations in both genes for 73% (66/90) of the combinations showed lower IQ than  
167 individuals with mutations in individual genes of the same combination, with 39/90 remaining  
168 significant after multiple testing correction (**Supp. Table 2; Supp. Figure 4**). These results  
169 provide evidence for synergistic effects of deleterious mutations within specific pairs of genes  
170 towards ID phenotypes.

171 We also applied *RareComb* to identify gene triplets associated with intellectual disability  
172 using the two cohorts and repeated the simulations to identify 1,593 significant combinations in  
173 the SPARK cohort. We selected 570 high-confidence triplets (with  $\geq 90\%$  statistical power at 5%  
174 significance threshold; **Supp. Table 3**) and found that 79 probands in the SSC cohort carried at

175 least one of these deleterious triplets. The average IQ score of individuals carrying significant  
176 gene triplets (average IQ score=73) was significantly lower than a distribution of average IQ  
177 scores from 10,000 draws of 79 SSC probands (average IQ score=82.5; min=72, max=94;  
178 empirical  $p=0.0011$ ; see **Supp. Figure 5**). This result reiterated that carriers of mutations in the  
179 significant gene combinations have lower IQ than a random group of probands. Our results also  
180 demonstrate the ability of the framework to identify higher order combinations of mutations that  
181 are significantly associated with specific phenotypes in individuals with complex disorders.

182

### 183 **Oligogenic combinations are enriched for specific inheritance patterns**

184 As individual variants can arise *de novo* or be inherited maternally or paternally, variants in pairs  
185 of genes can have six possible patterns of transmission (**Supp. Figure 6A**). We identified a total  
186 of 926 occurrences of the 148 pairs of mutated genes enriched among SPARK probands with ID  
187 ( $n=142$  probands), of which inheritance could be determined without ambiguity for 887  
188 instances. We found that one variant occurred *de novo* and the other variant was inherited from  
189 the mother in 244/887 instances (27.5%). Similarly, both mutated genes were inherited from the  
190 mother in 226/887 instances (25.4%) or occurred *de novo* in 221/887 instances (24.9%), while  
191 the remaining fraction (~22%) of variant pairs were either inherited from both parents, inherited  
192 from the father, or transmitted *de novo* and paternally. To assess the significance of our  
193 observations, we performed simulations to establish a baseline expectation of proportions for  
194 each category of parental inheritance pattern. We selected 926 pairs of genes in 1000 random  
195 draws of all possible mutated gene pairs among SPARK probands and calculated the fraction of  
196 instances that fell into each of the six transmission categories. The observed proportion was  
197 higher than the simulated proportions for instances when both variants occurred *de novo* (24.9%  
198 versus 17%, empirical  $p=0$ ) and when one variant was *de novo* and the other was inherited  
199 maternally (27.5% versus 25%,  $p=0.028$ ) (**Figure 3A**). We repeated this analysis for 7,596  
200 children affected with autism in the SPARK cohort compared to 11,740 unaffected parents and  
201 identified 110 gene pairs significantly associated with autism (**Supp. Table 4**). Similar to the  
202 results obtained for the ID phenotype, we found that both variants of a gene pair were more  
203 likely to occur *de novo* (24% versus 18%, empirical  $p=0$ ) or one variant occurring *de novo* and  
204 the other inherited maternally (33% versus 26%,  $p=0$ ) than expected based on simulation studies  
205 (**Supp. Figure 7**). The enrichment of *de novo* or maternally inherited variants for significant



206 gene pairs aligns with published reports that severely affected children tend to carry multiple *de*  
207 *novo* mutations or inherit pathogenic rare variants from mildly affected or unaffected carrier  
208 mothers<sup>16,23,24</sup>.

209 We then assessed whether the mutated gene pairs associated with ID were also found in  
210 siblings of carrier probands. Restricting our analysis to families with unaffected siblings whose  
211 probands had mutations in ID-enriched gene pairs, we found that both variants were present in  
212 the corresponding sibling for only 53/219 (24.2%) instances of gene pairs, while 102/219  
213 (46.6%) had variants in only one of the two genes and 64/219 (29.2%) instances had no variants  
214 in either of the genes in the siblings (**Supp. Figure 6B**). Using simulations, we found a  
215 significantly higher proportion of instances with only one of the two variants present in siblings  
216 compared to the expected values (46.6% versus 38.5%,  $p=0.007$ ). Furthermore, the proportion of  
217 observed instances with neither of the variants present in siblings (29.2% versus 33.1%,  
218 empirical  $p=0.098$ ) or both variants present in siblings (24.2% versus 28.4%,  $p=0.079$ ) was  
219 lower than expected (**Figure 3B**). The observation that only a small fraction of unaffected  
220 siblings carried both mutated gene pairs suggests a strong association of these gene pairs with ID  
221 phenotypes. These results suggest that mutations in pairs of genes significantly associated with a  
222 severe phenotype in probands are more likely to occur individually than simultaneously in  
223 unaffected siblings of the same family.

224

### 225 **Genes forming oligogenic combinations are distinct from canonical autism genes**

226 We expanded our analysis to include all 16,556 mutated genes in the SPARK cohort, as opposed  
227 to genes with mutations present in both the SPARK and SSC cohorts, and identified 52  
228 significant gene pairs (**Supp. Table 5**) and 230 triplets associated with the ID phenotype (with  
229  $\geq 90\%$  statistical power at 1% significance threshold; **Supp. Table 6**). Due to the expanded  
230 search space, the mutated gene pairs showed more significant p-values from the binomial tests  
231 when compared to those obtained from the more restricted set of variants overlapping both  
232 SPARK and SSC cohorts (**Supp. Figure 8**). Mutated genes within these combinations included  
233 several genes related to nervous system development, such as *NIN*, *HDC*, *NGF*, and *BRD8*.  
234 Furthermore, 5/52 pairs and 59/230 triplets contained at least one gene associated with autism in  
235 the SFARI database, including *FGFR1*, associated with multiple disorders including Kallmann  
236 syndrome<sup>25</sup> and Pfeiffer syndrome<sup>26</sup>; *RELN*, associated with temporal lobe epilepsy<sup>27</sup>; *SYNE1*,

237 associated with spinocerebellar ataxia<sup>28,29</sup>; and *PNPLA7*, associated with autism and ID<sup>30</sup>. Thus,  
238 most genes forming combinations are not involved in canonical autism or ID disorders,  
239 suggesting synergistic effects of these genes without prior association to disease.

240 We also performed gene ontology enrichment analysis for genes within the combinations  
241 and identified seven out of nine significantly enriched GO terms to be exclusively associated  
242 with nervous system-related functions, including synthesis and metabolism of catecholamines,  
243 axon/neuron regeneration, and neuron generation and differentiation (**Supp. Figure 9**)<sup>31</sup>.  
244 Furthermore, the differences in the type and specificity of GO terms enriched for significant  
245 pairs versus triplets were apparent, with genes forming pairs involved in nervous system function  
246 and genes forming triplets associated with both nervous system as well as other biological  
247 processes. We next assessed the enrichment and depletion of Human Phenotype Ontology (HPO)  
248 terms for genes forming significant pairs towards ID phenotypes<sup>32</sup>. *First*, we calculated the  
249 fraction of all 4,484 genes within the HPO database associated with each HPO term. *For*  
250 *example*, 30% (1,366/4,484) of all genes in HPO were associated with ID. We compared these  
251 expected values calculated for each HPO term with the corresponding fractions observed within  
252 the 95 genes forming 52 ID-associated pairs using binomial tests. Interestingly, genes associated  
253 with HPO terms related to neurodevelopmental phenotypes, such as ID, global developmental  
254 delay, seizure, and microcephaly, were significantly depleted within the set of 95 genes forming  
255 gene pairs (**Supp. Table 7**). *Next*, we evaluated whether genes within each of the 52 significant  
256 pairs shared one or more common HPO phenotype or disease. Of the 52 pairs, only one pair  
257 (*DNASE1* & *MTR*) shared an HPO phenotype (“epilepsy”). This was significantly lower than the  
258 expected value obtained from the distribution of the number of shared HPO phenotypes between  
259 all possible pairs of genes in the HPO database (1/52, 1.9% ID gene pairs compared to 31.5% of  
260 all HPO gene pairs shared one HPO phenotype,  $p=2.2\times 10^{-16}$ ; one-sided binomial test) (**Supp.**  
261 **Figure 10; Supp. Table 8**). We note that the 4,484 genes within HPO are potentially biased  
262 towards well-studied disorders, making pairs of genes drawn from HPO more likely to share  
263 phenotypes than random pairs of genes from the genome. Overall, GO and HPO analyses show  
264 that genes forming oligogenic combinations are involved in neuronal processes but have not  
265 been previously connected to neurodevelopmental phenotypes, indicating the novelty of the  
266 associations between these genes and ID phenotypes.

267

268

## 269 **Identifying variant combinations towards specific patterns of comorbid phenotypes**

270 We adapted our framework to identify significant associations of two or more genotypes with  
271 multiple comorbid phenotypes. To identify novel comorbid associations, we eliminated  
272 phenotypes that were highly correlated with each other, such as ADHD and reading disorder<sup>33</sup>.  
273 We analyzed variant profiles of 6,189 autism probands from the SPARK cohort with records of  
274 comorbid features, including 1,215 individuals with ID, 1,825 with anxiety and depression, and  
275 332 with schizophrenia features. We assessed for significant co-occurrences of two or more  
276 mutated genes with two or more of the above phenotypes (**Figure 4**). Using one-tailed binomial  
277 tests to compare the observed frequency of combinations of genotypes and phenotypes to the  
278 expected frequency, we first identified 169 significant associations between pairs of mutated  
279 genes and two comorbid phenotypes as well as 82 combinations of three mutated genes and two  
280 comorbid phenotypes (**Supp. Tables 9 & 10**). As some of these significant genotype-phenotype  
281 combinations can be confounded by high degree of co-occurrence of mutated genes, we next  
282 calculated genotype-only p-values using binomial tests for all significant genotype-phenotype  
283 associations. For 32/169 combinations of two mutated genes and two comorbid phenotypes and  
284 5/82 combinations of three mutated genes and two comorbid phenotypes, the composite  
285 genotype-phenotype p-values were significant while genotype-only p-values were not  
286 significant, suggesting stronger associations between these variant combinations and phenotypes.  
287 *For example*, even when variants in genes *COL28A1* and *MFSD2B* did not co-occur more  
288 frequently than expected under the assumption of independence, these mutated genes co-  
289 occurred more frequently than expected among probands with ID and schizophrenia phenotypes.  
290 Loss-of-function and rare missense mutations in *COL28A1* have been reported in individuals  
291 with autism<sup>34,35</sup>, and *MFSD2A*, a paralog of *MFSD2B*, has been directly implicated in an  
292 autosomal recessive disorder associated with progressive microcephaly, spasticity and brain  
293 imaging abnormalities<sup>36</sup>. Similarly, we found *ARVCF* and *FATI* to be significantly associated  
294 with ID and schizophrenia, with *ARVCF* mapping within the 22q11.2 DiGeorge syndrome  
295 region<sup>37</sup>, while rare *de novo* mutations in *FATI* being associated with autism and  
296 schizophrenia<sup>6,38</sup>. Finally, we found that the mutations in genes *ABCA4*, *DNAH10* and *MC1R*  
297 significantly co-occurred in individuals with ID and anxiety/depression phenotypes. These

298 results demonstrate the utility of identifying higher-order associations between genotypes and  
299 phenotypes in complex disorders such as autism.

300

## 301 **DISCUSSION**

302 Current rare variant analysis strategies are geared towards either searching for individual variants  
303 of high effect size whose influence on the phenotype is evident, such as *de-novo* gene-disruptive  
304 mutations, or comparing rare variant burden to explain collective effects on phenotypes<sup>7,39,40</sup>.

305 The wider space between these two extremes of the analysis spectrum that involves  
306 combinations of rare variants has largely remained understudied. Although digenic diseases and  
307 multi-hit models of complex diseases have been used to provide post-hoc explanations for an  
308 observed phenomenon, they are not equipped to serve as a framework to actively search and  
309 identify rare variant combinations that fit oligogenic models for specific phenotypes<sup>9,12,13</sup>. While  
310 machine learning has become the de-facto approach for disease outcome predictions, the lack of  
311 holy-grail predictors and reduced interpretability due to data sparsity makes it less fit to detect  
312 combinatorial effects<sup>17</sup>. In addition, the common practice of evaluating feature importance  
313 metrics of machine learning classifiers falls short of the objective to identify combinations of  
314 features that exert higher effect on the phenotype than evident from their independent effects<sup>17,18</sup>.  
315 Furthermore, prior studies to assess combinatorial effects have been inherently biased due to  
316 their need to minimize the search space by restricting the analysis to only a subset of genes  
317 chosen based on *a priori* knowledge<sup>41-43</sup>. Here, we provide a proof-of-concept analytical  
318 framework that remains agnostic to prior evidence and performs exhaustive searches to identify  
319 combinatorial effects among rare variants while retaining high granularity of data and  
320 interpretability of results.

321 We use our framework to identify gene pairs and triplets significantly associated with  
322 intellectual disability and show that several constituent genes are associated with nervous system  
323 processes. These mutated gene combinations are more likely to be inherited maternally or occur  
324 *de novo*, are depleted in unaffected siblings from the same family, and are less likely to involve  
325 canonical autism or ID genes, suggesting that genes forming significant combinations are less  
326 deleterious on their own but manifest effects only when combined with other similar genes  
327 carrying rare mutations. While previous studies have linked aggregate rare variant burden  
328 towards intellectual disability<sup>44,45</sup>, our results fine map the association to specific combinations

329 of constituent genes contributing to the burden. We propose a novel paradigm for dissecting the  
330 complexity of genetic disorders, where an affected individual carries multiple combinations of  
331 rare variants, and each combination contributes to either the same phenotype or distinct  
332 phenotypes at varying effect sizes (**Figure 5**). A limitation of our method is that it tends to be  
333 biased towards genes that are mutated frequently enough to be observed in a combination. This  
334 limitation can be addressed by fixing specific primary variants of interest irrespective of their  
335 frequency and screening for “second-hit” modifiers that significantly co-occur with the primary  
336 variant, such as the co-occurrence of *RBM8A* variants in distal 1q21.1 deletion carriers  
337 manifesting thrombocytopenia-absent-radius syndrome and *TBX6* variants in 16p11.2 deletion  
338 carriers with scoliosis<sup>46,47</sup>.

339 Our method is fast and scalable, allows for fine-tuning combinatorial searches based on  
340 frequency, statistical power, and multiple testing criteria, and can be adapted to enable  
341 computational approximations to further improve run time and assess higher-order combinations  
342 beyond triplets. While larger sample sizes are generally required for detecting smaller frequency  
343 differences, we note that our framework achieves reliable statistical power even with modest  
344 sample sizes, implying that our framework could be applied to exome sequencing studies of  
345 other neurodevelopmental disorders that have not been explored for combinatorial effects. This  
346 approach can also be used to address a variety of research questions involving rare event  
347 combinations, including searching for protective effects of rare variants where simultaneous  
348 mutations are enriched in controls but not in cases, and finding combinations that exhibit specific  
349 enrichment or depletion patterns in more than two phenotypic groups. In summary, we provide a  
350 conceptual framework and the necessary tools to identify the oligogenic basis for complex  
351 disorders such as autism and intellectual disability, which hitherto was restricted to the analysis  
352 of canonical disorders such as Hirschsprung disease<sup>48</sup> and Bardet-Biedl syndrome<sup>12</sup>.

353

354

## 355 MATERIALS AND METHODS

356 We developed *RareComb* to address computational and statistical challenges associated with  
357 combinatorial analysis of rare variants. *RareComb* first uses the apriori algorithm to efficiently  
358 count the frequencies of co-occurring variant combinations. It then uses one-tailed binomial tests  
359 to compare the observed frequency of each variant combination to the expected frequency  
360 derived under the assumption of independence among the constituent variants within each  
361 combination (**Figure 1**). This method can be applied to identify variant combinations that are  
362 significantly enriched in cases but not in controls. In studies involving multiple comorbid  
363 phenotypes, this method can also be used to detect associations between specific combinations of  
364 variants and one or more (comorbid) phenotypes (see **Supplementary Note**). The general  
365 principles of our method, built using the basic axioms of probability theory, can be easily  
366 extended to a variety of problems involving rare higher-order combinations (**Supp. Figure 11**).  
367

### 368 Identifying frequencies of rare variant combinations

369 *RareComb* utilizes the apriori algorithm to efficiently calculate frequencies of variant  
370 combinations from sparse Boolean matrices (of 0s and 1s) (**Supp. Figure 12A**). The apriori  
371 algorithm has been successfully applied to analyze consumer behavior, where identifying  
372 products frequently purchased together could benefit a company<sup>49,50</sup>. While an algorithm that is  
373 used to derive insights from patterns within highly frequent events (i.e. frequent itemset mining)  
374 might not seem like a good fit to analyze rare variant combinations, its ability to perform  
375 disciplined search based on both built-in and user-specified constraints makes it an ideal  
376 counting tool. *For example*, the apriori algorithm avoids enumerating each of the 50 million pairs  
377 or 167 billion triplets from just 10,000 variants, and instead prunes the search-space based on  
378 user-defined criteria such as minimum frequency threshold and size of combinations (pairs,  
379 triplets, etc.) (**Supp. Figure 12B**). *RareComb* applies an additional constraint to the algorithm to  
380 limit its search to co-occurring events, which further reduces the search space (see  
381 **Supplementary Note**). *For example*, when considering variants A and B, only the frequency of  
382 the presence of both variants (A=1 & B=1) is counted, and not absence of either or both variants  
383 (A=1 & B=0; A=0 & B=1; or A=0 & B=0).

384

385



## 386 **Statistical Inference**

387 *RareComb* utilizes the p-values of one-tailed binomial tests to establish the magnitude of  
388 enrichment for each rare variant combination (**Figure 1**). For each combination, *RareComb*  
389 formulates null and alternate hypotheses for the binomial test by considering the event of  
390 observing all constituent variants together within a group of individuals as success and all other  
391 possibilities as failure in a binomial trial:

$$392 \quad H_0 : \pi = \pi_0$$

$$393 \quad H_a : \pi > \pi_0$$

394 where,

395  $\pi$  = Probability of *observing* all constituent rare variants of a combination together within  
396 a cohort, i.e.,  $P(A=1 \ \& \ B=1)$

397  $\pi_0$  = *Expected* probability derived from the frequency of individual variants of a  
398 combination, under the assumption of independence, i.e.,  $P(A=1) * P(B=1)$ .

399 *RareComb* then compares the null binomial distribution derived using the sample size of the  
400 group (n) and the expected probability ( $\pi_0$ ) (i.e.,  $X \sim \text{Binom}(n, p = \pi_0)$ ) with the observed  
401 probability ( $\pi$ ), and calculates the probability of observing rare variants occurring together at  
402 least as frequently as they were observed within the cohort (i.e. p-value).

403 In case-control analyses, this method is applied independently to each group, and the p-  
404 values between them are compared. The combinations exhibiting enrichment in both cases and  
405 controls, likely due to proximity of variants in linkage disequilibrium, are eliminated, following  
406 which the p-values in cases are adjusted for multiple-testing to identify statistically significant  
407 combinations that exhibit enrichment in cases but not in controls. Finally, the effect sizes are  
408 calculated using Cohen's d and the statistical power is measured using 2-sample 2-proportion  
409 tests, as additional metrics to prioritize the final set of significant rare variant combinations. In  
410 genotype-comorbid phenotype association analyses, the method is applied just once to the entire  
411 cohort, with multiple-testing adjusted p-values serving as a sufficient metric to identify high  
412 quality associations between genotypes and two or more co-occurring phenotypes.

413

## 414 **Statistical power and computational performance of the method**

415 We measured the relationship between sample size and statistical power for both binomial and 2-  
416 sample 2-proportion tests used in the framework. It took 1,356 samples for the binomial test to

417 achieve a statistical power of 80% to establish statistical enrichment between expected and  
418 observed co-occurrence frequencies of 0.1% and 0.5% (**Supp. Figure 13**). This number  
419 increased to 6,469 when the test needed to be more sensitive to compare frequencies of 0.3% and  
420 0.5%. Similarly, it took 7,840 samples for the 2-sample 2-proportion test to achieve 80% power  
421 to establish statistical difference between co-occurrence frequencies of 2% and 0.5% observed in  
422 two groups (**Supp. Figure 14**). The sample size requirement increased to 14,633 to differentiate  
423 frequencies of 1.5% and 0.5% at 80% statistical power. These results align with the known  
424 relationship between sample size and statistical power, and indicate that our method can be  
425 reliably applied to analyze reasonably modest-size cohorts.

426 We also measured the run times for the case-control analysis to identify significant pairs  
427 and triplets of mutated genes using simulated data of three discrete sizes of samples (5,000,  
428 10,000, and 50,000 individuals) and genes (5,000, 10,000, and 15,000 genes). The apriori  
429 algorithm was run on single-core CPUs with 256 GB memory and was constrained to analyze  
430 combinations observed in at least 0.15% of the samples. Given the memory-intensive nature of  
431 the apriori algorithm implemented in the ‘arules’ package, 256 GB was chosen to maintain  
432 uniformity<sup>51</sup>. However, smaller input files could be processed successfully using much less  
433 memory. As expected, the runtimes were proportional to the size of the combination (pairs  
434 versus triplets) and the number of input variables (**Supp. Figure 15**). While the increase in run  
435 time with the increase in sample size is apparent for pairs, lower runtimes observed with running  
436 50,000 samples compared to 5,000 samples for triplets can be attributed to stochasticity of the  
437 input data. Overall, the analysis of gene pairs took between one minute and 12 minutes while  
438 triplets took between two minutes and 150 minutes. Since several factors influence the runtime  
439 of the method, a trial-and-error approach to determine an optimal minimum frequency threshold  
440 for co-occurring events can help identify relevant combinations without resulting in insufficient  
441 memory due to combinatorial explosion.

442

### 443 **Samples**

444 We used whole exome sequencing data from 6,189 affected males from the Simons Foundation  
445 Powering Autism Research (SPARK)<sup>21</sup> and 1,878 affected males from 2,247 simplex families  
446 from the Simons Simplex Collection (SSC)<sup>52</sup> cohort from the Simons Foundation Autism  
447 Research Initiative (SFARI)<sup>53</sup>. We selected only male probands for our analysis to avoid any



448 confounding effect due to gender or ascertainment bias<sup>54,55</sup>. While diagnosis information for  
449 intellectual disability (ID), anxiety, attention deficit hyperactivity disorders (ADHD),  
450 schizophrenia, language and sleep disorders were encoded as binary variables for the SPARK  
451 samples, full-scale intelligence quotient (IQ) scores were available for the SSC cohort.

452

### 453 **Data preparation and quality control**

454 Variant Call Format (VCF) files obtained from exome sequencing data were annotated using  
455 ANNOVAR<sup>56</sup> for rsID information and variant frequency using ExAC<sup>57</sup> and gnomAD<sup>58</sup>. To  
456 overcome the limitations of using a single method to predict pathogenicity, the effects of non-  
457 synonymous mutations were annotated using 11 prediction methods: SIFT<sup>59</sup>, Polyphen2<sup>60</sup>  
458 (HDIV), Polyphen2 (HVAR), LRT<sup>61</sup>, MutationTaster<sup>62</sup>, MutationAssessor<sup>63</sup>, FATHMM<sup>64</sup>,  
459 MetaSVM<sup>65</sup>, PROVEAN<sup>66</sup>, REVEL<sup>67</sup>, and CADD<sup>68</sup>. Briefly, all missense, stop-loss/gain, and  
460 start-loss/gain variants within exonic, 3', and 5' UTR regions with minor allele frequencies  $\leq 1\%$   
461 identified based on both ExAC and gnomAD databases were selected. Then, variants with allele  
462 depth of  $\geq 15$  and allele balance between 25% and 75% for heterozygous variants and  $> 90\%$  for  
463 homozygous variants were selected as high-quality variants. Deleteriousness of the variants were  
464 measured and reported differently by each prediction method. REVEL provided a score between  
465 0 and 1, with higher scores indicating higher level of deleteriousness, while Polyphen2 and  
466 MutationAssessor classified variants into one of three categories. *For example*, Polyphen2  
467 classified variants as 'Deleterious', 'Possibly damaging', or 'Tolerated', while MutationAssessor  
468 classified variants as 'High', 'Medium', or 'Low'. The other nine methods classified variants as  
469 either 'Deleterious' or 'Tolerated'. Pathogenicity reported by each tool was encoded as a binary  
470 variable, with the categories 'Possibly damaging' and 'Medium' encoded as 0.5. Thus, the  
471 composite pathogenicity score derived from the 10 tools could range between 0 and 10. Missense  
472 variants with a cumulative score of  $\geq 4$  and stop-loss/gain predicted as 'deleterious' either based  
473 on CADD score (CADD phred  $> 30$ ) or MutationTaster were considered deleterious for all  
474 analyses. Indels and other smaller structural variants were not considered, as their functional  
475 impact could not be easily assessed.

476

477

478

## 479 **Gene Ontology (GO) and Human Phenotype Ontology (HPO) enrichment analyses**

480 Gene Ontology term enrichment analyses were performed using the ‘Gene Ontology API’  
481 accessed using the ‘post’ command of the python package ‘requests’ (python version 3.7)<sup>31</sup>. All  
482 analyses were performed using parameters for *homo sapiens* (organism = ‘9606’) to identify  
483 biological processes enrichment (annotDataSet = ‘GO:0008150’) using binomial tests. HPO  
484 enrichment analyses were performed using data from the ‘genes\_to\_phenotype’ file obtained  
485 from the HPO website<sup>32</sup>. Since enrichment of phenotypes is not automatically evaluated by HPO,  
486 we used customized R scripts to derive baseline expectations that could be compared against the  
487 actual observations to determine significance using the p-values from binomial tests.

488

## 489 **Statistical analysis**

490 All statistical analyses were performed using R v3.6.1 (R Foundation for Statistical Computing,  
491 Vienna, Austria)<sup>69</sup> and Python (v3.7)<sup>70</sup>. All data-related plots were generated using the R  
492 package ggplot2<sup>71</sup>.

493

## 494 **Software Availability**

495 *RareComb* is available as an open-source (<https://github.com/girirajanlab/RareComb>) R package  
496 that can be downloaded from the Comprehensive R Archive Network (CRAN) repository<sup>72</sup>. It  
497 can also be installed into development environments via interfaces such as Rstudio<sup>73</sup> using the  
498 command `install.packages('RareComb')`. The tool provides several functionalities that allow  
499 users to run the types of analyses described in this manuscript. The functionalities are as follows:  
500 (1) Identify rare event combinations statistically enriched within a single group; (2) Identify rare  
501 event combinations statistically enriched in cases but not in controls; (3) Identify rare event  
502 combinations enriched in cases but depleted in controls; (4) Identify statistically enriched rare  
503 event combinations that include at least one element from an user-supplied list; and (5) Identify  
504 genotypes statistically enriched within individuals manifesting two or more comorbid  
505 phenotypes. Each functionality takes a Boolean matrix as input and provides a set of user-  
506 adjustable parameters to customize the analysis, and delivers the results in a tabular format as csv  
507 files. Detailed instructions on the available functionalities and parameters built into *RareComb*  
508 and their usage can be found on the GitHub page or CRAN website. A shiny app illustrating the  
509 ideas behind *RareComb* is available online at <https://girirajanlab.shinyapps.io/RareComb/><sup>74</sup>.

510 **DECLARATIONS**

511 **Ethics approval and consent to participate**

512 As these data were de-identified, all our samples were exempt from IRB review and conformed  
513 to the Helsinki Declaration. No other approvals were needed for the study.

514

515 **Consent for publication**

516 All authors agree and consent for publication of the manuscript.

517

518 **Competing interests**

519 The authors declare that no competing interests exist in relation to this work.

520

521 **Authors' contributions**

522 VK and SG conceived the project. VK performed the analyses, generated the plots/images, and  
523 wrote and revised the manuscript; SG supervised the research and wrote and revised the  
524 manuscript. All authors read and approved the final draft of the manuscript.

525

526 **Acknowledgements**

527 We thank Naomi Altman, Yifei Huang, Dajiang Liu, Matthew Jensen, and Corrine Smolen for  
528 constructive comments on the manuscript. This work was supported by R01-MH107431, R01-  
529 GM121907, Seed Grants program from the Institute of Computational and Data Sciences at Penn  
530 State, and resources from the Huck Institutes of the Life Sciences (to SG). The funding bodies  
531 had no role in data collection, analysis, and interpretation. The authors are grateful to all the  
532 families who participated in the SSC and SPARK consortia, as well as the principal  
533 investigators, clinical sites, and staff for the consortia. The authors appreciate obtaining access to  
534 genetic and phenotypic data for SPARK and SSC through the Simons Foundation Autism  
535 Research Initiative (SFARI) Base. Approved researchers can obtain the SSC and SPARK  
536 population datasets described in this study by applying at <https://base.sfari.org>.

537

538

539

## 540 References

- 541 1. Coventry, A. *et al.* Deep resequencing reveals excess rare recent variants consistent with  
542 explosive population growth. *Nat. Commun.* **1**, 131 (2010).
- 543 2. Keinan, A. & Clark, A. G. Recent explosive human population growth has resulted in an  
544 excess of rare genetic variants. *Science*. **336**, 740–743 (2012).
- 545 3. Tennessen, J. A. *et al.* Evolution and functional impact of rare coding variation from deep  
546 sequencing of human exomes. *Science*. **337**, 64–69 (2012).
- 547 4. McClellan, J. & King, M. C. Genetic heterogeneity in human disease. *Cell* **141**, 210–217  
548 (2010).
- 549 5. Wilfert, A. B. *et al.* Recent ultra-rare inherited variants implicate new autism candidate  
550 risk genes. *Nat. Genet.* **53**, 1125–1134 (2021).
- 551 6. Iossifov, I. *et al.* The contribution of de novo coding mutations to autism spectrum  
552 disorder. *Nature* **515**, 216–221 (2014).
- 553 7. Sebat, J. *et al.* Strong Association of De Novo Copy Number Mutations with Autism.  
554 *Science*. **316**, 445–449 (2007).
- 555 8. Badano, J. L. & Katsanis, N. Beyond mendel: An evolving view of human genetic disease  
556 transmission. *Nat. Rev. Genet.* **3**, 779–789 (2002).
- 557 9. Gifford, C. *et al.* Oligogenic inheritance of a human heart disease involving a genetic  
558 modifier. *Science*. **364**, 865–870 (2019).
- 559 10. Pizzo, L. *et al.* Rare variants in the genetic background modulate cognitive and  
560 developmental phenotypes in individuals carrying disease-associated variants. *Genet.*  
561 *Med.* **21**, 816–825 (2019).
- 562 11. Girirajan, S. *et al.* A recurrent 16p12.1 microdeletion supports a two-hit model for severe  
563 developmental delay. *Nat. Genet.* **42**, 203–209 (2010).
- 564 12. Badano, J. L. *et al.* Dissection of epistasis in oligogenic Bardet-Biedl syndrome. *Nature*  
565 **439**, 326–330 (2006).
- 566 13. Leblond, C. S. *et al.* Genetic and functional analyses of SHANK2 mutations suggest a  
567 multiple hit model of autism spectrum disorders. *PLoS Genet.* **8**, e1002521 (2012).
- 568 14. Uricchio, L. H., Zaitlen, N. A., Ye, C. J., Witte, J. S. & Hernandez, R. D. Selection and  
569 explosive growth alter genetic architecture and hamper the detection of causal rare  
570 variants. *Genome Res.* **26**, 863–873 (2016).

- 571 15. Halvorsen, M. *et al.* Increased burden of ultra-rare structural variants localizing to  
572 boundaries of topologically associated domains in schizophrenia. *Nat. Commun.* **11**, 1842  
573 (2020).
- 574 16. Krumm, N. *et al.* Excess of rare, inherited truncating mutations in autism. *Nat. Genet.* **47**,  
575 582–588 (2015).
- 576 17. Murdoch, W. J., Singh, C., Kumbier, K., Abbasi-Asl, R. & Yu, B. Definitions, methods,  
577 and applications in interpretable machine learning. *Proc. Natl. Acad. Sci. U. S. A.* **116**,  
578 22071–22080 (2019).
- 579 18. Molnar, C., Casalicchio, G. & Bischl, B. Interpretable Machine Learning – A Brief  
580 History, State-of-the-Art and Challenges. *Commun. Comput. Inf. Sci.* **1323**, 417–431  
581 (2020).
- 582 19. Agarwal, R. & Srikant, R. Fast algorithms for mining association rules. *Proc. 20th VLDB*  
583 *Conf.*
- 584 20. Agrawal, Rakesh; Ramakrishnan, S. Fast Algorithms for Mining Association Rules. *Proc.*  
585 *20th VLDB Conf.* 487–499 (1994).
- 586 21. The SPARK Consortium. SPARK: A US Cohort of 50,000 Families to Accelerate Autism  
587 Research. *Neuron* **97**, 488–493 (2018).
- 588 22. Fischbach, G. D. & Lord, C. The simons simplex collection: A resource for identification  
589 of autism genetic risk factors. *Neuron* **68**, 192–195 (2010).
- 590 23. Girirajan, S. *et al.* Phenotypic Heterogeneity of Genomic Disorders and Rare Copy-  
591 Number Variants. *N. Engl. J. Med.* **367**, 1321–1331 (2012).
- 592 24. Turner, T. N. *et al.* Genomic Patterns of De Novo Mutation in Simplex Autism. *Cell* **171**,  
593 710–722 (2017).
- 594 25. Dodé, C. *et al.* Loss-of-function mutations in FGFR1 cause autosomal dominant Kallmann  
595 syndrome. *Nat. Genet.* **33**, 463–465 (2003).
- 596 26. Schell, U. *et al.* Mutations in FGFR1 and FGFR2 cause familial and sporadic pfeiffer  
597 syndrome. *Hum. Mol. Genet.* **4**, 323–328 (1995).
- 598 27. Dazzo, E. *et al.* Heterozygous Reelin Mutations Cause Autosomal-Dominant Lateral  
599 Temporal Epilepsy. *Am. J. Hum. Genet.* **96**, 992–1000 (2015).
- 600 28. Yoshinaga, T. *et al.* A novel frameshift mutation of SYNE1 in a Japanese family with  
601 autosomal recessive cerebellar ataxia type 8. *Hum. Genome Var.* **4**, 17052 (2017).

- 602 29. Synofzik, M. *et al.* SYNE1 ataxia is a common recessive ataxia with major non-cerebellar  
603 features: A large multi-centre study. *Brain* **139**, 1378–1393 (2016).
- 604 30. Prasad, A. *et al.* A Discovery resource of rare copy number variations in individuals with  
605 autism spectrum disorder. *G3 Genes, Genomes, Genet.* **2**, 1665–1685 (2012).
- 606 31. Mi, H., Muruganujan, A., Ebert, D., Huang, X. & Thomas, D. PANTHER version 14 :  
607 more genomes , a new PANTHER GO-slim and improvements in enrichment analysis  
608 tools. *Nucleic Acids Res.* **47**, 419–426 (2019).
- 609 32. Köhler, S. *et al.* The human phenotype ontology in 2021. *Nucleic Acids Res.* **49**, D1207–  
610 D1217 (2021).
- 611 33. Gilger, J. W., Pennington, B. F. & DeFries, J. C. A Twin Study of the Etiology of  
612 Comorbidity: Attention-deficit Hyperactivity Disorder and Dyslexia. *J. Am. Acad. Child*  
613 *Adolesc. Psychiatry* **31**, 343–348 (1992).
- 614 34. Krumm, N. *et al.* Transmission disequilibrium of small CNVs in simplex autism. *Am. J.*  
615 *Hum. Genet.* **93**, 595–606 (2013).
- 616 35. Guo, H. *et al.* Genome-wide copy number variation analysis in a Chinese autism spectrum  
617 disorder cohort. *Sci. Rep.* **7**, 44155 (2017).
- 618 36. Guemez-Gamboa, A. *et al.* Inactivating mutations in MFSD2A, required for omega-3 fatty  
619 acid transport in brain, cause a lethal microcephaly syndrome. *Nat. Genet.* **47**, 809–813  
620 (2015).
- 621 37. Sanders, A. R. *et al.* Haplotypic association spanning the 22q11.21 genes COMT and  
622 ARVCF with schizophrenia. *Mol. Psychiatry* **10**, 353–365 (2005).
- 623 38. Kenny, E. M. *et al.* Excess of rare novel loss-of-function variants in synaptic genes in  
624 schizophrenia and autism spectrum disorders. *Mol. Psychiatry* **19**, 872–879 (2014).
- 625 39. Zheng, G. X. Y. *et al.* Haplotyping germline and cancer genomes with high-throughput  
626 linked-read sequencing. *Nat. Biotechnol.* **34**, 303–311 (2016).
- 627 40. Girirajan, S. *et al.* Relative burden of large CNVs on a range of neurodevelopmental  
628 phenotypes. *PLoS Genet.* **7**, e1002334 (2011).
- 629 41. Papadimitriou, S. *et al.* Predicting disease-causing variant combinations. *Proc. Natl. Acad.*  
630 *Sci. U. S. A.* **116**, 11878–11887 (2019).
- 631 42. Kerner, G. *et al.* A genome-wide case-only test for the detection of digenic inheritance in  
632 human exomes. *Proc. Natl. Acad. Sci. U. S. A.* **117**, 19367–19375 (2020).

- 633 43. Schaaf, C. P. *et al.* Oligogenic heterozygosity in individuals with high-functioning autism  
634 spectrum disorders. *Hum. Mol. Genet.* **20**, 3366–3375 (2011).
- 635 44. Singh, T. *et al.* The contribution of rare variants to risk of schizophrenia in individuals  
636 with and without intellectual disability. *Nat. Genet.* **49**, 1167–1173 (2017).
- 637 45. Fitzgerald, T. W. *et al.* Large-scale discovery of novel genetic causes of developmental  
638 disorders. *Nature* **519**, 223–228 (2015).
- 639 46. Albers, C. A. *et al.* Compound inheritance of a low-frequency regulatory SNP and a rare  
640 null mutation in exon-junction complex subunit RBM8A causes TAR syndrome. *Nat.*  
641 *Genet.* **44**, 435–439 (2012).
- 642 47. Yang, N. *et al.* TBX6 compound inheritance leads to congenital vertebral malformations  
643 in humans and mice. *Hum. Mol. Genet.* **28**, 539–547 (2019).
- 644 48. Gabriel, S. B. *et al.* Segregation at three loci explains familial and population risk in  
645 Hirschsprung disease. *Nat. Genet.* **31**, 89–93 (2002).
- 646 49. Brijs, T., Swinnen, G., Vanhoof, K. & Wets, G. Using association rules for product  
647 assortment decisions. *ACM* (1999). doi:10.1145/312129.312241
- 648 50. Glance, N. *et al.* Deriving marketing intelligence from online discussion. *Proc. ACM*  
649 *SIGKDD Int. Conf. Knowl. Discov. Data Min.* 419–428 (2005).  
650 doi:10.1145/1081870.1081919
- 651 51. Hahsler, M., Grun, B. & Hornik, K. arules – A Computational Environment for Mining  
652 Association Rules and Frequent Item Sets. *J. Stat. Softw.* **14**, 1–6 (2005).
- 653 52. Sanders, S. J. *et al.* Insights into Autism Spectrum Disorder Genomic Architecture and  
654 Biology from 71 Risk Loci. *Neuron* **87**, 1215–1233 (2015).
- 655 53. SFARI Base. <https://www.sfari.org/resource/sfari-base/>
- 656 54. Jacquemont, S. *et al.* A higher mutational burden in females supports a ‘female protective  
657 model’ in neurodevelopmental disorders. *Am. J. Hum. Genet.* **94**, 415–425 (2014).
- 658 55. Polyak, A., Rosenfeld, J. A. & Girirajan, S. An assessment of sex bias in  
659 neurodevelopmental disorders. *Genome Med.* **7**, 94 (2015).
- 660 56. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: Functional annotation of genetic  
661 variants from high-throughput sequencing data. *Nucleic Acids Res.* **38**, e164 (2010).
- 662 57. Exome aggregation Consortium. Analysis of protein-coding genetic variation in 60,706  
663 humans. *Nature* **536**, 285–291 (2016).



- 664 58. Genome Aggregation Database Consortium. The mutational constraint spectrum  
665 quantified from variation in 141,456 humans. *Nature* **581**, 434–443 (2020).
- 666 59. Ng, P. C. & Henikoff, S. SIFT: Predicting amino acid changes that affect protein function.  
667 *Nucleic Acids Res.* **31**, 3812–4 (2003).
- 668 60. Adzhubei, I. A. *et al.* A method and server for predicting damaging missense mutations.  
669 *Nat. Methods* **7**, 248–249 (2010).
- 670 61. Chun, S. & Fay, J. C. Identification of deleterious mutations within three human genomes.  
671 *Genome Res.* **19**, 1553–1561 (2009).
- 672 62. Schwarz, J. M., Rödelberger, C., Schuelke, M. & Seelow, D. MutationTaster evaluates  
673 disease-causing potential of sequence alterations. *Nat. Methods* **7**, 575–576 (2010).
- 674 63. Reva, B., Antipin, Y. & Sander, C. Predicting the functional impact of protein mutations:  
675 Application to cancer genomics. *Nucleic Acids Res.* **39**, 37–43 (2011).
- 676 64. Shihab, H. A. *et al.* Predicting the Functional, Molecular, and Phenotypic Consequences  
677 of Amino Acid Substitutions using Hidden Markov Models. *Hum. Mutat.* **34**, 57–65  
678 (2013).
- 679 65. Kim, S., Jhong, J. H., Lee, J. & Koo, J. Y. Meta-analytic support vector machine for  
680 integrating multiple omics data. *BioData Min.* **10**, 2 (2017).
- 681 66. Choi, Y. & Chan, A. P. PROVEAN web server: A tool to predict the functional effect of  
682 amino acid substitutions and indels. *Bioinformatics* **31**, 2745–2747 (2015).
- 683 67. Ioannidis, N. M. *et al.* REVEL: An Ensemble Method for Predicting the Pathogenicity of  
684 Rare Missense Variants. *Am. J. Hum. Genet.* **99**, 877–885 (2016).
- 685 68. Rentzsch, P., Witten, D., Cooper, G. M., Shendure, J. & Kircher, M. CADD: Predicting  
686 the deleteriousness of variants throughout the human genome. *Nucleic Acids Res.* **47**,  
687 D886–D894 (2019).
- 688 69. R Core Team. R: A language and environment for statistical computing. (2019).
- 689 70. Van Rossum, Guido; Drake, F. L. Python 3 Reference Manual. (2009).
- 690 71. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*. (Springer-Verlag New York,  
691 2016).
- 692 72. CRAN Repository. Available at: [https://cran.r-](https://cran.r-project.org/web/packages/available_packages_by_name.html)  
693 [project.org/web/packages/available\\_packages\\_by\\_name.html](https://cran.r-project.org/web/packages/available_packages_by_name.html).
- 694 73. RStudio Team. RStudio: Integrated Development for R. (2020).



695 74. Chang, W., Cheng, J., Allaire, J., Xie, Y. & McPherson, J. shiny: Web Application  
696 Framework for R. (2020).

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

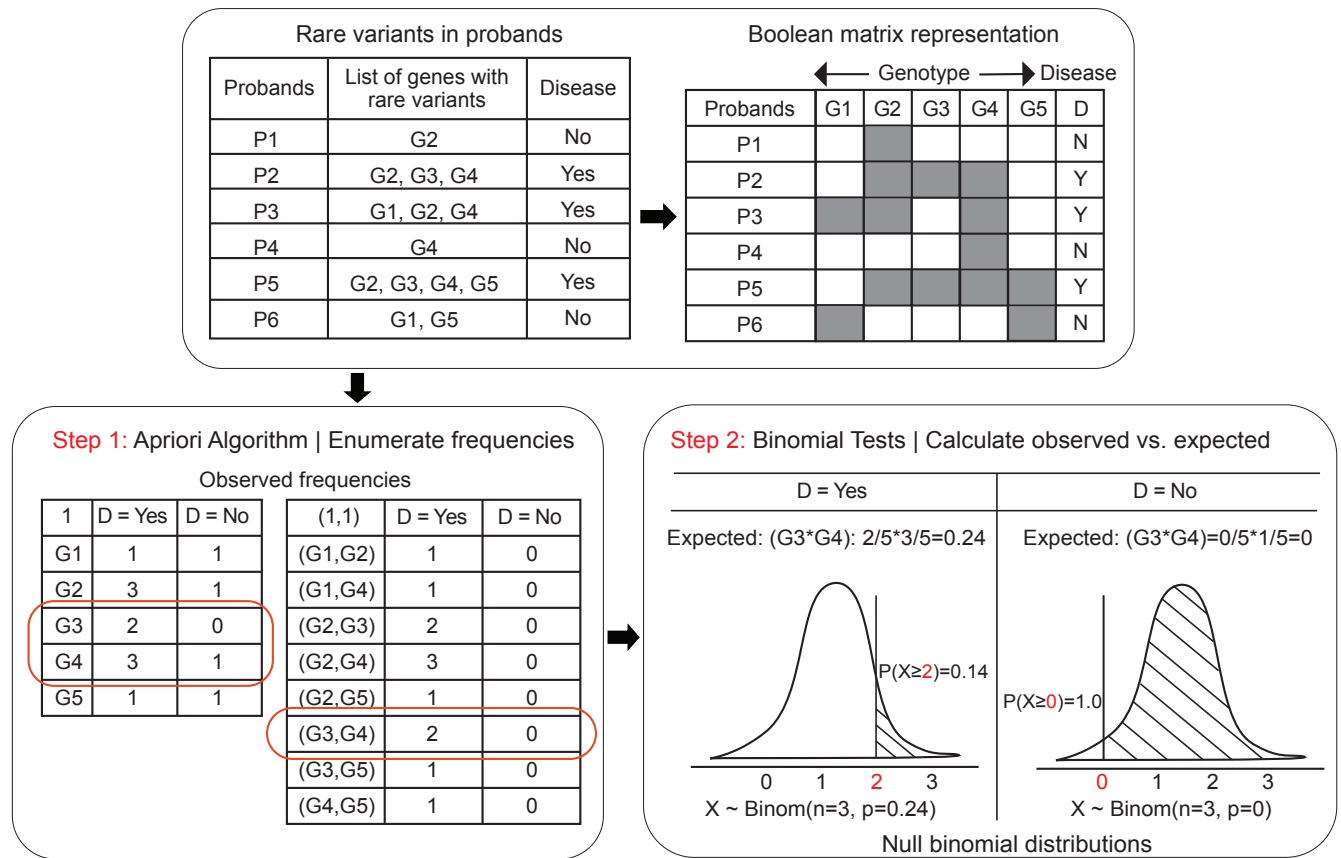
722

723

724

725

726 MAIN FIGURES



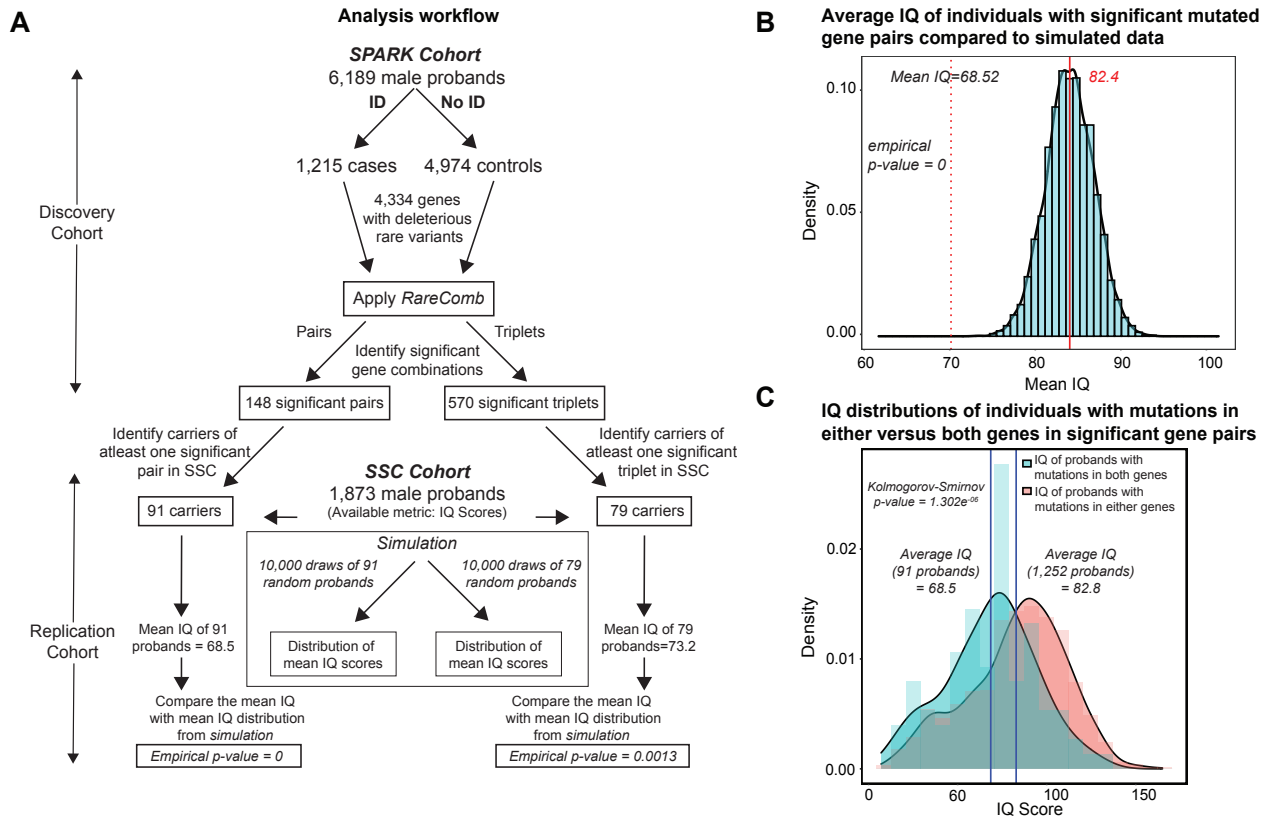
727

728 **Figure 1: Conceptual overview of combinatorial analyses using *RareComb*.** A Boolean  
 729 representation of genotype (mutated genes, G1, G2, etc) and disease status for probands (P1, P2,  
 730 etc) is shown. In *step 1*, the apriori algorithm is applied to the Boolean input matrix to calculate  
 731 the frequencies of individual (for example, G1) and simultaneous occurrences of events (G1 and  
 732 G2) that meet the user-specified criteria, including the size of combinations (pairs, triplets, etc.)  
 733 and minimum frequency threshold of simultaneous occurrences. In *step 2*, independently in case  
 734 and control groups, for each combination, the binomial test is applied to compare the observed  
 735 frequency of simultaneous occurrence of events with its corresponding null binomial distribution  
 736 of the expected frequencies calculated under the assumption of independence. Binomial test for  
 737 gene pair G3 and G4 is shown as an example.

738

739

740



741

742 **Figure 2: Combinations of rare variants contributing to intellectual disability (ID)**

743 **phenotype. (A)** An outline of the approach used to identify and validate mutated gene pairs and

744 triplets enriched in probands with ID is shown. We tested whether mutated gene pairs identified

745 as significant in one cohort (SPARK) are also associated with severe phenotypes in an

746 independent cohort (SSC). To test this, we obtained the mean IQ score of individuals from the

747 SSC cohort carrying significant combinations identified from the SPARK cohort. Empirical p-

748 values were then calculated based on the deviation of the mean IQ from the distribution of mean

749 IQ scores obtained from 10,000 random draws in the simulation. **(B)** The mean IQ of individuals

750 with mutated gene pairs in the SSC cohort was significantly lower (empirical p-value=0) when

751 compared to the distribution of mean IQ scores obtained from the simulation. **(C)** Histogram

752 shows the distributions of IQ scores of SSC probands who carried mutations in either genes

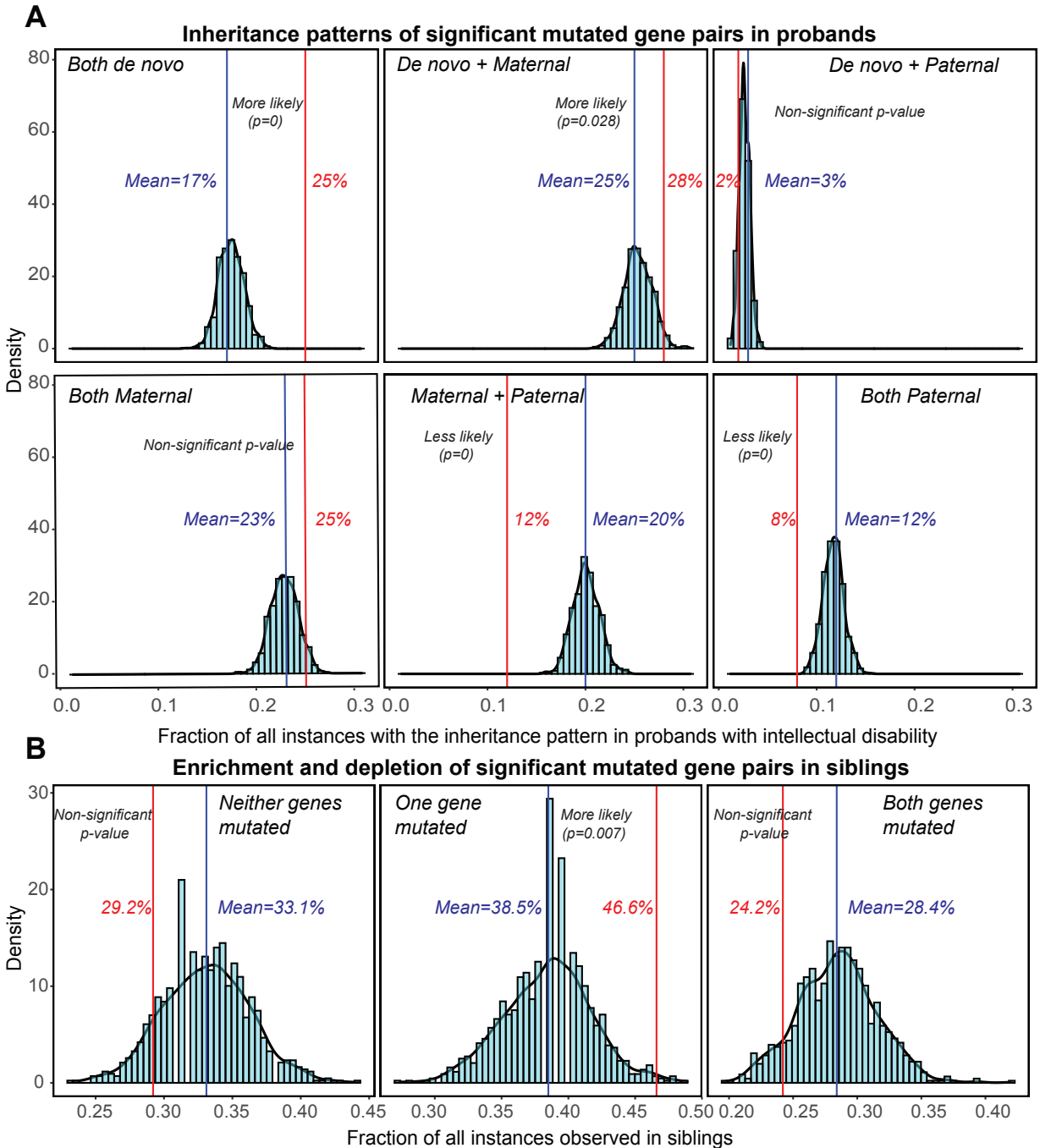
753 versus both constituent genes of the significant gene pairs. The distributions were significantly

754 different from each other (p-value =  $1.302 \times 10^{-6}$ , Kolmogorov-Smirnov test).

755

756

757



758

759 **Figure 3: Analysis of parental and sibling inheritance patterns of significant gene pairs**

760 **associated with ID. (A)** Fraction of all instances of significant gene pairs observed within each

761 of the six possible parental inheritance patterns (red) compared against 1,000 simulations is

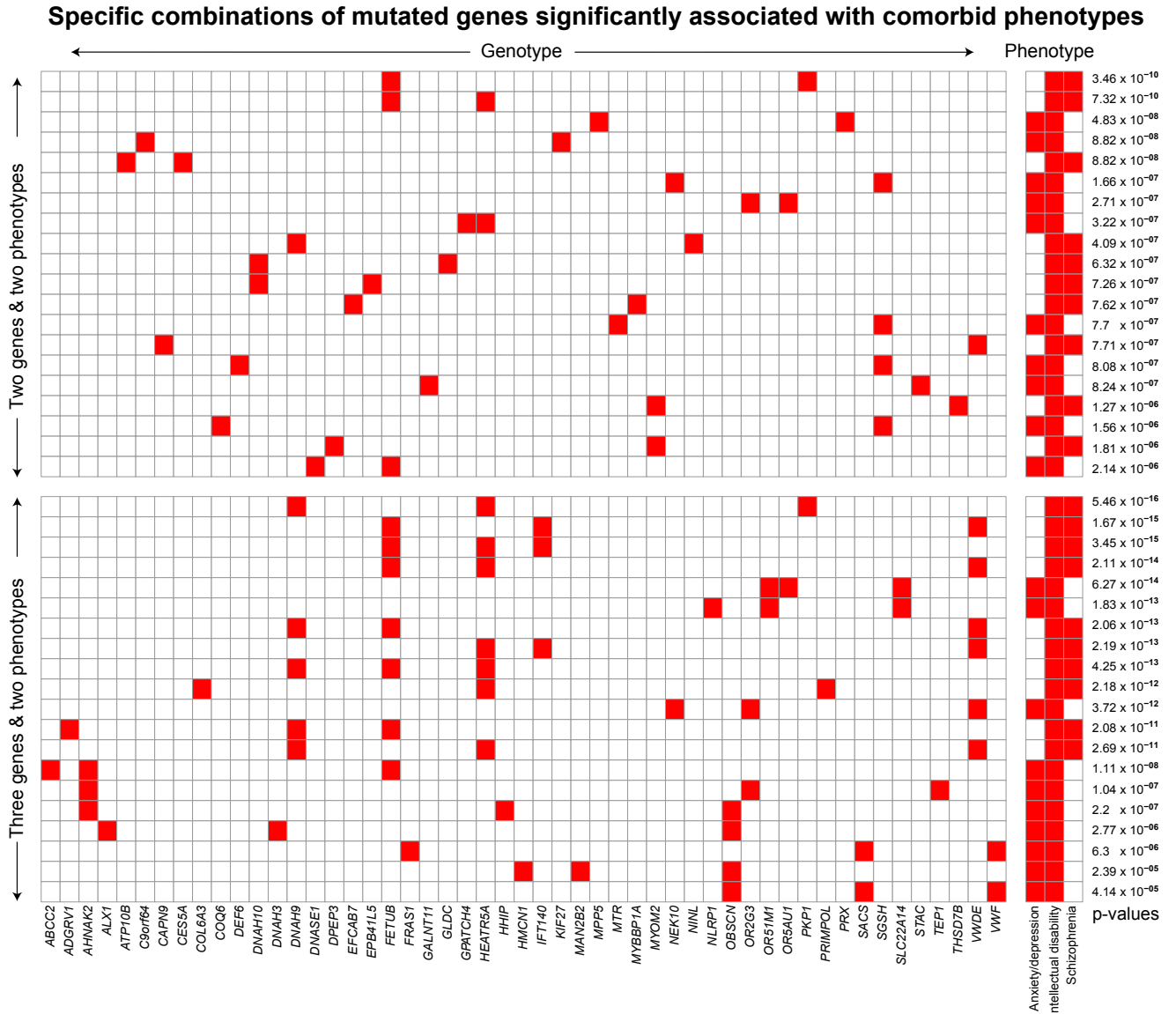
762 shown (blue). During each simulation, random mutated gene pairs from the SSC cohort were

763 selected, the inheritance status of the mutations was identified, and the fraction of those instances

764 belonging to one of the six pre-defined categories was calculated. Comparing the observed  
765 fractions with the simulated fractions indicate statistical enrichment for two specific inheritance  
766 patterns based on empirical p-values: both variants being *de novo*, and one variant being *de novo*  
767 and the other transmitted from the mother. **(B)** Histograms show the carrier status of significant  
768 gene pairs in siblings of carrier probands (red) compared against 1,000 simulations (blue).  
769 Among significant pairs, both genes were mutated in only 24.2% of all siblings (compared to  
770 28.4% in simulations), whereas one of the two genes was mutated in 46.6% of all siblings  
771 (compared to 38.5% in simulations). These results show that mutations are more likely to be  
772 observed in just one of the two genes within the gene pairs and are less likely to be observed  
773 simultaneously in siblings of carrier probands.

774

775



776

777 **Figure 4: Analysis of comorbid phenotypes using *RareComb*.** We analyzed the genotypes of

778 probands with anxiety/depression, ID, or schizophrenia. The heatmap shows combinations of

779 two or three mutated genes that were significantly enriched in individuals with specific patterns

780 of comorbid phenotypes compared to the expected frequency under the assumption of

781 independence.

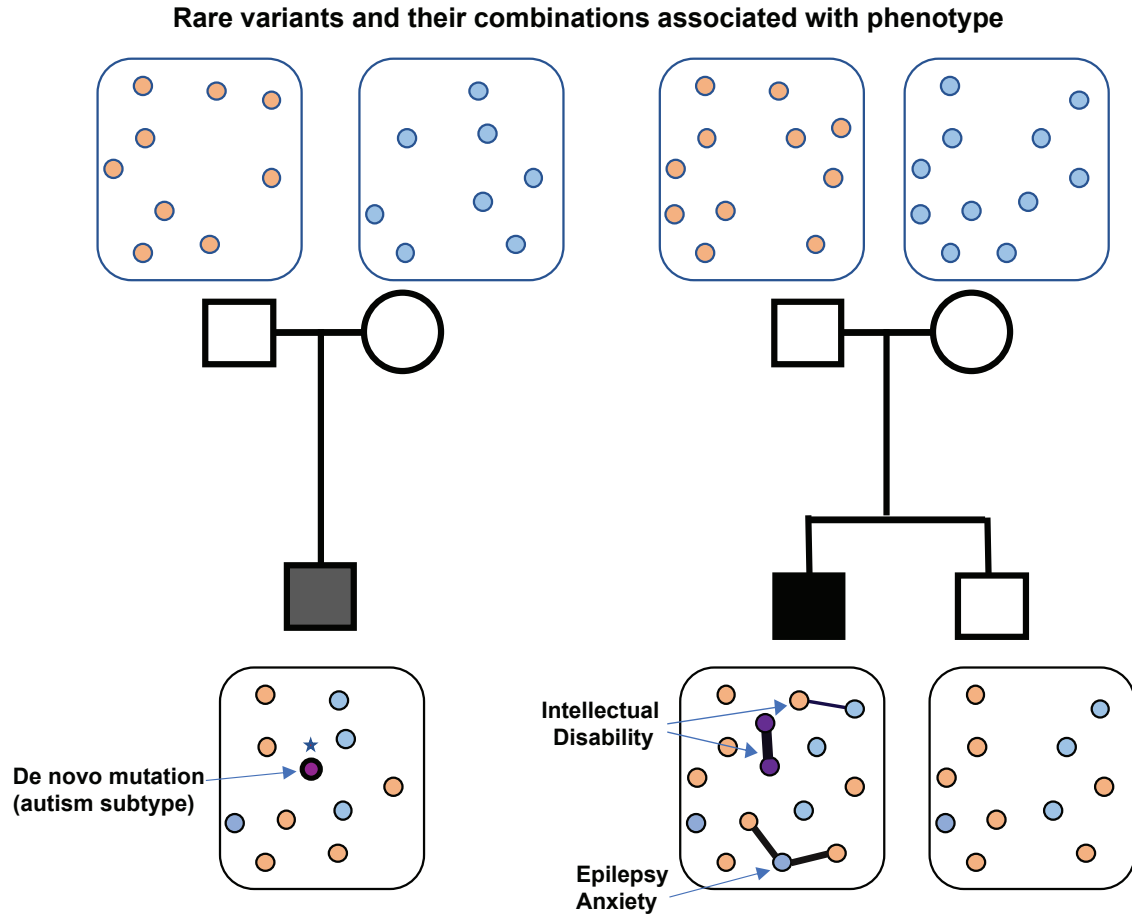
782

783

784

785

786



787

788 **Figure 5: Rare variant models for complex disorders.** The schematic shows two models for  
789 the genetic etiology of complex disorders. Circles represent rare variants present that are either  
790 *de novo* or inherited from a parent. On the left, individual high-effect *de novo* variants are  
791 strongly associated with a phenotype of interest. On the right, rare variants within an individual  
792 combine in multiple ways and contribute towards distinct phenotypes. The thickness of the  
793 connecting lines denotes effect sizes, and an affected individual can carry multiple oligogenic  
794 combinations of rare variants, each of which contributes to the same or distinct phenotypes. This  
795 extension of the oligogenic model enables further dissection of the genetic architecture of  
796 complex disorders.

797