



A general framework for identifying oligogenic combinations of rare variants in complex disorders

Vijay Kumar Pounraja and Santhosh Girirajan

Genome Res. published online March 17, 2022

Access the most recent version at doi:[10.1101/gr.276348.121](https://doi.org/10.1101/gr.276348.121)

P<P	Published online March 17, 2022 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Creative Commons License	This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see https://genome.cshlp.org/site/misc/terms.xhtml). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **A general framework for identifying oligogenic**
2 **combinations of rare variants in complex disorders**

3
4 Vijay Kumar Pounraja^{1,2} and Santhosh Girirajan^{1,2,3}

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
14
15
16 **Correspondence:**

17 Santhosh Girirajan

18 205A Life Sciences Building

19 Pennsylvania State University

20 University Park, PA 16802

21 E-mail: sxg47@psu.edu

22 **ABSTRACT**

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

43 INTRODUCTION

44 Recent human population growth has led to a rapid increase in the load of rare variants affecting
45 functionally important regions of the genome (Coventry et al. 2010; Keinan and Clark 2012;
46 Tennessen et al. 2012). Thus, rare variants are collectively more abundant in the population
47 compared to common variants, many of which confer significant risk for neurodevelopmental
48 disorders such as autism and intellectual disability (McClellan and King 2010; The 1000
49 Genomes Project Consortium 2015; Taliun et al. 2021; Backman et al. 2021). In fact, recent
50 studies have directly implicated rare damaging mutations that are very recent or *de novo* in more
51 than one hundred genes towards neurodevelopmental disorders (Wilfert et al. 2021; Iossifov et
52 al. 2014; Sebat et al. 2007). The ability to establish robust associations between rare variants of
53 high effect size and complex disease has made this class of variants the primary focus of recent
54 studies. However, a much larger class of rare and variably expressive variants that are
55 individually less deleterious but, in combination, exert large effects towards disease is often
56 overlooked. Variants in this category are often transmitted across generations without adverse
57 effects on their carriers until they encounter other similar variants that, when combined, lead to
58 genetic interactions conferring a higher risk for disease than their individual risks (Badano and
59 Katsanis 2002; Gifford et al. 2019). While this phenomenon underpins oligogenic models
60 proposed over the years, studies so far have not focused on detecting combinatorial effects of
61 specific sets of rare variants towards disease phenotypes (Pizzo et al. 2019; Girirajan et al. 2010;
62 Badano et al. 2006; Leblond et al. 2012).

63 Identifying the effects of specific combinations of rare variants towards disease etiology
64 has been challenging for many reasons. First, combinations of rare variants are rarer, and
65 extremely large cohorts are required to observe even a few recurrent instances of specific variant
66 combinations (Uricchio et al. 2016). Prior studies of oligogenic models for rare variants evaded
67 this problem by aggregating variant information at the sample level and comparing the overall
68 burden of rare variants between groups of individuals (such as cases and controls) (Sebat et al.
69 2007; Halvorsen et al. 2020; Krumm et al. 2015; Iossifov et al. 2014). Second, the combinatorial
70 explosion resulting from even a small set of rare variants makes it difficult to exhaustively
71 evaluate all combinations. While sophisticated frameworks such as network analysis and
72 machine learning provide powerful tools to model the composite effects of thousands of
73 variables on a complex system and predict emergent behaviors and quantitative outcomes,

74 adapting them to exhaustively search and delineate the effects of specific combinations of
75 variables is daunting (Murdoch et al. 2019; Molnar et al. 2020). Furthermore, incorporating an
76 efficient search tool into these frameworks and extending them to detect higher-order
77 combinatorial effects would be nearly impossible. Third, even when all combinations of rare
78 variants could be exhaustively evaluated within a large cohort, there is a lack of methods that are
79 sensitive enough to detect small differences between comparison groups to establish statistical
80 significance. Therefore, an alternate approach that is highly flexible, scalable, and sensitive is
81 necessary to address computational and statistical challenges associated with assessing rare
82 variant combinations.

83 Here, we present a combinatorial framework called RareComb that couples the *a priori*
84 algorithm with binomial tests to overcome the limitations of data sparsity and high
85 dimensionality, and systematically analyzes patterns of rare variants between groups of interest
86 to identify specific combinations that are significantly associated with phenotypes (Agrawal,
87 Rakesh; Ramakrishnan 1994). We demonstrate the utility and adaptability of our framework by
88 identifying mutated gene combinations significantly associated with one or more phenotypes
89 among children with autism. Our generalizable and modular framework does not depend on *a*
90 *priori* knowledge and can detect rare patterns from high-dimensional genetic data to generate
91 interpretable results, making it readily applicable for analyzing cohorts of all size ranges to
92 dissect the genetic basis of complex disorders.

93 **RESULTS**

94 We hypothesized that two or more genes disrupted simultaneously by rare deleterious mutations
95 contribute to a highly penetrant phenotype, as in an oligogenic model, or lead to a more severe
96 phenotype than when each of the same genes are disrupted individually. We developed
97 RareComb as a framework that combines data mining and statistical analysis to identify specific
98 combinations (such as pairs, triplets, etc.) of rare variants that show significant associations with
99 one or more phenotypes. RareComb analyzes an ' $n \times p$ ' sparse Boolean matrix with ' p ' genes in
100 ' n ' individuals in two discrete steps (**Figure 1**). First, it applies the apriori algorithm
101 independently in cases and controls to enumerate the frequency of all simultaneously mutated
102 combinations that meet a pre-set minimum frequency threshold (**Supplemental Fig. S1**).
103 Second, for each qualifying combination of variants, the method derives the expected frequency
104 of simultaneously observing mutations in the constituent genes under the assumption of
105 independence. It then independently quantifies the magnitude of deviation of the observed from
106 the expected frequencies using binomial tests in cases and controls and uses multiple-testing
107 adjusted p-values to identify combinations that are statistically enriched in cases but not in
108 controls. Finally, the method calculates effect sizes using Cohen's d and statistical power at 1%
109 and 5% significance thresholds, to enable prioritization of a high confidence set of combinations
110 that contribute to the phenotype in an oligogenic manner.

111

112 **RareComb identifies oligogenic combinations associated with ID and autism**

113 We sought to identify pairs and triplets of mutated genes that are significantly associated with
114 intellectual disability/cognitive impairment (ID) phenotypes by analyzing 6,189 affected males
115 from the Simons Foundation Powering Autism Research or SPARK (The SPARK Consortium
116 2018) cohort for discovery and 1,878 affected males from the Simons Simplex Collection or SSC
117 (Fischbach and Lord 2010) cohort for validation. To facilitate cross-cohort comparison, we
118 identified 10,217 rare variants ($MAF \leq 1\%$) that were predicted to be deleterious by multiple
119 methods and observed in both cohorts and aggregated these variants to genes for the analysis
120 (see **Methods**). For this study, we first categorized 1,215 probands from the SPARK cohort
121 diagnosed with ID/cognitive impairment as "cases" and 4,974 probands without ID as "controls"
122 (**Figure 2A**). We then applied RareComb to cases after constraining it to only evaluate those
123 gene combinations in which simultaneous mutations are observed in at least five probands (i.e.,

124 minimum frequency threshold). We identified 25,602 pairs involving 1,956 mutated genes in
125 cases that were observed at a higher frequency than expected under the assumption of
126 independence.

127 Similarly, analyzing the controls using only the 1,956 genes mutated in cases, RareComb
128 identified 148 pairs of mutated genes that were significantly enriched in cases but not in controls
129 (**Supplemental Table S1**), with moderate to high effect sizes (Cohen's d , 0.08-0.15) and
130 adequate statistical power (70%-100% at 5% significance threshold) (**Supplemental Fig. S2**).
131 These 148 gene pairs belonged to 142 probands, with 74% (105/142) of them carrying more than
132 one significant pair. These observations suggest that an individual can carry multiple
133 combinations, each contributing to the same phenotype at varying effect sizes (**Supplemental**
134 **Fig. S3**). To identify enrichment for specific variant types within combinations, we examined the
135 148 significant gene pairs by mutation type, including missense, stop-loss, and stop-gain
136 mutations. We identified 871 instances of variant pairs, of which 95.64% (833/871) contained a
137 missense mutation in both genes, 4.36% (38/871) contained a missense in one of the genes and a
138 stop-gain in the other. We found no instances of pairs of genes with missense/stop-loss, stop-
139 loss/stop-loss, stop-gain/stop-loss, or stop-gain/stop-gain mutations. To evaluate the statistical
140 significance of these observed proportions, we generated all possible variant pairs from all male
141 probands and calculated the expected proportion for each possible pair of variant types
142 (**Supplemental Table S2**). Out of all possible variant pairs, 93.9% were missense/missense
143 variant pairs and 5.89% were missense/stop-gain pairs. One-tailed binomial tests showed that the
144 95.64% observed in our data for missense/missense pairs was higher than the expected 93.9% (p -
145 value = 0.015) but the 4.36% observed for missense/stop-gain pairs was not significantly lower
146 than the expected 5.89% (p -value = 0.58). These results suggest that there is a higher propensity
147 for missense mutations to form combinations than the relatively higher impact stop-gain or stop-
148 loss mutations.

149 We next sought to validate the association of these 148 mutated gene pairs towards
150 intellectual disability. We hypothesized that if the association of the gene pairs with ID in the
151 SPARK cohort were truly significant, carriers of mutations in those gene pairs would tend to
152 have lower than average IQ scores in the independent SSC cohort. We found that 90 of the 148
153 significant pairs identified in the SPARK cohort were observed in at least one proband in the
154 SSC cohort. These 90 mutated gene pairs were carried by 91 unique probands, whose average

155 full-scale IQ scores (average IQ=68.52) were lower than those of all ascertained probands in the
156 SSC cohort (average IQ=86). To assess the significance of this result, we performed 10,000
157 random draws of 91 probands from the SSC cohort to generate a simulated distribution of their
158 average IQ scores. The average IQ of carriers of mutated gene pairs (average IQ=68.52) was
159 significantly lower than the overall distribution of average IQ derived from simulations (average
160 IQ ranged from 73 to 92; empirical $p=0$) (**Figure 2B**). Furthermore, the average IQ of the 91
161 SSC probands with both mutated genes was significantly lower than the average IQ of 1,252
162 carriers of mutations in only one of the two genes (68.5 versus 82.8; Kolmogorov-Smirnov $p =$
163 1.302×10^{-16}) (**Figure 2C**). When each of the 90 combinations was evaluated individually,
164 carriers of mutations in both genes for 73% (66/90) of the combinations showed lower IQ than
165 individuals with mutations in individual genes of the same combination, with 39/90 remaining
166 significant after multiple testing correction (**Supplemental Table S3; Supplemental Fig. S4**).
167 These results provide evidence for synergistic effects of deleterious mutations within specific
168 pairs of genes towards ID phenotypes. We note that this analysis only considered combinations
169 that were enriched in cases but not in controls for multiple testing. Therefore, we repeated the
170 analysis using a more conservative approach that considered all combinations that met the
171 frequency threshold in cases for multiple-testing correction (see **Methods**), and obtained 115
172 significant pairs belonging to 79 probands (**Supplemental Table S4**). The average IQ of carriers
173 of mutated gene pairs (average IQ=69.11) remained significantly lower than the overall
174 distribution of average IQ derived from simulations (average IQ ranged from 73 to 92; empirical
175 $p=0$) (**Supplemental Fig. S5**). We also conducted the analysis on the entire cohort by combining
176 6,189 male and 1,528 female probands together and identified 199 gene pairs belonging to 82
177 males and 14 female probands (**Supplemental Table S5**). Our results held true even when both
178 male and female probands were considered together, with the average IQ of these 96 probands
179 (average IQ=69.46) being significantly lower than the simulated distribution of average IQ
180 (average IQ ranged from 71 to 92; empirical p -value = 0) (**Supplemental Fig. S6**).

181 Next, we applied RareComb to identify gene triplets associated with intellectual disability
182 using the two cohorts of affected males and repeated the simulations to identify 1,593 significant
183 combinations in the SPARK cohort. We selected 570 high-confidence triplets (with $\geq 90\%$
184 statistical power at 5% significance threshold; **Supplemental Table S6**) and found that 79
185 probands in the SSC cohort carried at least one of these deleterious triplets. The average IQ score

186 of individuals carrying significant gene triplets (average IQ score=73) was significantly lower
187 than a distribution of average IQ scores from 10,000 draws of 79 SSC probands (average IQ
188 ranged from 72 to 94; empirical $p=0.0011$; see **Supplemental Fig. S7**). This result reiterated that
189 carriers of mutations in the significant gene combinations have lower IQ than a random group of
190 probands. Our results also demonstrate the ability of the framework to identify higher order
191 combinations of mutations that are significantly associated with specific phenotypes in
192 individuals with complex disorders.

193

194 **Oligogenic combinations are enriched for specific inheritance patterns**

195 As individual variants can arise *de novo* or be inherited maternally or paternally, variants in pairs
196 of genes can have six possible patterns of transmission (**Supplemental Fig. S8A**). We identified
197 a total of 926 occurrences of the 148 pairs of mutated genes enriched among SPARK probands
198 with ID ($n=142$ probands), of which inheritance could be determined without ambiguity for 887
199 instances. We found that one variant occurred *de novo* and the other variant was inherited from
200 the mother in 244/887 instances (27.5%). Similarly, both mutated genes were inherited from the
201 mother in 226/887 instances (25.4%) or occurred *de novo* in 221/887 instances (24.9%), while
202 the remaining fraction (~22%) of variant pairs were either inherited from both parents, inherited
203 from the father, or transmitted *de novo* and paternally. To assess the significance of our
204 observations, we performed simulations to establish a baseline expectation of proportions for
205 each category of parental inheritance pattern. We selected 926 pairs of genes in 1000 random
206 draws of all possible mutated gene pairs among SPARK probands and calculated the fraction of
207 instances that fell into each of the six transmission categories. Unaffected siblings were not
208 considered for this simulation. The observed proportion was higher than the simulated
209 proportions for instances when both variants occurred *de novo* (24.9% versus 17%, empirical
210 $p=0$) and when one variant was *de novo* and the other was inherited maternally (27.5% versus
211 25%, $p=0.028$) (**Figure 3A**). We note that the depletions observed in categories ‘Maternal +
212 Paternal’ and ‘Both Paternal’ could simply be due to the numerical offset resulting from
213 enrichment of other categories. We repeated this analysis for 7,596 children affected with autism
214 in the SPARK cohort compared to 11,740 unaffected parents and identified 110 gene pairs
215 significantly associated with autism (**Supplemental Table S7**). Similar to the results obtained
216 for the ID phenotype, we found that both variants of a gene pair were more likely to occur *de*

217 *novo* (24% versus 18%, empirical $p=0$) or one variant occurring *de novo* and the other inherited
218 maternally (33% versus 26%, $p=0$) than expected based on simulation studies (**Supplemental**
219 **Fig. S9**). The enrichment of *de novo* or maternally inherited variants for significant gene pairs
220 aligns with published reports that severely affected children tend to carry multiple *de novo*
221 mutations or inherit pathogenic rare variants from mildly affected or unaffected carrier mothers
222 (Girirajan et al. 2012; Turner et al. 2017; Krumm et al. 2015).

223 We then assessed whether the mutated gene pairs associated with ID were also found in
224 siblings of carrier probands. Restricting our analysis to families with unaffected siblings whose
225 probands had mutations in ID-enriched gene pairs, we found that both variants were present in
226 the corresponding sibling for only 53/219 (24.2%) instances of gene pairs, while 102/219
227 (46.6%) had variants in only one of the two genes and 64/219 (29.2%) instances had no variants
228 in either of the genes in the siblings (**Supplemental Fig. S8B**). Using simulations, we found a
229 significantly higher proportion of instances with only one of the two variants present in siblings
230 compared to the expected values (46.6% versus 38.5%, $p=0.007$). Furthermore, the proportion of
231 observed instances with neither of the variants present in siblings (29.2% versus 33.1%,
232 empirical $p=0.098$) or both variants present in siblings (24.2% versus 28.4%, $p=0.079$) was
233 lower than expected (**Figure 3B**). The observation that only a small fraction of unaffected
234 siblings carried both mutated gene pairs suggests a strong association of these gene pairs with ID
235 phenotypes. These results suggest that mutations in pairs of genes significantly associated with a
236 severe phenotype in probands are more likely to occur individually than simultaneously in
237 unaffected siblings of the same family.

238

239 **Genes forming oligogenic combinations are distinct from canonical autism genes**

240 We expanded our analysis to include all 16,556 mutated genes in the SPARK male cohort, as
241 opposed to genes with mutations present in both the SPARK and SSC male cohorts, and
242 identified 52 significant gene pairs (**Supplemental Table S8**) and 230 triplets associated with
243 the ID phenotype (with $\geq 90\%$ statistical power at 1% significance threshold; **Supplemental**
244 **Table S9**). Due to the expanded search space, the 52 mutated gene pairs showed more significant
245 p-values from the binomial tests when compared to those obtained from the more restricted set of
246 variants overlapping both SPARK and SSC cohorts (**Supplemental Fig. S10**). Mutated genes
247 within the 52 combinations included several genes related to nervous system development, such

248 as *NIN*, *HDC*, *NGF*, and *BRD8*. Furthermore, 5/52 pairs and 59/230 triplets contained at least
249 one gene associated with autism in the SFARI database, including *FGFR1*, associated with
250 multiple disorders including Kallmann syndrome (Dodé et al. 2003) and Pfeiffer syndrome
251 (Schell et al. 1995); *RELN*, associated with temporal lobe epilepsy (Dazzo et al. 2015); *SYNE1*,
252 associated with spinocerebellar ataxia (Yoshinaga et al. 2017; Synofzik et al. 2016); and
253 *PNPLA7*, associated with autism and ID (Prasad et al. 2012). Thus, most genes forming
254 combinations are not involved in canonical autism or ID disorders, suggesting synergistic effects
255 of these genes without prior association to disease. We also analyzed 14,708 variants from 1,528
256 female probands (375 probands with ID and 1,153 without ID) and identified 19 significant pairs
257 associated with ID, indicating that significant combinations can be identified even when the
258 sample sizes are small (**Supplemental Table S10**).

259 We performed Gene Ontology enrichment analysis for genes within the 52 combinations
260 and identified seven out of nine significantly enriched GO terms to be exclusively associated
261 with nervous system-related functions, including synthesis and metabolism of catecholamines,
262 axon/neuron regeneration, and neuron generation and differentiation (**Supplemental Fig. S11**)
263 (Mi et al. 2019). Enrichment of several annotations related to growth and maintenance of brain
264 cells such as ‘axon development’, ‘neurogenesis’, ‘axon regeneration’, ‘neuron differentiation’,
265 ‘neuron projection regeneration’, and ‘response to axon injury’ indicate the physiological
266 relevance of the genes identified by our method. Furthermore, the differences in the type and
267 specificity of GO terms enriched for significant pairs versus triplets were apparent, with genes
268 forming pairs involved in nervous system function and genes forming triplets associated with
269 both nervous system as well as other biological processes. We next assessed the enrichment and
270 depletion of Human Phenotype Ontology (HPO) terms for genes forming significant pairs
271 towards ID phenotypes (Köhler et al. 2021). First, we calculated the fraction of all 4,484 genes
272 within the HPO database associated with each HPO term. For example, 30% (1,366/4,484) of all
273 genes in HPO were associated with ID. We compared these expected values calculated for each
274 HPO term with the corresponding fractions observed within the 95 genes forming 52 ID-
275 associated pairs using binomial tests. Genes associated with HPO terms related to
276 neurodevelopmental phenotypes, such as ID, global developmental delay, seizure, and
277 microcephaly, were significantly depleted within the set of 95 genes forming gene pairs
278 (**Supplemental Table S11**). Next, we evaluated whether genes within each of the 52 significant

279 pairs shared one or more common HPO phenotype or disease. Of the 52 pairs, only one pair
280 (*DNASE1* & *MTR*) shared an HPO phenotype (“epilepsy”). This was significantly lower than the
281 expected value obtained from the distribution of the number of shared HPO phenotypes between
282 all possible pairs of genes in the HPO database (1/52, 1.9% ID gene pairs compared to 31.5% of
283 all HPO gene pairs shared one HPO phenotype, $p=2.2\times 10^{-16}$; one-sided binomial test)
284 (**Supplemental Fig. S12; Supplemental Table S12**). We note that the 4,484 genes within HPO
285 are potentially biased towards well-studied disorders, making pairs of genes drawn from HPO
286 more likely to share phenotypes than random pairs of genes from the genome. Overall, GO and
287 HPO analyses show that genes forming oligogenic combinations are involved in neuronal
288 processes but have not been previously connected to neurodevelopmental phenotypes, indicating
289 the novelty of the associations between these genes and ID phenotypes.

290

291 **Identifying variant combinations towards specific patterns of comorbid phenotypes**

292 We adapted our framework to identify significant associations of two or more genotypes with
293 multiple comorbid phenotypes. To identify novel comorbid associations, we eliminated
294 phenotypes that were highly correlated with each other, such as ADHD and reading disorder
295 (Gilger et al. 1992). We analyzed variant profiles of 6,189 autism probands from the SPARK
296 cohort with records of comorbid features, including 1,215 individuals with ID, 1,825 with
297 anxiety and depression, and 332 with schizophrenia features. We assessed for significant co-
298 occurrences of two or more mutated genes with two or more of the above phenotypes (**Figure 4**).
299 Using one-tailed binomial tests to compare the observed frequency of combinations of genotypes
300 and phenotypes to the expected frequency, we first identified 169 significant associations
301 between pairs of mutated genes and two comorbid phenotypes as well as 82 combinations of
302 three mutated genes and two comorbid phenotypes (**Supplemental Tables S13 & S14**). As some
303 of these significant genotype-phenotype combinations can be confounded by high degree of co-
304 occurrence of mutated genes, we next calculated genotype-only p-values using binomial tests for
305 all significant genotype-phenotype associations. For 32/169 combinations of two mutated genes
306 and two comorbid phenotypes and 5/82 combinations of three mutated genes and two comorbid
307 phenotypes, the composite genotype-phenotype p-values were significant while genotype-only p-
308 values were not significant, suggesting stronger associations between these variant combinations
309 and phenotypes. For example, even when variants in genes *COL28A1* and *MFSD2B* did not co-

310 occur more frequently than expected under the assumption of independence, these mutated genes
311 co-occurred more frequently than expected among probands with ID and schizophrenia
312 phenotypes. Loss-of-function and rare missense mutations in *COL28A1* have been reported in
313 individuals with autism (Krumm et al. 2013; Guo et al. 2017), and *MFSD2A*, a paralog of
314 *MFSD2B*, has been directly implicated in an autosomal recessive disorder associated with
315 progressive microcephaly, spasticity and brain imaging abnormalities (Guemez-Gamboa et al.
316 2015). Similarly, we found *ARVCF* and *FATI* to be significantly associated with ID and
317 schizophrenia, with *ARVCF* mapping within the 22q11.2 DiGeorge syndrome region (Sanders et
318 al. 2005), while rare *de novo* mutations in *FATI* being associated with autism and schizophrenia
319 (Iossifov et al. 2014; Kenny et al. 2014). Finally, we found that the mutations in genes *ABCA4*,
320 *DNAH10* and *MC1R* significantly co-occurred in individuals with ID and anxiety/depression
321 phenotypes. These results demonstrate the utility of identifying higher-order associations
322 between genotypes and phenotypes in complex disorders such as autism.

323

324 **DISCUSSION**

325 Current rare variant analysis strategies are geared towards either searching for individual variants
326 of high effect size whose influence on the phenotype is evident, such as *de-novo* gene-disruptive
327 mutations, or comparing rare variant burden to explain collective effects on phenotypes (Sebat et
328 al. 2007; Zheng et al. 2016; Girirajan et al. 2011). The wider space between these two extremes
329 of the analysis spectrum that involves combinations of rare variants has largely remained
330 understudied. Although digenic diseases and multi-hit models of complex diseases have been
331 used to provide post-hoc explanations for an observed phenomenon, they are not equipped to
332 serve as a framework to actively search and identify rare variant combinations that fit oligogenic
333 models for specific phenotypes (Gifford et al. 2019; Badano et al. 2006; Leblond et al. 2012).
334 While machine learning has become the de-facto approach for disease outcome predictions, the
335 lack of holy-grail predictors and reduced interpretability due to data sparsity makes it less fit to
336 detect combinatorial effects (Murdoch et al. 2019). In addition, the common practice of
337 evaluating feature importance metrics of machine learning classifiers falls short of the objective
338 to identify combinations of features that exert higher effect on the phenotype than evident from
339 their independent effects (Molnar et al. 2020; Murdoch et al. 2019). Even if this black-box nature
340 of machine learning could be overcome, identifying even a handful of truly oligogenic variants

341 to use as ‘ground truth’ for training a classifier can be challenging. While few studies managed
342 to use a small number of examples as training sets effectively, those approaches were limited to
343 digenic models (Papadimitriou et al. 2019). Furthermore, prior studies to assess combinatorial
344 effects have been inherently biased due to their need to minimize the search space by restricting
345 the analysis to only a subset of genes chosen based on *a priori* knowledge (Papadimitriou et al.
346 2019; Kerner et al. 2020; Schaaf et al. 2011). Here, we provide a proof-of-concept analytical
347 framework that remains agnostic to prior evidence and performs exhaustive searches to identify
348 combinatorial effects among rare variants while retaining high granularity of data and
349 interpretability of results.

350 Here, we use our framework to identify gene pairs and triplets significantly associated
351 with intellectual disability and show that several constituent genes are associated with nervous
352 system processes. These mutated gene combinations are more likely to be inherited maternally or
353 occur *de novo*, are depleted in unaffected siblings from the same family, and are less likely to
354 involve canonical autism or ID genes, suggesting that genes forming significant combinations
355 are less deleterious on their own but manifest effects only when combined with other similar
356 genes carrying rare mutations. While previous studies have linked aggregate rare variant burden
357 towards intellectual disability (Singh et al. 2017; Fitzgerald et al. 2015), our results fine map the
358 association to specific combinations of constituent genes contributing to the burden. Based on
359 these observations, we propose a novel paradigm for dissecting the complexity of genetic
360 disorders, where an affected individual carries multiple combinations of rare variants, and each
361 combination contributes to either the same phenotype or distinct phenotypes at varying effect
362 sizes (**Figure 5**).

363 A limitation of our method is that it tends to be biased towards genes that are mutated
364 frequently enough to be observed in a combination, and therefore variant types such as large
365 structural variants were not included in our analysis. This limitation can be addressed by fixing
366 specific primary variants of interest irrespective of their frequency and screening for “second-
367 hit” modifiers that significantly co-occur with the primary variant, such as the co-occurrence of
368 *RBM8A* variants in proximal 1q21.1 deletion carriers manifesting thrombocytopenia-absent-
369 radius syndrome, and *TBX6* variants in proximal 16p11.2 deletion carriers with scoliosis (Albers
370 et al. 2012; Yang et al. 2019). Another limitation of our method is that it does not take
371 population substructure into account. Because allele frequencies of rare variants vary across

372 populations and ancestries, our results are likely true for individuals of European descent that
373 make up a majority of the cohort than other subgroups. Future studies applying our method to
374 more heterogeneous populations should consider taking the population substructure into account
375 prior to making inferences. Even though our analyses did not consider prior functional
376 knowledge of genes from co-expression or protein interaction networks, future studies can refine
377 the results and infer biological significance during post-processing, depending on the research
378 questions and contexts. Alternatively, if the objective is to only find mutation combinations
379 within specific pathways, functions, or interaction networks, our method will still be effective if
380 the input is pre-processed to include a specific set of mutated genes based on prior knowledge.

381 Our method is fast and scalable, allows for fine-tuning combinatorial searches based on
382 frequency, statistical power, and multiple testing criteria, and can be adapted to enable
383 computational approximations to further improve run time and assess higher-order combinations
384 beyond triplets. While larger sample sizes are generally required for detecting smaller frequency
385 differences, we note that our framework achieves reliable statistical power even with modest
386 sample sizes, implying that our framework could be applied to exome sequencing studies of
387 other neurodevelopmental disorders that have not been explored for combinatorial effects. This
388 approach can also be used to address a variety of research questions involving rare event
389 combinations, including searching for protective effects of rare variants where simultaneous
390 mutations are enriched in controls but not in cases, and finding combinations that exhibit specific
391 enrichment or depletion patterns in more than two phenotypic groups. In summary, we provide a
392 conceptual framework and the necessary tools to identify the oligogenic basis for complex
393 disorders such as autism and intellectual disability, which hitherto was restricted to the analysis
394 of canonical disorders such as Hirschsprung disease (Gabriel et al. 2002) and Bardet-Biedl
395 syndrome (Badano et al. 2006).

396 **METHODS**

397 We developed RareComb to address computational and statistical challenges associated with
398 combinatorial analysis of rare variants. RareComb first uses the apriori algorithm to efficiently
399 count the frequencies of co-occurring variant combinations. It then uses one-tailed binomial tests
400 to compare the observed frequency of each variant combination to the expected frequency
401 derived under the assumption of independence among the constituent variants within each
402 combination (**Figure 1**). This method can be applied to identify variant combinations that are
403 significantly enriched in cases but not in controls. In studies involving multiple comorbid
404 phenotypes, this method can also be used to detect associations between specific combinations of
405 variants and one or more (comorbid) phenotypes (see **Supplemental Material**). The general
406 principles of our method, built using the basic axioms of probability theory, can be easily
407 extended to a variety of problems involving rare higher-order combinations (**Supplemental Fig.**
408 **S13**).

410 **Identifying frequencies of rare variant combinations**

411 RareComb utilizes the apriori algorithm to efficiently calculate frequencies of variant
412 combinations from sparse Boolean matrices (of 0s and 1s) (**Supplemental Fig. S14A**). The
413 apriori algorithm has been successfully applied to analyze consumer behavior, where identifying
414 products frequently purchased together could benefit a company (Brijs et al. 1999; Glance et al.
415 2005). While an algorithm that is used to derive insights from patterns within highly frequent
416 events (i.e., frequent itemset mining) might not seem like a good fit to analyze rare variant
417 combinations, its ability to perform disciplined search based on both built-in and user-specified
418 constraints makes it an ideal counting tool. For example, the apriori algorithm avoids
419 enumerating each of the 50 million pairs or 167 billion triplets from just 10,000 variants, and
420 instead prunes the search-space based on user-defined criteria such as minimum frequency
421 threshold and size of combinations (pairs, triplets, etc.) (**Supplemental Fig. S14B**). RareComb
422 applies an additional constraint to the algorithm to limit its search to co-occurring events, which
423 further reduces the search space (see **Supplemental Material**). For example, when considering
424 variants A and B, only the frequency of the presence of both variants ($A=1 \ \& \ B=1$) is counted,
425 and not absence of either or both variants ($A=1 \ \& \ B=0$; $A=0 \ \& \ B=1$; or $A=0 \ \& \ B=0$).

426

427 **Statistical Inference**

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

433
$$H_0: \pi = \pi_0$$

434
$$H_a: \pi > \pi_0$$

435 where,

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

438 π_0 = *Expected* probability derived from the frequency of individual variants of a
 439 combination, under the assumption of independence, i.e., $P(A=1) * P(B=1)$.

440

441 RareComb then compares the null binomial distribution derived using the sample size of the
 442 group (n) and the expected probability (π_0) (i.e., $X \sim \text{Binom}(n, p = \pi_0)$) with the observed
 443 probability (π), and calculates the probability of observing rare variants occurring together at
 444 least as frequently as they were observed within the cohort (i.e. p-value).

445

446 In case-control analyses, binomial test is applied independently to each group, and the p-
 447 values between them are compared. The combinations exhibiting enrichment in both cases and
 448 controls, likely due to proximity of variants in linkage disequilibrium, are eliminated, following
 449 which the p-values in cases are adjusted for multiple-testing to identify statistically significant
 450 combinations that exhibit enrichment in cases but not in controls. For a more conservative
 451 approach, multiple testing adjustment is applied earlier in the RareComb pipeline by considering
 452 the total number of combinations that meet the minimum frequency threshold in cases as the
 453 total number of tests. Once adjusted for multiple testing, combinations with significant p-values
 454 in cases but not in controls are selected as significant. While Bonferroni corrections were used
 455 for all our analysis, our method provides users the flexibility to use Benjamini-Hochberg
 456 corrections as well. Finally, the effect sizes are calculated using Cohen's d and the statistical
 457 power is measured using 2-sample 2-proportion tests, as additional metrics to prioritize the final

458 set of significant rare variant combinations. In genotype-comorbid phenotype association
459 analyses, the method is applied just once to the entire cohort, with multiple-testing adjusted p-
460 values serving as a sufficient metric to identify high quality associations between genotypes and
461 two or more co-occurring phenotypes.

462

463 **Statistical power and computational performance of the method**

464 We measured the relationship between sample size and statistical power for both binomial and 2-
465 sample 2-proportion tests used in the framework. It took 1,356 samples for the binomial test to
466 achieve a statistical power of 80% to establish statistical enrichment between expected and
467 observed co-occurrence frequencies of 0.1% and 0.5% (**Supplemental Fig. S15**). This number
468 increased to 6,469 when the test needed to be more sensitive to compare frequencies of 0.3% and
469 0.5%. Similarly, it took 7,840 samples for the 2-sample 2-proportion test to achieve 80% power
470 to establish statistical difference between co-occurrence frequencies of 2% and 0.5% observed in
471 two groups (**Supplemental Fig. S16**). The sample size requirement increased to 14,633 to
472 differentiate frequencies of 1.5% and 0.5% at 80% statistical power. Furthermore, increasing the
473 size of the control cohort alone could increase the statistical power to identify significant
474 differences in proportions between cases and controls (**Supplemental Fig. S17**). For example, if
475 a particular variant combination is observed 5 times in 500 cases and 500 controls, the statistical
476 power available for the 2-sample 2-proportion test is 1%, but the power increases to 64% if the
477 combination is instead observed 5 times in 5,000 controls. These results align with the known
478 relationship between sample size and statistical power and indicate that our method can be
479 reliably applied to analyze reasonably modest-size cohorts.

480 We also measured the run times for the case-control analysis to identify significant pairs
481 and triplets of mutated genes using simulated data of three discrete sizes of samples (5,000,
482 10,000, and 50,000 individuals) and genes (5,000, 10,000, and 15,000 genes). The apriori
483 algorithm was run on single-core CPUs with 256 GB memory and was constrained to analyze
484 combinations observed in at least 0.15% of the samples. Given the memory-intensive nature of
485 the apriori algorithm implemented in the ‘arules’ package, 256 GB was chosen to maintain
486 uniformity (Hahsler et al. 2005). However, smaller input files could be processed successfully
487 using much less memory. As expected, the runtimes were proportional to the size of the
488 combination (pairs versus triplets) and the number of input variables (**Supplemental Fig. S18**).

489 While the increase in run time with the increase in sample size is apparent for pairs, lower
490 runtimes observed with running 50,000 samples compared to 5,000 samples for triplets can be
491 attributed to stochasticity of the input data. Overall, the analysis of gene pairs took between one
492 minute and 12 minutes while triplets took between two minutes and 150 minutes. Since several
493 factors influence the runtime of the method, a trial-and-error approach to determine an optimal
494 minimum frequency threshold for co-occurring events can help identify relevant combinations
495 without resulting in insufficient memory due to combinatorial explosion.

496

497 **Samples**

498 We used whole exome sequencing data from 6,189 affected males from the Simons Foundation
499 Powering Autism Research (SPARK) (The SPARK Consortium 2018) and 1,878 affected males
500 from 2,247 simplex families from the Simons Simplex Collection (SSC) (Sanders et al. 2015)
501 cohort from the Simons Foundation Autism Research Initiative (SFARI). We also used whole
502 exome sequencing data from 1,528 affected females from the SPARK cohort for the female-
503 specific analysis, and the entire SPARK cohort of 7,717 affected males and females were
504 considered for the combined analysis. For the inheritance analysis, 121 samples with low-
505 confidence autism diagnosis were removed and only 7,596 samples were considered. While
506 clinical diagnosis information for intellectual disability (ID), anxiety, attention deficit
507 hyperactivity disorders (ADHD), schizophrenia, language and sleep disorders were encoded as
508 binary variables for the SPARK samples, full-scale intelligence quotient (IQ) scores were
509 available for the SSC cohort. Although the entire SPARK cohort is composed of individuals with
510 autism diagnosis, for the purposes of this study, individuals with a clinical diagnosis of
511 intellectual disability/cognitive impairment were labeled as *cases* and those without the ID
512 diagnosis as *controls*. Continuous variables such as IQ scores were not used for case/control
513 classification of the SPARK cohort.

514

515 **Data preparation and quality control**

516 All SPARK exome sequencing samples were aligned using the hg38 reference genome. Variant
517 Call Format (VCF) files for these samples were annotated using ANNOVAR (Wang et al. 2010)
518 for rsID information and variant frequency using ExAC (Exome aggregation Consortium 2016)
519 and gnomAD (Genome Aggregation Database Consortium 2020). To overcome the limitations of

520 using a single method to assess the effect of non-synonymous mutations, pathogenicity predicted
521 by the following 11 methods were collectively obtained from dbNSFPv3.0a (Liu et al. 2016) and
522 annotated using ANNOVAR: SIFT (Ng and Henikoff 2003), PolyPhen-2 (Adzhubei et al. 2010)
523 (HDIV), PolyPhen-2 (HVAR), LRT (Chun and Fay 2009), MutationTaster (Schwarz et al. 2010),
524 MutationAssessor (Reva et al. 2011), FATHMM (Shihab et al. 2013), MetaSVM (Kim et al.
525 2017), PROVEAN (Choi and Chan 2015), REVEL (Ioannidis et al. 2016), and CADD v1.3
526 (Rentzsch et al. 2019). We also note that the limitations associated with using different versions
527 of annotation tools (such as CADD v1.3 versus v1.6) are also overcome by our strategy of using
528 11 different pathogenicity predictors for our analysis. Briefly, all missense, stop-loss/gain, and
529 start-loss/gain variants within exonic regions with minor allele frequencies $\leq 1\%$ identified based
530 on both ExAC and gnomAD databases were selected. Then, variants with allele depth of ≥ 15
531 and allele balance between 25% and 75% for heterozygous variants and $> 90\%$ for homozygous
532 variants were selected as high-quality variants. Deleteriousness of the variants were measured
533 and reported differently by each prediction method. REVEL provided a score between 0 and 1,
534 with higher scores indicating higher level of deleteriousness, while PolyPhen-2 and
535 MutationAssessor classified variants into one of three categories. For example, PolyPhen-2
536 classified variants as ‘Deleterious’, ‘Possibly damaging’, or ‘Tolerated’, while MutationAssessor
537 classified variants as ‘High’, ‘Medium’, or ‘Low’. The other nine methods classified variants as
538 either ‘Deleterious’ or ‘Tolerated’. Pathogenicity reported by each tool was encoded as a binary
539 variable, with the categories ‘Possibly damaging’ and ‘Medium’ encoded as 0.5. Thus, the
540 composite pathogenicity score derived from the 11 tools could range between 0 and 11. Missense
541 variants with a cumulative score of ≥ 4 and stop-loss/gain predicted as ‘deleterious’ either based
542 on CADD score (CADD v1.0 Phred > 30) or MutationTaster were considered deleterious for all
543 analyses. Indels and other smaller structural variants were not considered, as their functional
544 impact could not be easily assessed.

545

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

547 Gene Ontology term enrichment analyses were performed using the ‘Gene Ontology API’
548 accessed using the ‘post’ command of the Python package ‘requests’ (Python version 3.7) (Mi et
549 al. 2019). All analyses were performed using parameters for *Homo sapiens* (organism = ‘9606’)
550 to identify biological processes enrichment (annotDataSet = ‘GO:0008150’) using binomial tests.

551 HPO enrichment analyses were performed using data from the ‘genes_to_phenotype’ file
552 obtained from the HPO website (Köhler et al. 2021). Since enrichment of phenotypes is not
553 automatically evaluated by HPO, we used customized R scripts to derive baseline expectations
554 that could be compared against the actual observations to determine significance using the p-
555 values from binomial tests.

556

557 **Statistical analysis**

558 All statistical analyses were performed using R v3.6.1 (R Foundation for Statistical Computing,
559 Vienna, Austria) (R Core Team 2019) and Python (v3.7) (Van Rossum, Guido; Drake 2009). All
560 data-related plots were generated using the R package ggplot2 (Wickham 2016).

561

562 **Software Availability**

563 RareComb is available as an open-source (<https://github.com/girirajanlab/RareComb>) R package
564 that can be downloaded from the Comprehensive R Archive Network (CRAN) repository
565 (<https://cran.r-project.org/web/packages/RareComb/index.html>). It can also be installed into
566 development environments via interfaces such as Rstudio (RStudio Team 2020) using the
567 command `install.packages('RareComb')`. The entire R script repository is also provided as a
568 supplemental material (Supplemental_Code.zip). The tool provides several functionalities that
569 allow users to run the types of analyses described in this manuscript. The functionalities are as
570 follows: (1) Identify rare event combinations statistically enriched within a single group; (2)
571 Identify rare event combinations statistically enriched in cases but not in controls; (3) Identify
572 rare event combinations enriched in cases but depleted in controls; (4) Identify statistically
573 enriched rare event combinations that include at least one element from an user-supplied list; and
574 (5) Identify genotypes statistically enriched within individuals manifesting two or more
575 comorbid phenotypes. Each functionality takes a Boolean matrix as input and provides a set of
576 user-adjustable parameters to customize the analysis, and delivers the results in a tabular format
577 as csv files. Detailed instructions on the available functionalities and parameters built into
578 RareComb and their usage can be found on the GitHub page or CRAN website. A shiny app
579 illustrating the ideas behind RareComb is available online at
580 <https://girirajanlab.shinyapps.io/RareComb/> (Chang et al. 2020).

581

582 **Ethics approval and consent to participate**

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

585

586

587 **DECLARATIONS**

588 **Consent for publication**

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

590

591 **Competing interests**

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

593

594 **Acknowledgements**

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

605

606 **Authors' contributions**

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

610

611 **References**

- 612 Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, Kondrashov AS,
613 Sunyaev SR. 2010. A method and server for predicting damaging missense mutations. *Nat*
614 *Methods* **7**: 248–249. <http://dx.doi.org/10.1038/nmeth0410-248>.
- 615 Agrawal, Rakesh; Ramakriahnan S. 1994. Fast Algorithms for Mining Association Rules. *Proc*
616 *20th VLDB Conf* 487–499.
- 617 Albers CA, Paul DS, Schulze H, Freson K, Stephens JC, Smethurst PA, Jolley JD, Cvejic A,
618 Kostadima M, Bertone P, et al. 2012. Compound inheritance of a low-frequency regulatory
619 SNP and a rare null mutation in exon-junction complex subunit RBM8A causes TAR
620 syndrome. *Nat Genet* **44**: 435–439. <http://dx.doi.org/10.1038/ng.1083>.
- 621 Backman JD, Li AH, Marcketta A, Sun D, Mbatchou J, Kessler MD, Benner C, Liu D, Locke
622 AE, Balasubramanian S, et al. 2021. Exome sequencing and analysis of 454,787 UK
623 Biobank participants. *Nature* **599**: 628–634. <http://dx.doi.org/10.1038/s41586-021-04103-z>.
- 624 Badano JL, Katsanis N. 2002. Beyond mendel: An evolving view of human genetic disease
625 transmission. *Nat Rev Genet* **3**: 779–789. <http://dx.doi.org/10.1038/nrg910>.
- 626 Badano JL, Leitch CC, Ansley SJ, May-Simera H, Lawson S, Lewis RA, Beales PL, Dietz HC,
627 Fisher S, Katsanis N. 2006. Dissection of epistasis in oligogenic Bardet-Biedl syndrome.
628 *Nature* **439**: 326–330. <http://dx.doi.org/10.1038/nature04370>.
- 629 Brijs T, Swinnen G, Vanhoof K, Wets G. 1999. Using association rules for product assortment
630 decisions. *ACM*. <http://dx.doi.org/10.1145/312129.312241>.
- 631 Chang W, Cheng J, Allaire J, Xie Y, McPherson J. 2020. shiny: Web Application Framework for
632 R.
- 633 Choi Y, Chan AP. 2015. PROVEAN web server: A tool to predict the functional effect of amino
634 acid substitutions and indels. *Bioinformatics* **31**: 2745–2747.
635 <http://dx.doi.org/10.1093/bioinformatics/btv195>.
- 636 Chun S, Fay JC. 2009. Identification of deleterious mutations within three human genomes.
637 *Genome Res* **19**: 1553–1561. <http://dx.doi.org/10.1101/gr.092619.109>.
- 638 Coventry A, Bull-Otterson LM, Liu X, Clark AG, Maxwell TJ, Crosby J, Hixson JE, Rea TJ,
639 Muzny DM, Lewis LR, et al. 2010. Deep resequencing reveals excess rare recent variants
640 consistent with explosive population growth. *Nat Commun* **1**.
641 <http://dx.doi.org/10.1038/ncomms1130>.

- 642 Dazzo E, Fanciulli M, Serioli E, Minervini G, Pulitano P, Binelli S, Di Bonaventura C, Luisi C,
643 Pasini E, Striano S, et al. 2015. Heterozygous Reelin Mutations Cause Autosomal-
644 Dominant Lateral Temporal Epilepsy. *Am J Hum Genet* **96**: 992–1000.
645 <http://dx.doi.org/10.1016/j.ajhg.2015.04.020>.
- 646 Dodé C, Levilliers J, Dupont JM, De Paepe A, Le Dû N, Soussi-Yanicostas N, Coimbra RS,
647 Delmaghani S, Compain-Nouaille S, Baverel F, et al. 2003. Loss-of-function mutations in
648 FGFR1 cause autosomal dominant Kallmann syndrome. *Nat Genet* **33**: 463–465.
649 <http://dx.doi.org/10.1038/ng1122>.
- 650 Exome aggregation Consortium. 2016. Analysis of protein-coding genetic variation in 60,706
651 humans. *Nature* **536**: 285–291. <http://dx.doi.org/10.1038/nature19057>.
- 652 Fischbach GD, Lord C. 2010. The simons simplex collection: A resource for identification of
653 autism genetic risk factors. *Neuron* **68**: 192–195.
654 <http://dx.doi.org/10.1016/j.neuron.2010.10.006>.
- 655 Fitzgerald TW, Gerety SS, Jones WD, Van Kogelenberg M, King DA, McRae J, Morley KI,
656 Parthiban V, Al-Turki S, Ambridge K, et al. 2015. Large-scale discovery of novel genetic
657 causes of developmental disorders. *Nature* **519**: 223–228.
658 <http://dx.doi.org/10.1038/nature14135>.
- 659 Gabriel SB, Salomon R, Pelet A, Angrist M, Amiel J, Fornage M, Attié-Bitach T, Olson JM,
660 Hofstra R, Buys C, et al. 2002. Segregation at three loci explains familial and population
661 risk in Hirschsprung disease. *Nat Genet* **31**: 89–93. <http://dx.doi.org/10.1038/ng868>.
- 662 Genome Aggregation Database Consortium. 2020. The mutational constraint spectrum
663 quantified from variation in 141,456 humans. *Nature* **581**: 434–443.
664 <http://dx.doi.org/10.1038/s41586-020-2308-7>.
- 665 Gifford C, Ranade S, Samarakoon R, Salunga H, T.Yvanka de S, Huang Y, Zhou P, Elfenbein
666 A, Wyman S, Bui Y, et al. 2019. Oligogenic inheritance of a human heart disease involving
667 a genetic modifier. *Science* **364**: 865–870. <http://dx.doi.org/10.1126/science.aat5056>.
- 668 Gilger JW, Pennington BF, DeFries JC. 1992. A Twin Study of the Etiology of Comorbidity:
669 Attention-deficit Hyperactivity Disorder and Dyslexia. *J Am Acad Child Adolesc Psychiatry*
670 **31**: 343–348. <http://dx.doi.org/10.1097/00004583-199203000-00024>.
- 671 Girirajan S, Brkanac Z, Coe BP, Baker C, Vives L, Vu TH, Shafer N, Bernier R, Ferrero GB,
672 Silengo M, et al. 2011. Relative burden of large CNVs on a range of neurodevelopmental

- 673 phenotypes. *PLoS Genet* **7**: e1002334. <http://dx.doi.org/10.1371/journal.pgen.1002334>.
- 674 Girirajan S, Rosenfeld JA, Coe BP, Parikh S, Friedman N, Goldstein A, Filipink RA, McConnell
675 JS, Angle B, Meschino WS, et al. 2012. Phenotypic Heterogeneity of Genomic Disorders
676 and Rare Copy-Number Variants. *N Engl J Med* **367**: 1321–1331.
677 <http://dx.doi.org/10.1056/nejmoa1200395>.
- 678 Girirajan S, Rosenfeld JA, Cooper GM, Antonacci F, Siswara P, Itsara A, Vives L, Walsh T,
679 McCarthy SE, Baker C, et al. 2010. A recurrent 16p12.1 microdeletion supports a two-hit
680 model for severe developmental delay. *Nat Genet* **42**: 203–209.
681 <http://dx.doi.org/10.1038/ng.534>.
- 682 Glance N, Siegler M, Hurst M, Stockton R, Nigam K, Tomokiyo T. 2005. Deriving marketing
683 intelligence from online discussion. *Proc ACM SIGKDD Int Conf Knowl Discov Data Min*
684 **419–428**. <http://dx.doi.org/10.1145/1081870.1081919>.
- 685 Guemez-Gamboa A, Nguyen LN, Yang H, Zaki MS, Kara M, Ben-Omran T, Akizu N, Rosti
686 RO, Rosti B, Scott E, et al. 2015. Inactivating mutations in MFSD2A, required for omega-3
687 fatty acid transport in brain, cause a lethal microcephaly syndrome. *Nat Genet* **47**: 809–813.
688 <http://dx.doi.org/10.1038/ng.3311>.
- 689 Guo H, Peng Y, Hu Z, Li Y, Xun G, Ou J, Sun L, Xiong Z, Liu Y, Wang T, et al. 2017. Genome-
690 wide copy number variation analysis in a Chinese autism spectrum disorder cohort. *Sci Rep*
691 **7**: 44155. <http://dx.doi.org/10.1038/srep44155>.
- 692 Hahsler M, Grun B, Hornik K. 2005. arules – A Computational Environment for Mining
693 Association Rules and Frequent Item Sets. *J Stat Softw* **14**: 1–6.
694 <https://doi.org/10.18637/jss.v014.i15>.
- 695 Halvorsen M, Huh R, Oskolkov N, Wen J, Netotea S, Giusti-Rodriguez P, Karlsson R, Bryois J,
696 Nystedt B, Ameer A, et al. 2020. Increased burden of ultra-rare structural variants
697 localizing to boundaries of topologically associated domains in schizophrenia. *Nat Commun*
698 **11**: 1842. <http://dx.doi.org/10.1038/s41467-020-15707-w>.
- 699 Ioannidis NM, Rothstein JH, Pejaver V, Middha S, McDonnell SK, Baheti S, Musolf A, Li Q,
700 Holzinger E, Karyadi D, et al. 2016. REVEL: An Ensemble Method for Predicting the
701 Pathogenicity of Rare Missense Variants. *Am J Hum Genet* **99**: 877–885.
702 <http://dx.doi.org/10.1016/j.ajhg.2016.08.016>.
- 703 Iossifov I, O’Roak BJ, Sanders SJ, Ronemus M, Krumm N, Levy D, Stessman HA, Witherspoon

- 704 KT, Vives L, Patterson KE, et al. 2014. The contribution of de novo coding mutations to
705 autism spectrum disorder. *Nature* **515**: 216–221. <http://dx.doi.org/10.1038/nature13908>.
- 706 Keinan A, Clark AG. 2012. Recent explosive human population growth has resulted in an excess
707 of rare genetic variants. *Science* **336**: 740–743. <http://dx.doi.org/10.1126/science.1217283>.
- 708 Kenny EM, Cormican P, Furlong S, Heron E, Kenny G, Fahey C, Kelleher E, Ennis S, Tropea D,
709 Anney R, et al. 2014. Excess of rare novel loss-of-function variants in synaptic genes in
710 schizophrenia and autism spectrum disorders. *Mol Psychiatry* **19**: 872–879.
711 <http://dx.doi.org/10.1038/mp.2013.127>.
- 712 Kerner G, Bouaziz M, Cobat A, Bigio B, Timberlake AT, Bustamante J, Lifton RP, Casanova
713 JL, Abel L. 2020. A genome-wide case-only test for the detection of digenic inheritance in
714 human exomes. *Proc Natl Acad Sci U S A* **117**: 19367–19375.
715 <http://dx.doi.org/10.1073/pnas.1920650117>.
- 716 Kim S, Jhong JH, Lee J, Koo JY. 2017. Meta-analytic support vector machine for integrating
717 multiple omics data. *BioData Min* **10**. <http://dx.doi.org/10.1186/s13040-017-0126-8>.
- 718 Köhler S, Gargano M, Matentzoglou N, Carmody LC, Lewis-Smith D, Vasilevsky NA, Danis D,
719 Balagura G, Baynam G, Brower AM, et al. 2021. The human phenotype ontology in 2021.
720 *Nucleic Acids Res* **49**: D1207–D1217. <http://dx.doi.org/10.1093/nar/gkaa1043>.
- 721 Krumm N, O’Roak BJ, Karakoc E, Mohajeri K, Nelson B, Vives L, Jacquemont S, Munson J,
722 Bernier R, Eichler EE. 2013. Transmission disequilibrium of small CNVs in simplex
723 autism. *Am J Hum Genet* **93**: 595–606. <http://dx.doi.org/10.1016/j.ajhg.2013.07.024>.
- 724 Krumm N, Turner TN, Baker C, Vives L, Mohajeri K, Witherspoon K, Raja A, Coe BP,
725 Stessman HA, He ZX, et al. 2015. Excess of rare, inherited truncating mutations in autism.
726 *Nat Genet* **47**: 582–588. <http://dx.doi.org/10.1038/ng.3303>.
- 727 Leblond CS, Heinrich J, Delorme R, Proepper C, Betancur C, Huguet G, Konyukh M, Chaste P,
728 Ey E, Rastam M, et al. 2012. Genetic and functional analyses of SHANK2 mutations
729 suggest a multiple hit model of autism spectrum disorders. *PLoS Genet* **8**: e1002521.
730 <http://dx.doi.org/10.1371/journal.pgen.1002521>.
- 731 Liu X, Wu C, Li C, Boerwinkle E. 2016. dbNSFP v3.0: A One-Stop Database of Functional
732 Predictions and Annotations for Human Nonsynonymous and Splice-Site SNVs. *Hum*
733 *Mutat* **37**: 235–241. <http://dx.doi.org/10.1002/humu.22932>.
- 734 McClellan J, King MC. 2010. Genetic heterogeneity in human disease. *Cell* **141**: 210–217.

- 735 <http://dx.doi.org/10.1016/j.cell.2010.03.032>.
- 736 Mi H, Muruganujan A, Ebert D, Huang X, Thomas D. 2019. PANTHER version 14: more
737 genomes, a new PANTHER GO-slim and improvements in enrichment analysis tools.
738 *Nucleic Acids Res* **47**: 419–426. <http://dx.doi.org/10.1093/nar/gky1038>.
- 739 Molnar C, Casalicchio G, Bischl B. 2020. Interpretable Machine Learning – A Brief History,
740 State-of-the-Art and Challenges. *Commun Comput Inf Sci* **1323**: 417–431.
741 http://dx.doi.org/10.1007/978-3-030-65965-3_28.
- 742 Murdoch WJ, Singh C, Kumbier K, Abbasi-Asl R, Yu B. 2019. Definitions, methods, and
743 applications in interpretable machine learning. *Proc Natl Acad Sci U S A* **116**: 22071–
744 22080. <http://dx.doi.org/10.1073/pnas.1900654116>.
- 745 Ng PC, Henikoff S. 2003. SIFT: Predicting amino acid changes that affect protein function.
746 *Nucleic Acids Res* **31**: 3812–4.
- 747 Papadimitriou S, Gazzo A, Versbraegen N, Nachtegael C, Aerts J, Moreau Y, Van Dooren S,
748 Nowé A, Smits G, Lenaerts T. 2019. Predicting disease-causing variant combinations. *Proc*
749 *Natl Acad Sci U S A* **116**: 11878–11887. <http://dx.doi.org/10.1073/pnas.1815601116>.
- 750 Pizzo L, Jensen M, Polyak A, Rosenfeld JA, Mannik K, Krishnan A, McCready E, Pichon O, Le
751 Caignec C, Van Dijck A, et al. 2019. Rare variants in the genetic background modulate
752 cognitive and developmental phenotypes in individuals carrying disease-associated variants.
753 *Genet Med* **21**: 816–825. <http://dx.doi.org/10.1038/s41436-018-0266-3>.
- 754 Prasad A, Merico D, Thiruvahindrapuram B, Wei J, Lionel AC, Sato D, Rickaby J, Lu C,
755 Szatmari P, Roberts W, et al. 2012. A Discovery resource of rare copy number variations in
756 individuals with autism spectrum disorder. *G3 Genes, Genomes, Genet* **2**: 1665–1685.
757 <http://dx.doi.org/10.1534/g3.112.004689>.
- 758 R Core Team. 2019. R: A language and environment for statistical computing.
- 759 Rentzsch P, Witten D, Cooper GM, Shendure J, Kircher M. 2019. CADD: Predicting the
760 deleteriousness of variants throughout the human genome. *Nucleic Acids Res* **47**: D886–
761 D894. <http://dx.doi.org/10.1093/nar/gky1016>.
- 762 Reva B, Antipin Y, Sander C. 2011. Predicting the functional impact of protein mutations:
763 Application to cancer genomics. *Nucleic Acids Res* **39**: 37–43.
764 <http://dx.doi.org/10.1093/nar/gkr407>.
- 765 RStudio Team. 2020. RStudio: Integrated Development for R. <http://www.rstudio.com/>.

- 766 Sanders AR, Rusu I, Duan J, Vander Molen JE, Hou C, Schwab SG, Wildenauer DB, Martinez
767 M, Gejman P V. 2005. Haplotypic association spanning the 22q11.21 genes COMT and
768 ARVCF with schizophrenia. *Mol Psychiatry* **10**: 353–365.
769 <http://dx.doi.org/10.1038/sj.mp.4001586>.
- 770 Sanders SJ, He X, Willsey AJ, Ercan-Sencicek AG, Samocha KE, Cicek AE, Murtha MT, Bal
771 VH, Bishop SL, Dong S, et al. 2015. Insights into Autism Spectrum Disorder Genomic
772 Architecture and Biology from 71 Risk Loci. *Neuron* **87**: 1215–1233.
773 <http://dx.doi.org/10.1016/j.neuron.2015.09.016>.
- 774 Schaaf CP, Sabo A, Sakai Y, Crosby J, Muzny D, Hawes A, Lewis L, Akbar H, Varghese R,
775 Boerwinkle E, et al. 2011. Oligogenic heterozygosity in individuals with high-functioning
776 autism spectrum disorders. *Hum Mol Genet* **20**: 3366–3375.
777 <http://dx.doi.org/10.1093/hmg/ddr243>.
- 778 Schell U, Hehr A, Feldman GJ, Robin NH, Zackai EH, De Die-smulders C, Viskochil DH,
779 Stewart JM, Wolff G, Ohashi H, et al. 1995. Mutations in FGFR1 and FGFR2 cause
780 familial and sporadic pfeiffer syndrome. *Hum Mol Genet* **4**: 323–328.
781 <http://dx.doi.org/10.1093/hmg/4.3.323>.
- 782 Schwarz JM, Rödelberger C, Schuelke M, Seelow D. 2010. MutationTaster evaluates disease-
783 causing potential of sequence alterations. *Nat Methods* **7**: 575–576.
784 <http://dx.doi.org/10.1038/nmeth0810-575>.
- 785 Sebat J, Lakshmi B, Malhotra D, Troge J, Lese-martin C, Walsh T, Yamrom B, Yoon S, Krasnitz
786 A, Kendall J, et al. 2007. Strong Association of De Novo Copy Number Mutations with
787 Autism. *Science* **316**: 445–449. <http://dx.doi.org/10.1126/science.1138659>.
- 788 Shihab HA, Gough J, Cooper DN, Stenson PD, Barker GLA, Edwards KJ, Day INM, Gaunt TR.
789 2013. Predicting the Functional, Molecular, and Phenotypic Consequences of Amino Acid
790 Substitutions using Hidden Markov Models. *Hum Mutat* **34**: 57–65.
791 <http://dx.doi.org/10.1002/humu.22225>.
- 792 Singh T, Walters JTR, Johnstone M, Curtis D, Suvisaari J, Torniainen M, Rees E, Iyegbe C,
793 Blackwood D, McIntosh AM, et al. 2017. The contribution of rare variants to risk of
794 schizophrenia in individuals with and without intellectual disability. *Nat Genet* **49**: 1167–
795 1173.
- 796 Synofzik M, Smets K, Mallaret M, Di Bella D, Gallenmüller C, Baets J, Schulze M, Magri S,

- 797 Sarto E, Mustafa M, et al. 2016. SYNE1 ataxia is a common recessive ataxia with major
798 non-cerebellar features: A large multi-centre study. *Brain* **139**: 1378–1393.
799 <http://dx.doi.org/10.1093/brain/aww079>.
- 800 Taliun D, Harris DN, Kessler MD, Carlson J, Szpiech ZA, Torres R, Taliun SAG, Corvelo A,
801 Gogarten SM, Kang HM, et al. 2021. Sequencing of 53,831 diverse genomes from the
802 NHLBI TOPMed Program. *Nature* **590**: 290–299. [http://dx.doi.org/10.1038/s41586-021-](http://dx.doi.org/10.1038/s41586-021-03205-y)
803 [03205-y](http://dx.doi.org/10.1038/s41586-021-03205-y).
- 804 Tennessen JA, Bigham AW, O'Connor TD, Fu W, Kenny EE, Gravel S, McGee S, Do R, Liu X,
805 Jun G, et al. 2012. Evolution and functional impact of rare coding variation from deep
806 sequencing of human exomes. *Science* **337**: 64–69.
807 <http://dx.doi.org/10.1126/science.1219240>.
- 808 The 1000 Genomes Project Consortium. 2015. A global reference for human genetic variation.
809 *Nature* **526**: 68–74. <http://dx.doi.org/10.1038/nature15393>.
- 810 The SPARK Consortium. 2018. SPARK: A US Cohort of 50,000 Families to Accelerate Autism
811 Research. *Neuron* **97**: 488–493. <http://dx.doi.org/10.1016/j.neuron.2018.01.015>.
- 812 Turner TN, Coe BP, Dickel DE, Hoekzema K, Nelson BJ, Zody MC, Kronenberg ZN,
813 Hormozdiari F, Raja A, Pennacchio LA, et al. 2017. Genomic Patterns of De Novo
814 Mutation in Simplex Autism. *Cell* **171**: 710–722.
- 815 Uricchio LH, Zaitlen NA, Ye CJ, Witte JS, Hernandez RD. 2016. Selection and explosive
816 growth alter genetic architecture and hamper the detection of causal rare variants. *Genome*
817 *Res* **26**: 863–873. <http://dx.doi.org/10.1101/gr.202440.115>.
- 818 Van Rossum, Guido; Drake FL. 2009. Python 3 Reference Manual.
- 819 Wang K, Li M, Hakonarson H. 2010. ANNOVAR: Functional annotation of genetic variants
820 from high-throughput sequencing data. *Nucleic Acids Res* **38**: e164.
- 821 Wickham H. 2016. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York
822 <https://ggplot2.tidyverse.org>.
- 823 Wilfert AB, Turner TN, Murali SC, Hsieh P, Sulovari A, Wang T, Coe BP, Guo H, Hoekzema
824 K, Bakken TE, et al. 2021. Recent ultra-rare inherited variants implicate new autism
825 candidate risk genes. *Nat Genet* **53**: 1125–1134. [http://dx.doi.org/10.1038/s41588-021-](http://dx.doi.org/10.1038/s41588-021-00899-8)
826 [00899-8](http://dx.doi.org/10.1038/s41588-021-00899-8).
- 827 Yang N, Wu N, Zhang L, Zhao Y, Liu J, Liang X, Ren X, Li W, Chen W, Dong S, et al. 2019.

828 TBX6 compound inheritance leads to congenital vertebral malformations in humans and
829 mice. *Hum Mol Genet* **28**: 539–547. <http://dx.doi.org/10.1093/hmg/ddy358>.

830 Yoshinaga T, Nakamura K, Ishikawa M, Yamaguchi T, Takano K, Wakui K, Kosho T, Yoshida
831 K, Fukushima Y, Sekijima Y. 2017. A novel frameshift mutation of SYNE1 in a Japanese
832 family with autosomal recessive cerebellar ataxia type 8. *Hum Genome Var* **4**: 17052.
833 <http://dx.doi.org/10.1038/hgv.2017.52>.

834 Zheng GXY, Lau BT, Schnall-Levin M, Jarosz M, Bell JM, Hindson CM, Kyriazopoulou-
835 Panagiotopoulou S, Masquelier DA, Merrill L, Terry JM, et al. 2016. Haplotyping germline
836 and cancer genomes with high-throughput linked-read sequencing. *Nat Biotechnol* **34**: 303–
837 311. <http://dx.doi.org/10.1038/nbt.3432>.

838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858

859 **FIGURE LEGENDS**

860

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

871

872 **FIGURE 2: Combinations of rare variants contributing to intellectual disability (ID)**
873 **phenotype. (A)** An outline of the approach used to identify and validate mutated gene pairs and
874 triplets enriched in probands with ID is shown. We tested whether mutated gene pairs identified
875 as significant in one cohort (SPARK) are also associated with severe phenotypes in an
876 independent cohort (SSC). To test this, we obtained the mean IQ score of individuals from the
877 SSC cohort carrying significant combinations identified from the SPARK cohort. Empirical p-
878 values were then calculated based on the deviation of the mean IQ from the distribution of mean
879 IQ scores obtained from 10,000 random draws in the simulation. **(B)** The mean IQ of individuals
880 with mutated gene pairs in the SSC cohort was significantly lower (empirical p-value=0) when
881 compared to the distribution of mean IQ scores obtained from the simulation. **(C)** Histogram
882 shows the distributions of IQ scores of SSC probands who carried mutations in either of the
883 genes versus both constituent genes of the significant gene pairs. The distributions were
884 significantly different from each other (p-value = 1.302×10^{-16} , Kolmogorov-Smirnov test).

885

886 **FIGURE 3: Analysis of parental and sibling inheritance patterns of significant gene pairs**
887 **associated with ID. (A)** Fraction of all instances of significant gene pairs observed within each
888 of the six possible parental inheritance patterns (red) compared against 1,000 simulations is
889 shown (blue). During each simulation, random mutated gene pairs from the SSC cohort were

890 selected, the inheritance status of the mutations was identified, and the fraction of those instances
891 belonging to one of the six pre-defined categories was calculated. Comparing the observed
892 fractions with the simulated fractions indicate statistical enrichment for two specific inheritance
893 patterns based on empirical p-values: both variants being de novo, and one variant being de novo
894 and the other transmitted from the mother. **(B)** Histograms show the carrier status of significant
895 gene pairs in siblings of carrier probands (red) compared against 1,000 simulations (blue).
896 Among significant pairs, both genes were mutated in only 24.2% of all siblings (compared to
897 28.4% in simulations), whereas one of the two genes was mutated in 46.6% of all siblings
898 (compared to 38.5% in simulations). These results show that mutations are more likely to be
899 observed in just one of the two genes within the gene pairs and are less likely to be observed
900 simultaneously in siblings of carrier probands.

901

902 **FIGURE 4: Analysis of comorbid phenotypes using RareComb.** We analyzed the genotypes
903 of probands with anxiety/depression, ID, or schizophrenia. The heatmap shows combinations of
904 two or three mutated genes that were significantly enriched in individuals with specific patterns
905 of comorbid phenotypes compared to the expected frequency under the assumption of
906 independence.

907

908 **FIGURE 5: Rare variant models for complex disorders.** The schematic shows two models for
909 the genetic etiology of complex disorders. Circles represent rare variants present that are either
910 de novo or inherited from a parent. On the left, individual high effect de novo variants are
911 strongly associated with a phenotype of interest. On the right, rare variants within an individual
912 combine in multiple ways and contribute towards distinct phenotypes. The thickness of the
913 connecting lines denotes effect sizes, and an affected individual can carry multiple oligogenic
914 combinations of rare variants, each of which contributes to the same or distinct phenotypes. This
915 extension of the oligogenic model enables further dissection of the genetic architecture of
916 complex disorders.

917

Rare variants in probands

Probands	List of genes with rare variants	Disease
P1	G2	No
P2	G2, G3, G4	Yes
P3	G1, G2, G4	Yes
P4	G4	No
P5	G2, G3, G4, G5	Yes
P6	G1, G5	No

Boolean matrix representation

Probands	Genotype					Disease
	G1	G2	G3	G4	G5	D
P1		■				N
P2		■	■	■		Y
P3	■	■		■		Y
P4				■		N
P5		■	■	■	■	Y
P6	■				■	N

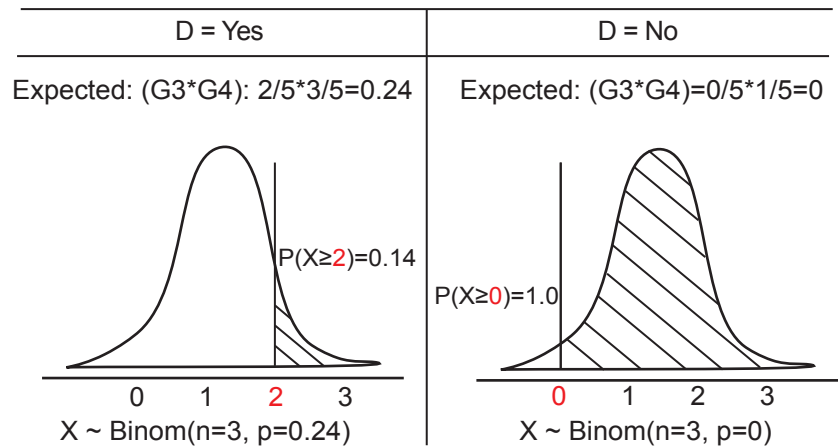
Step 1: Apriori Algorithm | Enumerate frequencies

Observed frequencies

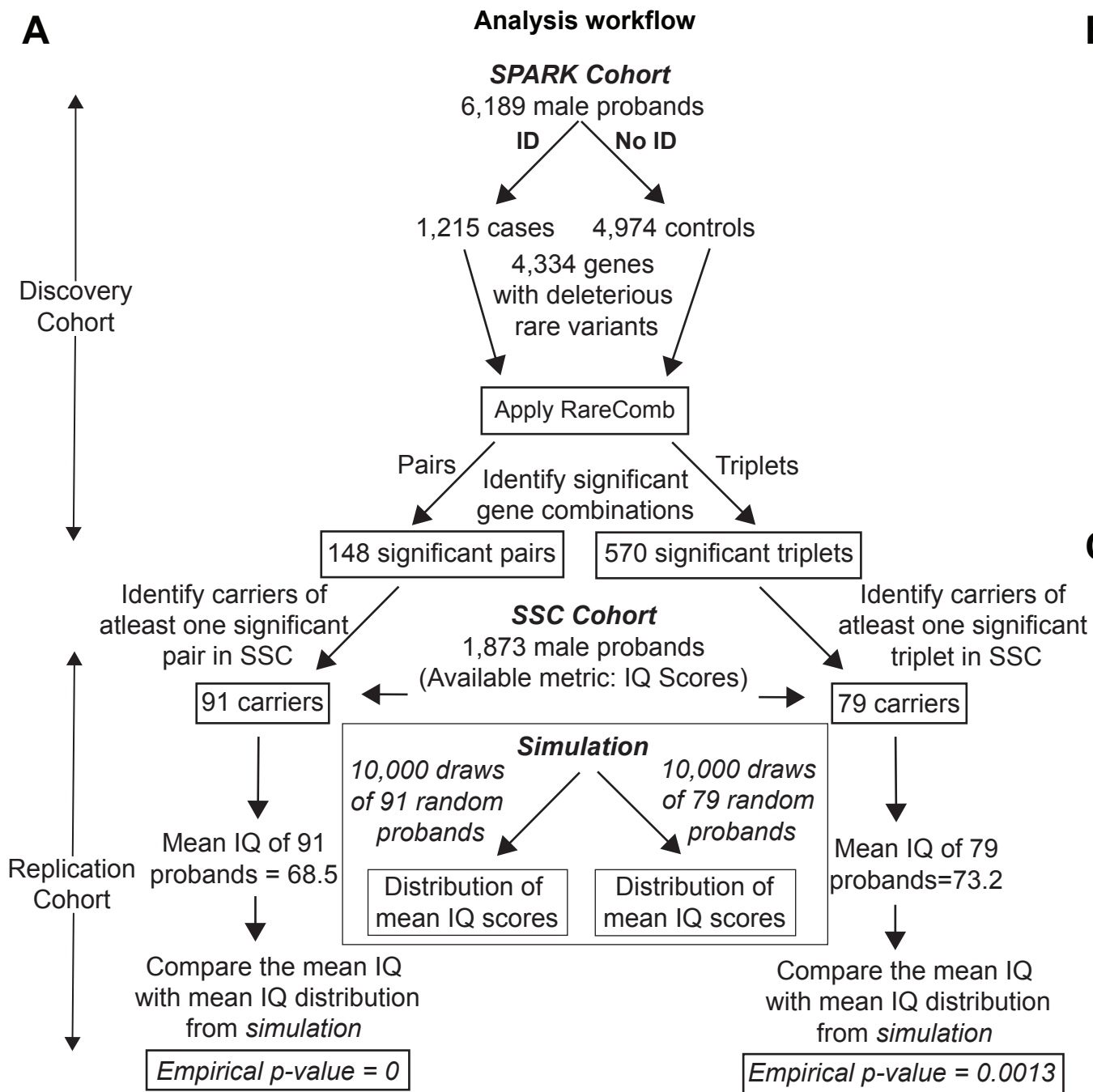
1	D = Yes	D = No
G1	1	1
G2	3	1
G3	2	0
G4	3	1
G5	1	1

(1,1)	D = Yes	D = No
(G1,G2)	1	0
(G1,G4)	1	0
(G2,G3)	2	0
(G2,G4)	3	0
(G2,G5)	1	0
(G3,G4)	2	0
(G3,G5)	1	0
(G4,G5)	1	0

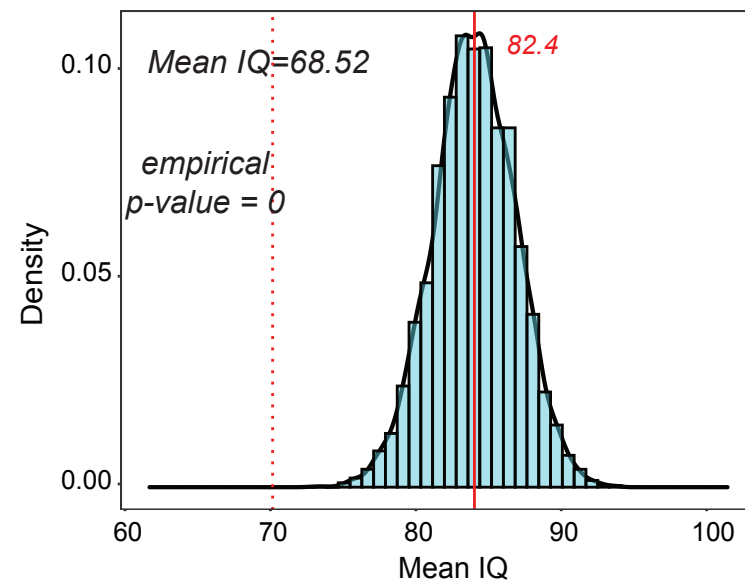
Step 2: Binomial Tests | Calculate observed vs. expected



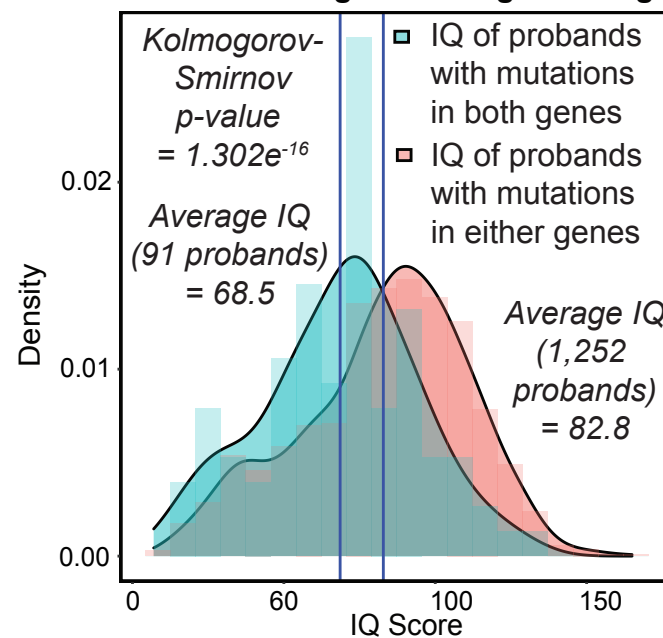
Null binomial distributions

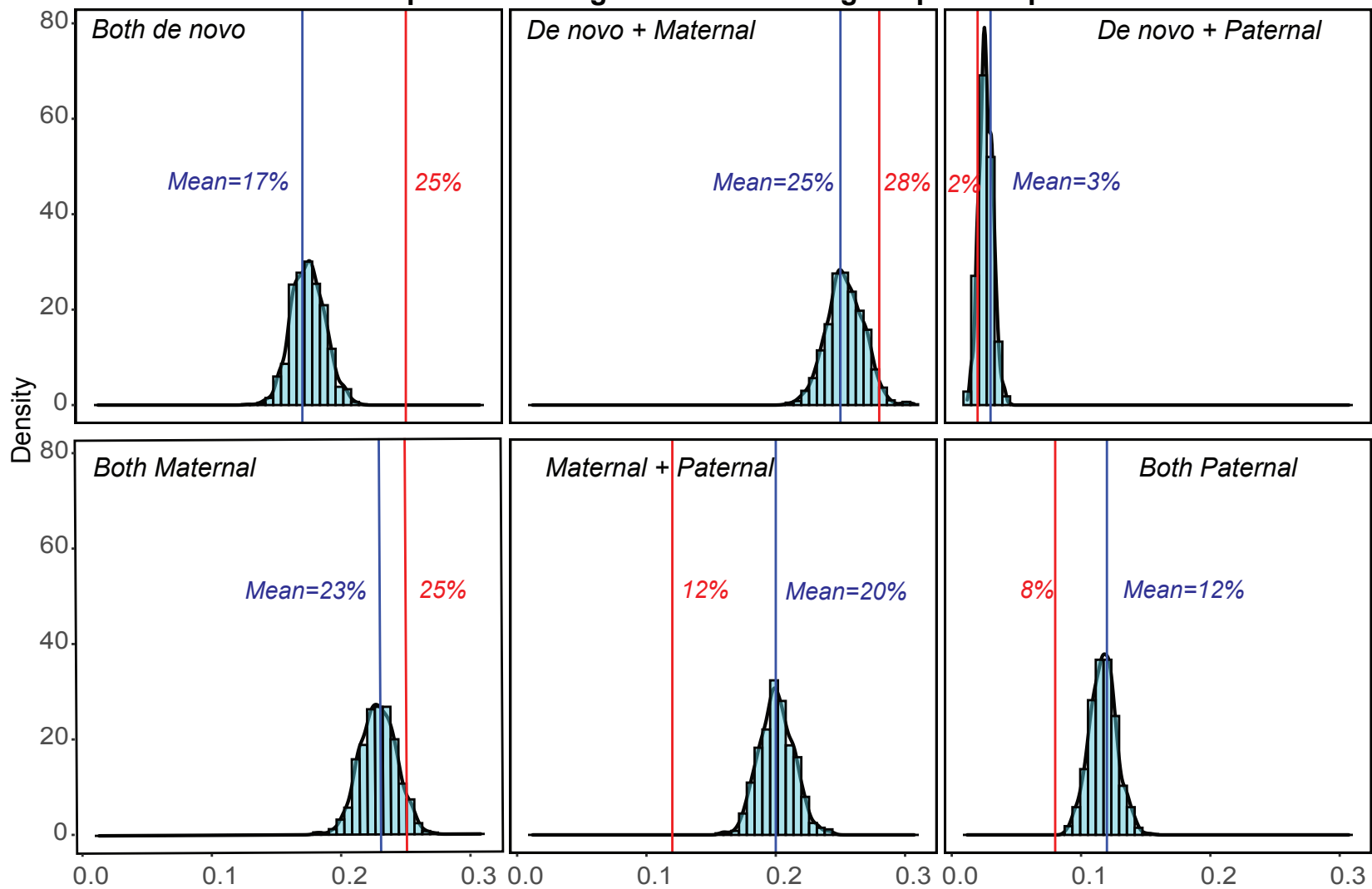
A**B**

Average IQ of individuals with significant mutated gene pairs compared to simulated data

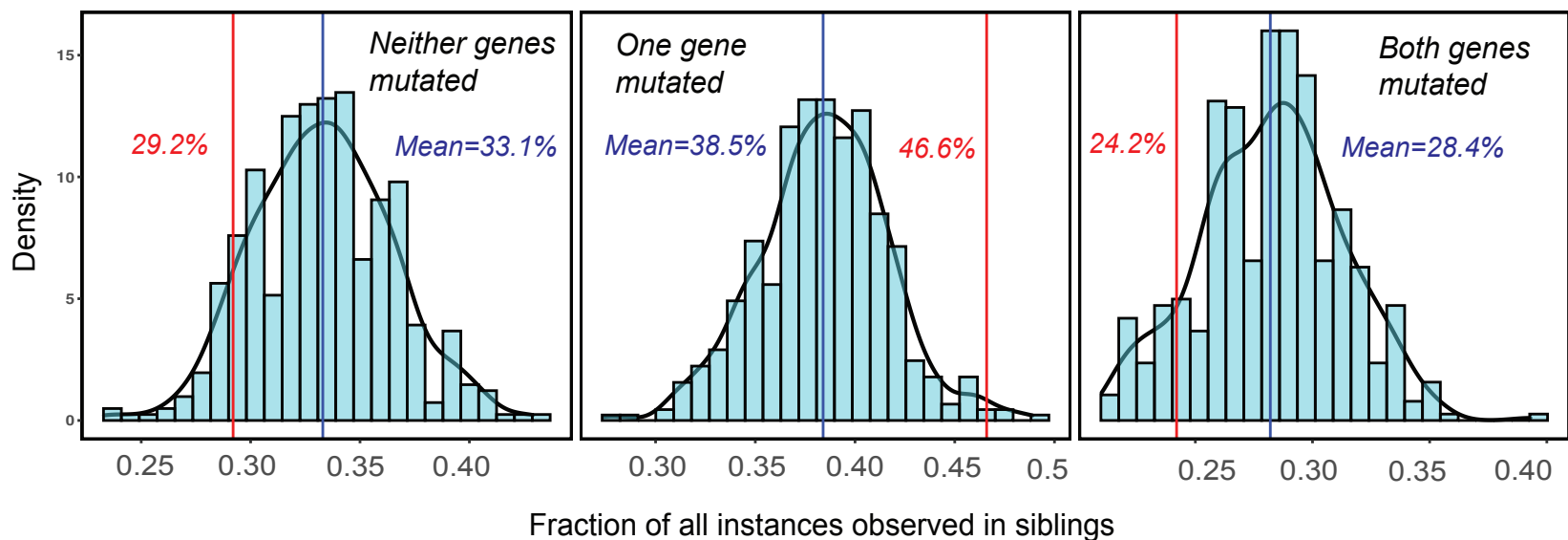
**C**

IQ distributions of individuals with mutations in either versus both genes in significant gene pairs



A**Inheritance patterns of significant mutated gene pairs in probands**

Fraction of all instances with the inheritance pattern in probands with intellectual disability

B**Enrichment and depletion of significant mutated gene pairs in siblings**

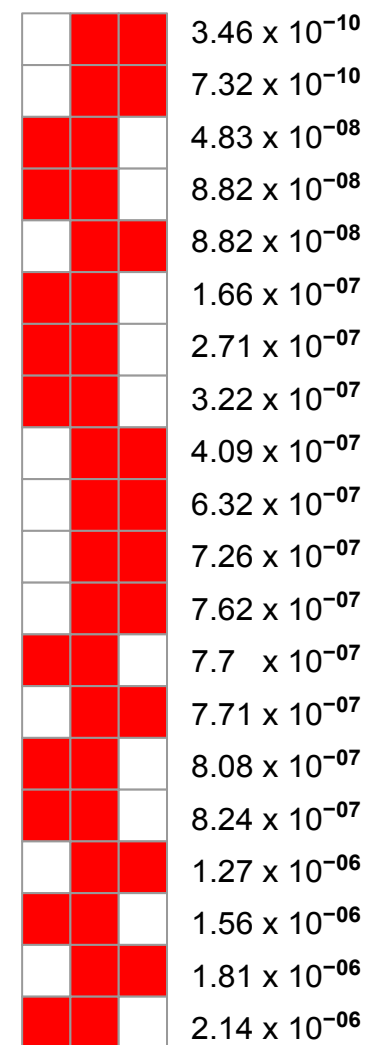
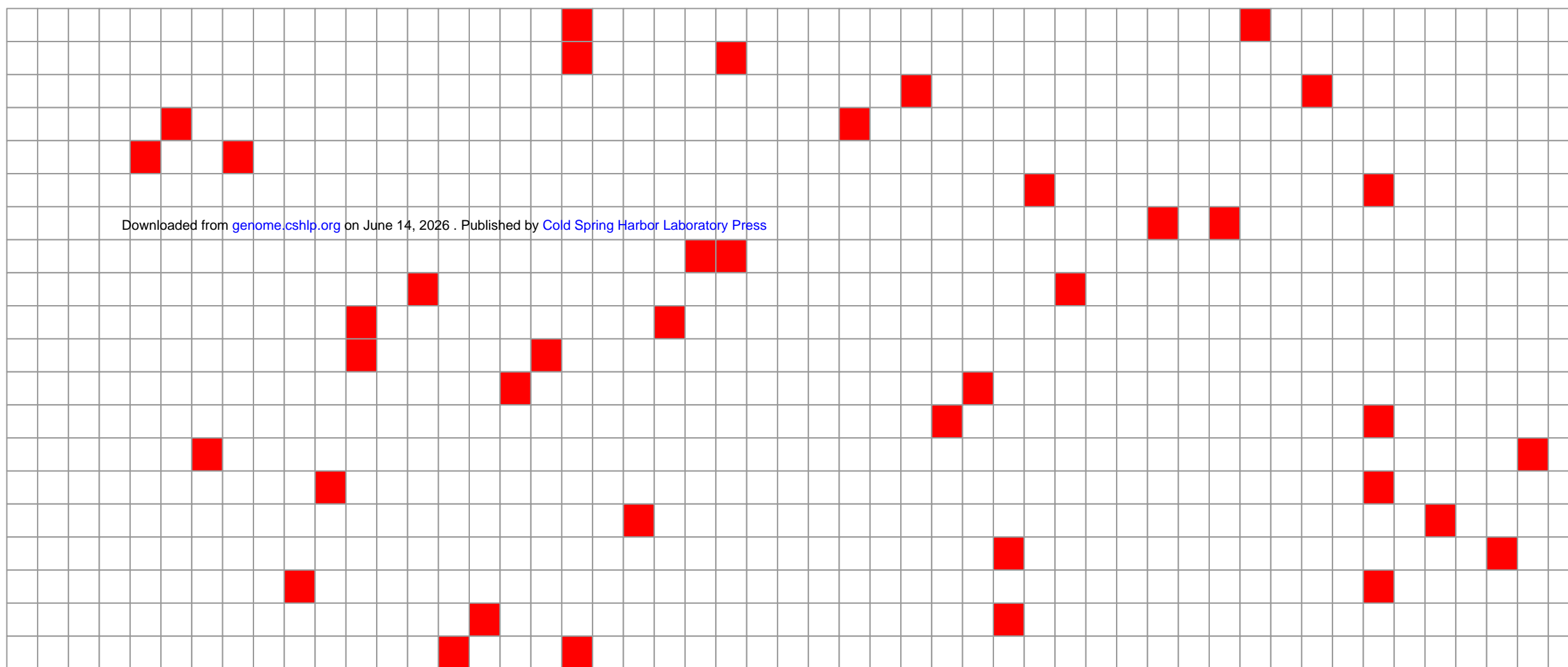
Fraction of all instances observed in siblings

Specific combinations of mutated genes significantly associated with comorbid phenotypes

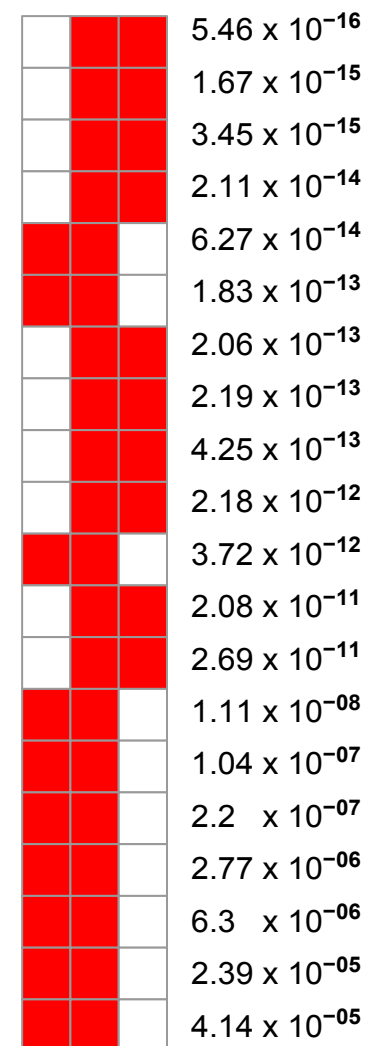
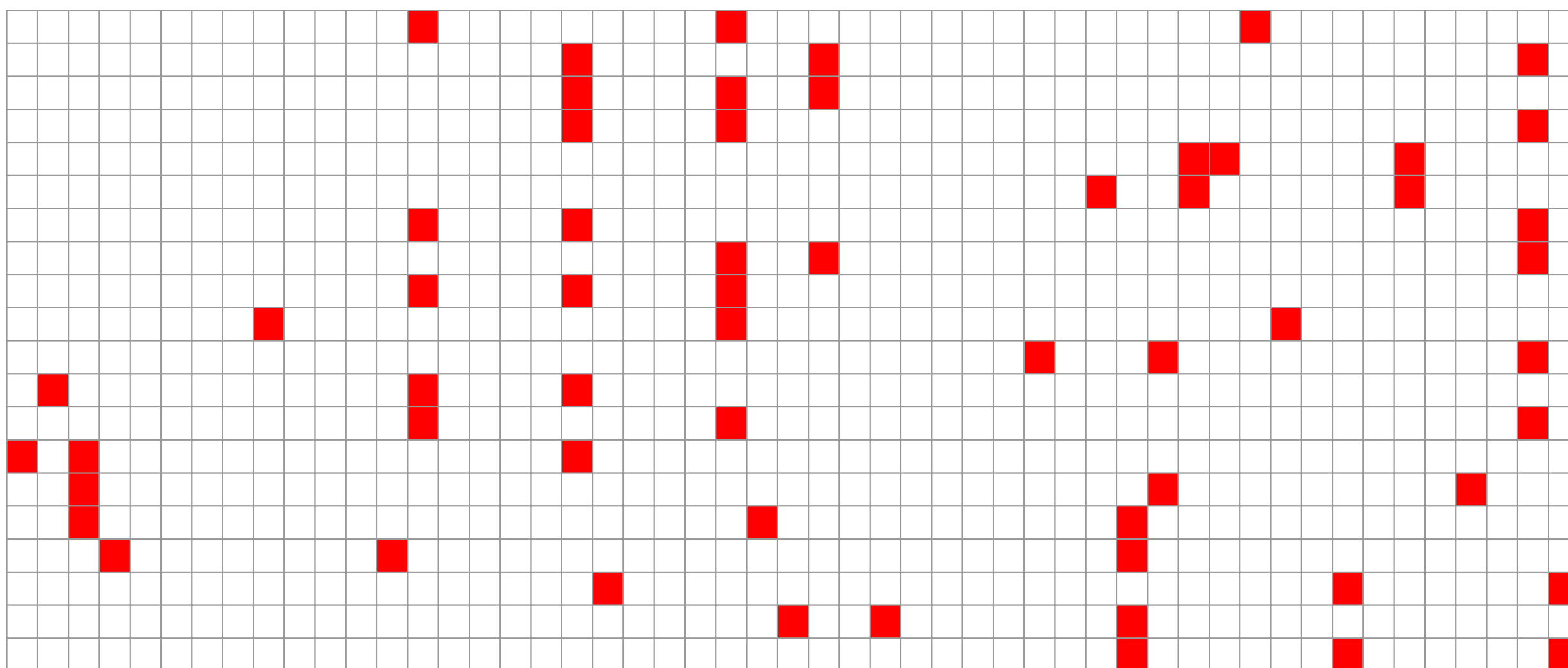
← Genotype →

Phenotype

↑
Two genes & two phenotypes
↓



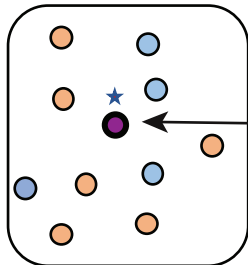
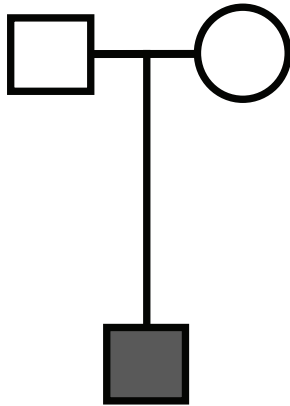
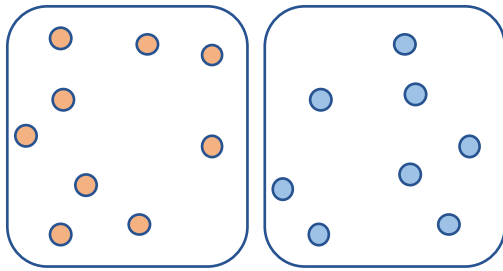
↑
Three genes & two phenotypes
↓



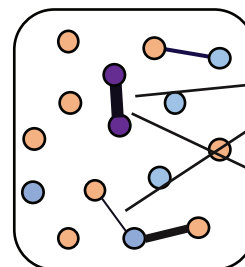
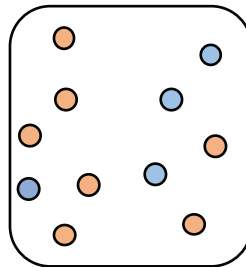
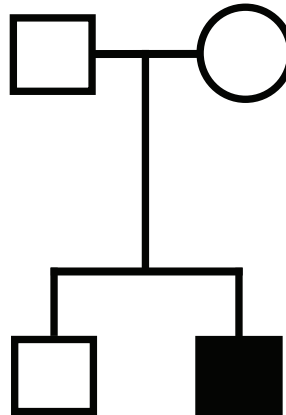
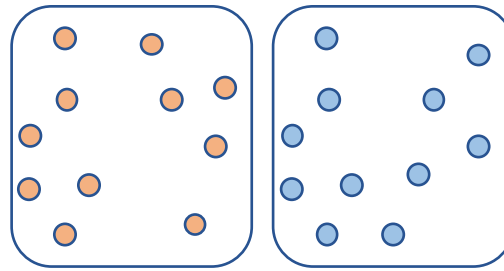
ABCC2
ADGRV1
AHNAK2
ALX1
ATP10B
C9orf64
CAPN9
CES5A
COL6A3
COQ6
DEF6
DNAH10
DNAH3
DNAH9
DNASE1
DPEP3
EFCAB7
EPB41L5
FETUB
FRAS1
GALNT11
GLDC
GPATCH4
HEATR5A
HHIP
HMCN1
IFT140
KIF27
MAN2B2
MPP5
MTR
MYBBP1A
MYOM2
NEK10
NINL
NLRP1
OBSCN
OR2G3
OR51M1
OR5AU1
PKP1
PRIMPOL
PRX
SACS
SGSH
SLC22A14
STAC
TEP1
THSD7B
VWDE
VWF

Anxiety/depression
Intellectual disability
Schizophrenia
p-values

Rare variants and their combinations associated with phenotype



De novo mutation
(autism subtype)



Variant Pairs



Phenotype

Int. Disability



Epilepsy

Effect Size



High



Low

Intellectual
Disability (ID)

Epilepsy